

Acknowledgment

First of all I would like to thank my two supervisors, Per Johan Brandvik and Øyvind Mikkelsen for help and guidance throughout this project. Special thanks to Per Johan for suggestions and comments on my work in this period, and for giving extremely valuable feedback at the end of this master period, even though I know it has been colliding with the moose hunting season!

Thank you to all the wonderful people at SINTEF Marine Environmental Technology for making me feel so welcome, and never minding me asking any questions! Special thanks to Marianne Unaas Rønsberg and Inger Kjersti Almås for help with laboratory work, and for running GC-FID and GC-MS samples. I would also like to thank Kjersti for teaching me how to use Chemstation and Cosiweb.

Special thanks my fiancé, for joining me on my fieldtrip to Sula, and for invaluable technical support with latex! Thank you for the support, interest and encouragement in this period. I would also like to thank my 15 month old daughter Oda for making sure that I would be the first student to arrive at SINTEF in the morning. Thank you for always being a reminder of what is important in life, words can't describe how much I look forward to spending more time with you!

Lastly I would like to thank my mom for always supporting my choices, and for always showing interest in what I do.

Trondheim, 15.10.16

Marie Myrstad

Abstract

This thesis aims to characterize 112 weathered oil samples collected on shore, at 18 islands, along the coastline of Mid-Norway during a time period of 2011-2015. Emphasis have been made on characterizing samples by three different multivariate methods; Principal component analysis (PCA), Partial least square-discriminant analysis (PLS-DA) and Hierarchical cluster analysis (HCA), however univariate methods have been applied as a starting point.

In multivariate data analysis, diagnostic ratios were calculated between biomarkers and PAH components and applied to look for interesting structures in the plot. The classification from univariate methods, were combined to identify the position of different oil types in the plots.

PCA, PLS-DA and HCA demonstrated their ability to categorize weathered samples, and identified samples that could not be identified by the traditional univariate method. The multivariate techniques were able to classify samples without some of the typical identifying biomarkers that are used in univariate oil spill forensics and indicates that multivariate techniques could be a promising method for identifying heavily weathered samples that often have inconclusive or missing measurements for typically used biomarkers and diagnostic ratios.

Selected samples were imported into an international oil spill database to identify matches to external samples from other projects and laboratories. Six samples in this study were a probable match to oil samples collected at the Shetland islands.

Sammendrag

Denne oppgaven har som mål å karakterisere 112 forvitrede oljeprøver som er samlet inn på 18 øyer langs Trøndelagskysten i løpet av tidsperioden 2011-2015. Prøvene har blitt karakterisert ved hjelp av tre ulike multivariate teknikker; Prinsipal komponent analyse (PCA), Delvis minstekvadrat-diskriminant analyse (PLS-DA) og Hierarkisk klyngeanalyse (HCA), og univariate statistiske metoder har blitt brukt som et utgangspunkt for dette.

I multivariat dataanalyse har diagnostiske ratioer blitt regnet ut mellom biomarkører og PAH komponenter, og blitt benyttet for å se etter interessante strukturer i plottet. Klassifiseringen fra den univariate metoden ble kombinert for å identifisere posisjonen til ulike oljetyper i plottet.

PCA, PLS-DA og HCA har vist at de kan kategorisere forvitrede prøver og identifisere prøver som ikke kunne bli identifisert av tradisjonelle univariate metoder. De multivariate metodene klarte å klassifisere prøver uten bruk av de typiske identifis-

erende biomarkørene som brukes i univariat oljesøl identifikasjon. Dette indikerer at multivariate teknikker kan være en lovende metode for å identifisere tungt forvitrede prøver som ofte har manglende målinger for de typiske brukte biomarkørene og diagnostiske ratioer.

Utvalgte prøver fra prøvesettet har blitt importert til en internasjonal oljesøl database for å identifiserte likheter mellom eksterne prøver og prøver fra dette datasettet. Seks prøver ble i dette prosjektet antatt å være en sannsynlig match med oljeprøver som er samlet inn fra Shetland øyene.

Contents

Acknowledgment	i
Abstract	ii
1 Introduction	2
1.1 Background	3
1.2 Objectives	5
1.3 Approach	6
1.4 Abbreviations	7
2 Theory	8
2.1 Factors affecting oil spill fingerprinting	8
2.2 Composition of Crude oil	9
2.2.1 Hydrocarbons	9
2.2.2 Non-hydrocarbons	9
2.3 Refinery products	10
2.4 Weathering of oils in water	10
2.4.1 Evaporation	11
2.4.2 Dissolution	12
2.4.3 Photo-Oxidation	13
2.4.4 Biodegradation	13
2.4.5 Sedimentation	13
2.4.6 Water in oil emulsification and natural dispersion	13
2.5 Biomarker	14
2.5.1 Characterization of biomarkers	14
2.6 Chromatography	17
2.6.1 Flame ionization detector	18
2.6.2 Gas Chromatography Mass Spectrometry	19
2.7 CEN methodology for oil spill identification	20
2.7.1 Level 1	22

2.7.2	Level 2	23
2.7.3	Level 3	23
2.8	Multivariate Statistics	24
2.8.1	Preprocessing	24
2.8.2	Principal Component Analysis	26
2.8.3	HCA	27
2.8.4	PLS-DA	28
2.9	Computerized oil spill identification database	29
3	Materials and Methods	31
3.1	Sample material	31
3.2	Sample material from SINTEF	31
3.2.1	Sample location	32
3.3	Sample material from Sula 2015	32
3.3.1	Inspection of samples	34
3.3.2	Experimental procedure	36
3.3.3	Experimental procedure of samples from Sula	36
3.3.4	Preparation of oil samples	36
3.3.5	GC-FID conditions	37
3.3.6	Solid phase extraction	37
3.3.7	GC-MS conditions	37
3.4	Data treatment	38
3.4.1	Integration	38
3.4.2	Noise	38
3.4.3	Inspection of chromatograms	39
3.4.4	Diagnostic Ratios	40
3.4.5	Statistical analysis	40
3.4.6	COSIWeb	41
4	Results	44
4.1	CEN	44
4.2	Pretreatment of raw data	44
4.2.1	Descriptive Statistics	50
4.3	Multivariate data analysis	52
4.3.1	Principal component analysis	52
4.3.2	Partial Least Square-Discriminant analysis	65
4.3.3	PLS-DA: Crude oils	69
4.3.4	Hierarchical Cluster Analysis	72

4.3.5	Color coding of groups in HCA and PCA based on results inspection of chromatograms	75
4.3.6	COSIWeb	75
5	Discussion	80
5.1	Sources of errors	80
5.2	Inspection of chromatograms	81
5.3	Principal component analysis	83
5.3.1	Combining results from inspection of chromatograms with results from PCA-2D score plot	83
5.3.2	Combining results from inspection of chromatograms with results from PCA-3D score plot	83
5.3.3	Biplott	84
5.4	PLS-DA	85
5.4.1	PLS-DA-for crude oil and bunker oil	85
5.4.2	PLS-DA-for crude oil and non north sea crude oil	86
5.5	HCA	86
5.6	Search in the international database COSIWeb	87
5.7	Comparing multivariate methods	88
6	Conclusion	89
7	Further work	90
	Bibliography	91
A	Description of Oil Samples	97
B	PAH and Biomarkers	101
C	Diagnostic ratios	104
D	GC-FID	107
E	COSIWeb	142

Chapter 1

Introduction

Oil spilled into the sea represents a threat to the environment and marine life. In addition, cleaning up procedures are costly, and will only be able to limit the damage an oil spill may cause. A majority of the oil spills covered by the media, are ships accidents and offshore releases. These types of oil spill are only a small part of the total releases worldwide. Nevertheless, they are typically few events that each is large in scale. They are given great attention because of the amount of oil that are released, fatalities, large local environmental problems and large clean-up costs [Fingas, 2010] [Wang and Stout, 2010]. Acute oil spills are divided into accidents and intentional operational discharges. Accidental oil spills may occur from pipeline spills, blowout from wells, collision between vessels, or if vessels run aground. Examples of intentional operational discharges include discharges from vessels, such a bilge water. Bilge water is oily wastewater from ships and contain a mixture of petroleum and other compounds [Wang and Stout, 2010]. In many ships it is customary to run bilge water through an oil-water separator, but it is suspected that poor maintenance may cause releases of oil with concentrations above the legal level. Equipment failure and human errors are the largest causes for accidental oil spills. Transportation of petroleum from oil fields to shore involves different transportation steps and transportation alternatives. These are factors that increase the chance of an oil spill to occur [Fingas, 2010].

Both the industry and governments have been working to prevent oil spills worldwide. This has lead to increased focus on operating and maintenance procedures. In addition, pollution protocols and protocols for operation of ships to prevent pollution have been established [Fingas, 2010] [International Maritime Organization, 2016]. According to the database provided by *The International Tanker Owners Pollution Federation (ITOPF)*, the number of large spills, which is defined as spills

over 700 tonnes, have decreased from 1970 to 2015, even though ship traffic from oil trade have increased in this period [ITOPF, 2016].

In oil spill forensic, chemical fingerprinting is used for the purpose of localizing the source of an oil spill by applying analytical techniques such as Gas Chromatography-Flame Ionisation Detector (GC-FID) and Gas Chromatography-Mass Spectrometry (GC-MS). This has to hold in court, thus making it possible to take legal actions against the accountable party [Wang and Stout, 2010]. Efforts have been made to make the chemical fingerprinting as defensible as possible. The CEN methodology was established in 2002 with the aim of being a forensic tool for comparison and identification of oil spills, by using different analytical techniques and gives legal support when actions are taken against the party responsible for the oil spill in the environment [CEN, 2012]. The Bonn agreement was established in the late 1970s between north sea countries and European union. Their aim is to work together and assist each other to detect and prevent pollution of oil [Wang and Stout, 2010] [Bonn Agreement, 2016]. Computerized oil spill identification (COSIWeb) database was created in 1999, and is a database that strengthens the work with identification of unknown oil samples. This database contain oil samples from accidental oil spills, with both known and unknown oil types and different crude oils from around the world. The database is available through any browser and allows the user to upload oil samples into the database which is then available for all other users. The user may then compare the uploaded oil samples with other samples in the database to investigate possible matches [COSIWeb, 2016] [Stout and Wang, 2016].

The CEN methodology, and COSIweb database are based on univariate statistics, however the use of multivariate data analysis is becoming more used within oil spill forensic. Multivariate analysis offer many advantages such as handling many variables at the same time, and give greater insight into the data by explaining hidden structures. PCA and PLS are examples of multivariate techniques in oil spill forensics [Christensen et al., 2004] [Nielsen et al., 2012] [Christensen and Tomasi, 2007] [Wang et al., 1999] [Stout et al., 2001].

1.1 Background

The coastline of Norway can be exposed to contamination of oil from ship traffic and from fields in the Norwegian sea. This has the potential to affect areas such the marine life, seabirds and fisheries [Norwegian Environment Agency, 2016]. There are examples of major oil spills that have polluted the Norwegian coastline, such as the Full City and Server, both shipping accidents. In 2009 Full city ran aground near Langesund in Telemark. Approximately 293 tonnes of heavy oil leaked out and 75

km of the coastline was contaminated. They were able to collect about 100 tonnes of oil. In 2007 the MS Server ran aground near Fedje in Hordaland. 388 tonnes of heavy bunker oil was released into the environment [S. Boitsov and Dolva, 2013].

Figure 1.1 shows a map of the oil fields collectively known as the Haltenbanken fields, and ship traffic around Mid-Norway.

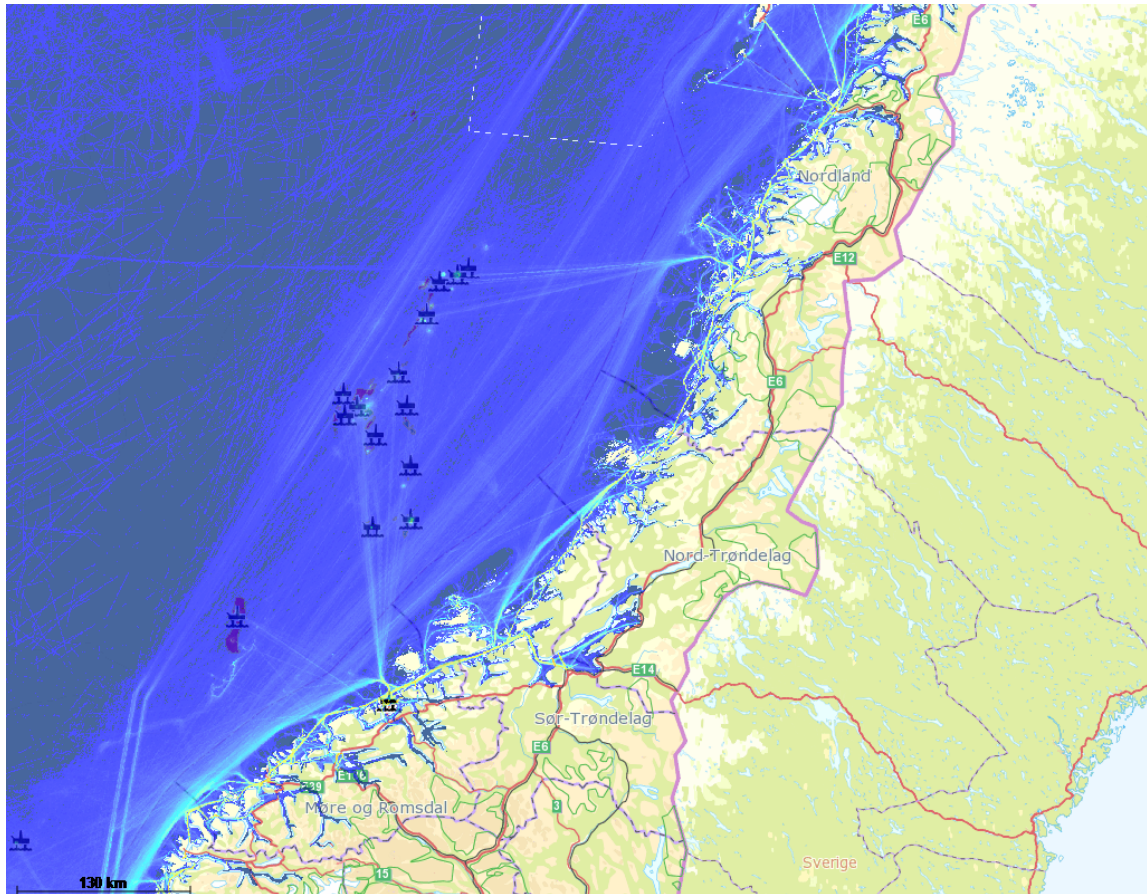


Figure 1.1: Overview over fairways and oil fields close to Mid-Norway [Norwegian Environment Agency, 2016]

A consequence analysis was published in 1998 for the oilfields in the Haltenbanken area. This report says that blowouts releases, leaks from floating storages, and accidents with shuttle tankers are accidents that most likely would cause the largest spills of oil in the Mid-Norwegian shelf, typically in the range between a few thousand tons to over 300,000 tons. The probability of these oil spills was very low and summed up to a release every 38 years in year 2000 and every 24 years in 2009. Calculation of

oil spill trajectories from the Mid-Norwegian shelf, presume that the oil will spread to the northeast and east due to the predominant wind direction being south to southwest, and due to the background currents that follows the coast northward. The coastline that may be affected from a blowout in the Mid-Norwegian shelf are areas from the south of Hitra to Kvaløya in Troms [Statoil, 1998].

The following quote is also written on the Norwegian Petroleum webpage [Norwegian Petroleum, 2016]:

The Norwegian petroleum industry has not been the cause of any major oil spills that have resulted in environmental damage. In the roughly 50 years since Norway's petroleum activities began, no oil spill from the industry has ever reached the shoreline. (Norwegian Petroleum webpage, 08.03.2016)

As a part of the course *KJ3050-Marine Organic Environmental Chemistry* provided by NTNU and SINTEF, students have searched for weathered samples of oil spills along the coastline of Mid-Norway during a time period of 2010 to 2015. A majority of these islands are a part of the Froan nature reserve, which an area were the aim is to protect living and nesting sites for birds and marine mammals, as well as flora and fauna [Visit Norway, 2016]. Approximately 400 samples have been collected during these fieldtrips and interpreted by students as a part of this course. Work with part of this material has also been presented by earlier master students [Henriksen, 2012] [Vike, 2014]. However, nobody have evaluated the complete dataset and analyzed it for the major trends.

1.2 Objectives

The objective of this thesis is to

- Identify trends in sample material collected by previous students in the period 2011-2015 from *KJ3050* by using multivariate techniques.
- Characterize samples according to oil type.
- Collect and characterize oil samples from Sula.
- Import selected sample material to COSIWeb with the aim of sharing results with other oil spill laboratories, and search for identifying samples in the database.

1.3 Approach

Sample material collected at Sula will be analysed by GC-FID and GC-MS. The entire sample material will be assessed by inspection of chromatograms according to the CEN procedure to identify the type of oil (Crude or Bunker). The sample material will be analyzed by different multivariate techniques to observe trends in the dataset. Different multivariate techniques include principal component analysis (PCA), Hierarchical Cluster Analysis (HCA), and Partial least squares discriminant analysis (PLS-DA). Part of the sample material will be included and analyzed in the COSIWeb database.

1.4 Abbreviations

AGNES	Agglomerative Nesting
CEN	European Committee for Standardization
CI	Chemical ionization
DCM	Dichloromethane
DR	Diagnostic Ratio
EI	Electron Ionization
FID	Flame Ionization Detector
GC	Gas Chromatography
HCA	Hierarchical Cluster Analysis
HFO	Heavy Fuel Oil
LFO	Light Fuel Oil
LOQ	Limit Of Quantification
MS	Mass Spectrometry
SIM	Selected Ion Monitoring
MVA	Multivariate Analysis
M/Z	Mass to Charge Ratio
NA	Not Available
NS	North sea
NTNU	Norwegian University of Science and Technology
PAH	Polyaromatic Hydrocarbon
PC	Principal Component
PCA	Principal Component Analysis
PLS-DA	Partial Least Square-Discriminant Analysis
RSD	Relative Standard Deviation
SD	Standard Deviation
UCM	Unresolved Complex Mixture

Chapter 2

Theory

2.1 Factors affecting oil spill fingerprinting

In oil spill forensic the idea is to characterize different properties of the oil to form a conclusion about the type of oil and origin of the oil spill. These properties will depend on factors both prior to the spill, such as the origin of crude oil, and refining processes of the oil. It will also depend on factors after the spill, such as weathering and mixing. Crude oil is naturally occurring petroleum (oil that is found underground), and is formed as a result of various geological processes. The composition of crude oil depends on the origins of the oil, hence the chemical and physical properties of crude oil will vary. This also means that when crude oil is spilt at sea they will behave differently [Wang and Stout, 2010].

Crude oil is refined into (manmade) petroleum products. Examples of refined petroleum products are fuel oil, lubricants and sludge products. The composition of these refined petroleum product is a mixture of the original chemical properties, and new chemical properties, since some alterations of the chemical compositions occur during the refinery process. It has for example been shown that distillation and thermal cracking may affect certain PAHs and biomarkers [Wang and Stout, 2010]. When oil is released into the sea there are both natural and anthropogenic sources that may mix with the oil spill. These sources contribute to a large amount of the total oil in the environment. Examples of anthropogenic sources are runoff and deposition of air pollution. Natural sources may come from natural oil seeps [Wang and Stout, 2010]. Moreover in the event of an oil spill, oil may mix with other oils if several oils are included in an oil spill (e.g. if tankers collide) which makes the fingerprinting even more demanding.

2.2 Composition of Crude oil

The composition of crude oil can be divided into two main groups; hydrocarbons and nonhydrocarbons. The majority of compounds are hydrocarbons. As the name implies, hydrocarbons consist of hydrogens (12-14 %) and carbons (84-87 %), and these vary in complexity from light, volatile compounds such as propane to heavy compounds such as waxes. They can be straight chained, branched or cyclic, and their chemical bonds can be saturated or unsaturated. In addition, crude oil contain elements of sulphur, nitrogen and oxygen as well as trace metals such as vanadium, nickel and chromium [Brandvik and Daling, 2014].

2.2.1 Hydrocarbons

Hydrocarbons can be further classified into paraffins, naphthenes and aromatics [Simanzhenkov and Idem, 2003]. Paraffins are saturated, straight-chain alkanes and isomers of alkanes. Straight chained alkanes up to four carbon atoms are in gaseous form (methane, ethane, propane and buthane), and straight chained alkanes with 5 to 17 carbon atoms are liquids. When the number of carbon atoms in a straight chain configuration reaches 18 carbons they are termed waxes. The occurrence of paraffins in crude oil can vary from 2-50 % [Simanzhenkov and Idem, 2003]. Naphthenes are saturated cyclic alkanes with one or more ring structure and may include paraffinic side chains [Speight, 2014], see figure 2.1. Aromatic hydrocarbons are unsaturated compounds containing one or more aromatic ring structure such as benzene, and these can be connected to paraffinic hydrocarbons or naphthene rings [Speight, 2014].

2.2.2 Non-hydrocarbons

Crude oil contain small portions of sulphur, nitrogen and oxygen and trace metals of vanadium, nickel and chromium. Asphaltenes and resins are examples of non-hydrocarbons. Resins are compounds that show polar character compared to hydrocarbons. Their molecular weight lies in the area between 700-1000. They often contain compounds such as carboxylic acids, sulphoxides and phenols and are dark in color [Simanzhenkov and Idem, 2003] [Brandvik and Daling, 2014] Asphaltenes are compounds with large molecular weights ranging from 1000 to 10 000. Their structures are poorly described in the litterature, and they are one of the most complex compounds in crude oils. They are red to brown in in color and are described as polycyclic aromatic compounds. [Simanzhenkov and Idem, 2003] [Brandvik and Daling, 2014].

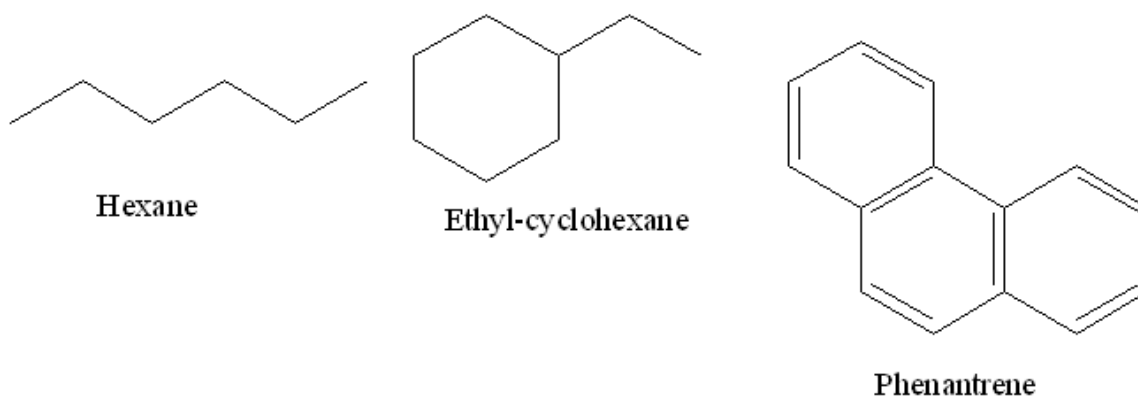


Figure 2.1: Examples of different hydrocarbons. From left to right: straight chained alkane (hexane), a naphthenic structure (Ethyl-cyclohexane) and polyaromatic hydrocarbon (Phenanthrene). Chemical structures are made in Chemdraw.

2.3 Refinery products

In this thesis bunker oil is used as a collective term for oil products formed from refinery processes and include fuel oil, lubricating oil and sludge. Light fuel oil represents distilled intermediates and fuels such as gas oil and marine diesel. Biomarkers with high boiling point are typically not present in these oils [CEN, 2012]. Heavy fuel oil are oils with high density and viscosity. The refinery processes often changes the methylphenanthrene pattern due to a catalytic cracker which alters some of the aromatic patterns. It increases the amount of methylanthracene and reduces the concentration of retene [CEN, 2012]. Sludge is a high fuel oil, or a mixture of HFO and lubricating oil (motor oil). It typically do not contain the biomarker retene [CEN, 2012].

2.4 Weathering of oils in water

When oil products are spilt at sea there are different weathering processes that will alter the physical and chemical properties of the oil. The oils degree of weathering will be dependent on the oils original physical and chemical properties as well as environmental conditions (waves, wind, temperature etc) and properties of the sea (currents, temperature, salinity, density, oxygen, bacteria etc). Weathering can further be divided into evaporation, dissolution, emulsification, redistribution of compo-

nents, biodegradation, chemical alteration and contamination [Brandvik and Daling, 2014] [CEN, 2012]. Figure 2.2 shows the different weathering processes of oil at sea.

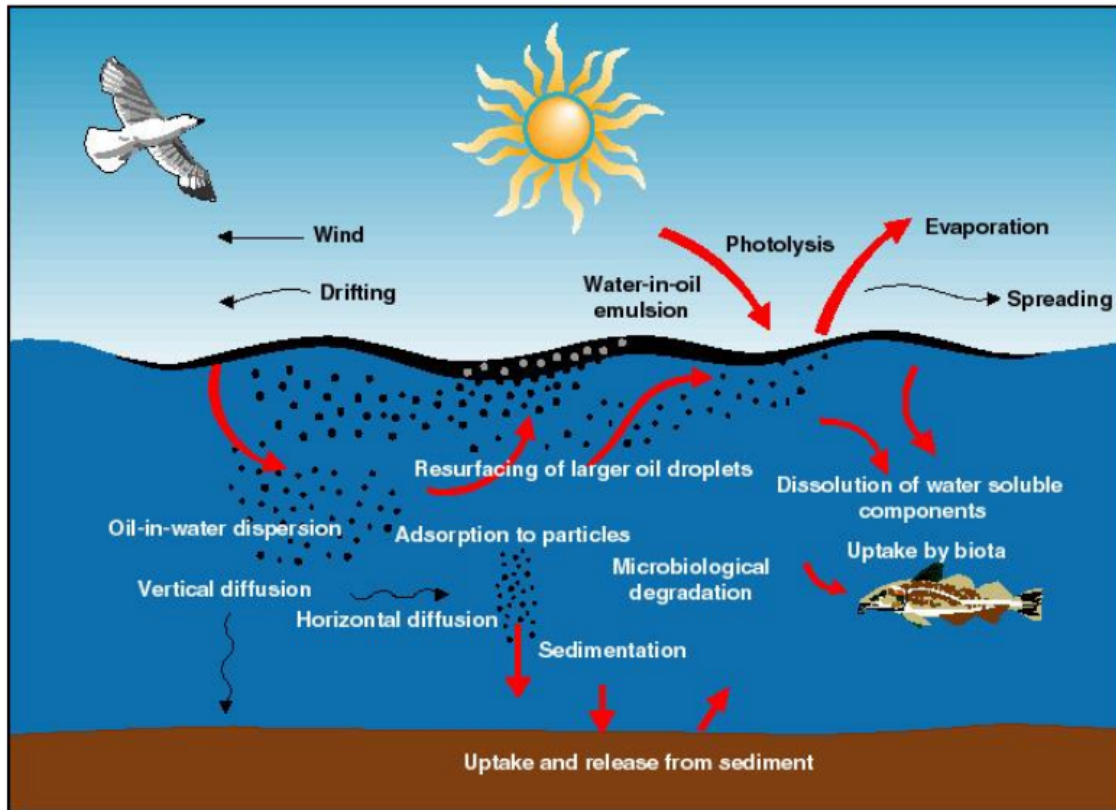


Figure 2.2: Weathering processes of oil spilt at sea. Image courtesy of SINTEF.

2.4.1 Evaporation

Evaporation is one of the first weathering processes that occurs when oil is spilt at sea, and is especially prominent for the first days as presented in figure 2.3. It is one of the most important processes that removes oil from water. In particular, components with low boiling point are exposed to this weathering process. As a rule of thumb, oil components with boiling point less than 200°C will evaporate within 12-24 hours. This correspond to n-alkanes up to approximately $n - C_{11}$. The degree of evaporation will depend on the amount of light compounds in the oil, as well as temperature in the sea, wind speed and the thickness of the oil slick [Brandvik and

Daling, 2014] [CEN, 2012].

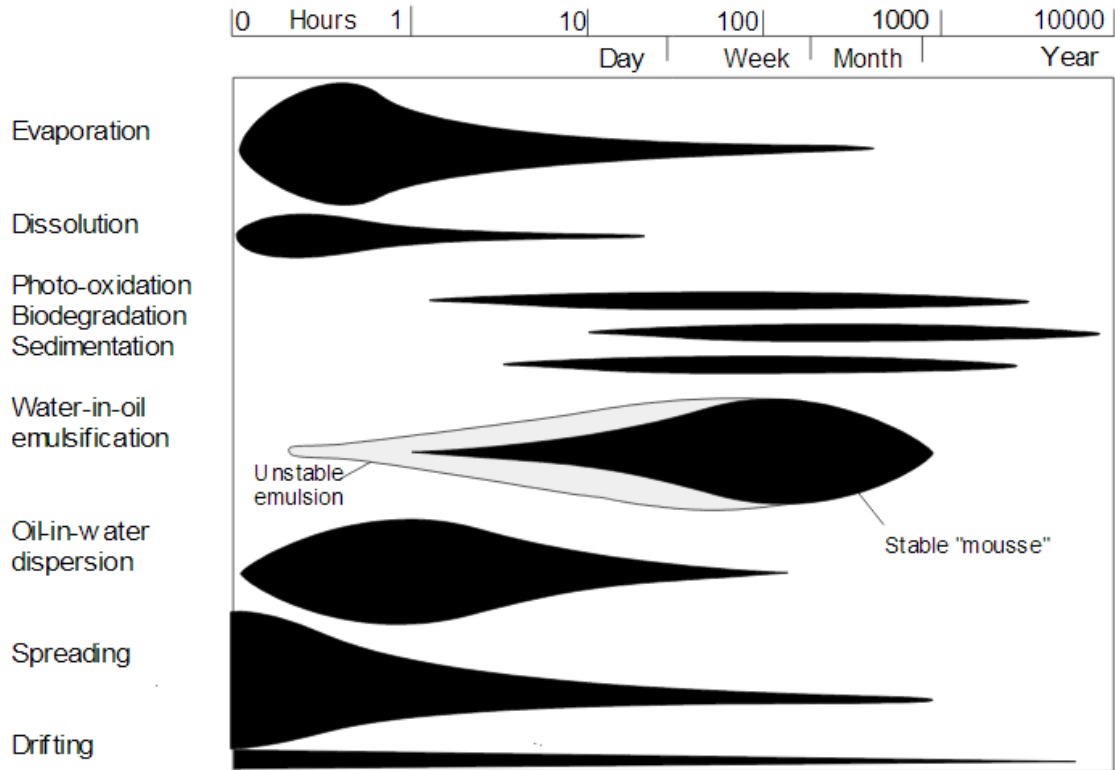


Figure 2.3: Dominance of weathering processes at different times. Image courtesy of SINTEF.

2.4.2 Dissolution

Dissolution is a process that removes the most soluble components in the oil and applies to smaller compounds such as hetero compounds and low substituted aromatic hydrocarbons. These compounds however, tend to evaporate more quickly and dissolution is not always a prominent factor for removing oil from the sea, except in those cases where there is a high degree natural dispersion of oil [Brandvik and Daling, 2014] [CEN, 2012].

2.4.3 Photo-Oxidation

Photo oxidation is a process that alters some of the chemical composition of oil due to the effect sunlight. Aromatics are especially prone to this effect, and they are oxidized to resins and eventually asphaltenes. This actually stabilizes water-in-oil emulsions and increases the oils persistence at sea. Because of this, slowly degradable tarballs can be formed and drift on the sea for a long time or drift at shore [Brandvik and Daling, 2014] [CEN, 2012].

2.4.4 Biodegradation

Microorganisms in the sea, such as bacteria, may use the oil to gain nutrition. Biodegradation is not prominent before approximately one to two weeks after oil is spilt at sea, see figure 2.3. It can occur in any type of oil component except for asphaltenes, but is especially evident in straight-chain hydrocarbons. The degree of biodegradation will depend on the amount of nutrients in the oil (nitrogen and phosphate) as well as oxygen and temperature. Biodegradation takes place at the water-to-oil interface which means that oil that has drifted on shored degrade slower. The degree of biodegradatoin can be estimated by the loss of n-C17 and n-C18 compared to pristane and phytane [Brandvik and Daling, 2014] [CEN, 2012].

2.4.5 Sedimentation

Crude oil residues will normally not have higher density than water and will not sink, except in cases where there are high concentrations of sediments, in which the sediments may cling or stick to the oil and sink. Recent developments in oil refineries makes the oil more dense and hence sinking of heavy fuel oils may become a larger problem in the future [Brandvik and Daling, 2014].

2.4.6 Water in oil emulsification and natural dispersion

There are two processes that happens simultaneously after an oil spill, namely water-in-oil emulsification and natural dispersion, also known as oil-in-water emulsification. Natural dispersion refers to oil that break into droplets and mix with water and the water column. This removes oil from the sea surface and leads to natural breakdown of the oil. This process occurs when there are breaking waves (typically when wind speed is above 5 m/s). Natural dispersion is one of the dominant processes the first days after an oil spill. Formation of water-in-oil emulsification means that water droplets are in a continuous oil phase. The oil will become more persistent and

remain at sea or on shore. Since the viscosity of an oil increases due to weathering, water-in-oil emulsification will eventually be the dominant process compared to natural dispersion [Brandvik and Daling, 2014].

2.5 Biomarker

In order to localize the source of an oil spill we have to look for identifying markers or fingerprints. The fingerprint of the oil is what we call a biomarker. Biomarkers are described as “complex molecules derived from formerly living organisms” [Wang et al., 2006] [Kao et al., 2015] [Wang and Stout, 2010]. They are valuable for the purpose of fingerprinting since the organic structures are resistant to environmental degradation and hence show no variation, or only small variation in structure from their parent organic molecule. This makes it possible to extract specific information regarding the source of the spilled oil. In addition, biomarkers are useful for differentiating and correlating oils, monitoring degradation processes and evaluate the weathering state of oils. They can be detected in low quantities by methods as gas chromatography-mass spectrometry (GC-MS). By performing chemical analysis of various biomarkers, the analyser can gain information about the chemical composition of the oil which is helpful in the puzzle for determining the source of the spill [Wang and Stout, 2010].

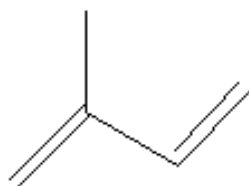
2.5.1 Characterization of biomarkers

The very basic structure of most biomarkers are the isoprene unit, which contains five carbon atoms, see figure 2.4. The chemical formula is C_5H_8 [Peters et al., 2005] [Wang and Stout, 2010].

Compounds that consists of these subunits are termed terpenoids or isoprenoids. The terpenoids consists of a large group of biomarkers with either cyclic or acyclic configuration. They can be saturated, have double bonds or hold other elements besides carbon and hydrogen. The subunits can be linked in a regular manner (head-to-tail) or irregular manner (tail-to-tail or head-to-head) [Peters et al., 2005] [Wang and Stout, 2010]. Examples of terpenoids that are important for oil spill identification are terpanes, steranes and aromatic steranes [Wang and Stout, 2010].

Terpanes

Terpanes that are characteristic in crude oil include Sesqui- (C15), di- (C20) and triterpanes (C30). As the name implies the sesquiterpanes, diterpanes and triterpanes consists of three, four and six isoprene subunits respectively [Wang and Stout,



Isoprene subunit

Figure 2.4: Isoprene subunit. Chemical structure is drawn in ChemDraw

2010]. Pristane, phytane and squalene are examples of acyclic terpanes and their chemical structures are presented in figure 2.5.

Squalene is classified as an acyclic triterpanes and are linked in a irregular manner, with six isoprene subunits and one tail-to-tail linkage. Pristane and phytane belong the the acyclic diterpanes. They are linked head-to-tail and stem from phytol which is the side chain of chlorofyll. Generally, acyclic isoprenoids are more resistant to biodegradation than the n-alkanes. Hence ratios between n-alkanes and isoprenoids are useful for indicating the extent of biodegradation. Examples of such ratios include pristane/n-c17 and phytane/n-c18 where the ratio will decrease along with increased weathering [Waples, 1985]. Hopanes belong to the group of pentacyclic triterpanes. They exist in a naphthenic structure with four six-membered rings and one five membered ring. They consist of three stereoisomeric series; $17\alpha, 21\beta$ and $17\beta, 21\beta$ and $17\alpha, 21\alpha$ where α and β indicate whether hydrogen atoms are below or above the plane of the rings [Peters et al., 2005]. $17\alpha(H), 21\beta(H)$ -configuration of hopane ranging from 27 to 35 carbon atoms are the most abundant in petroleum compared to the other configurations and helpful for oil identification due to its ther-

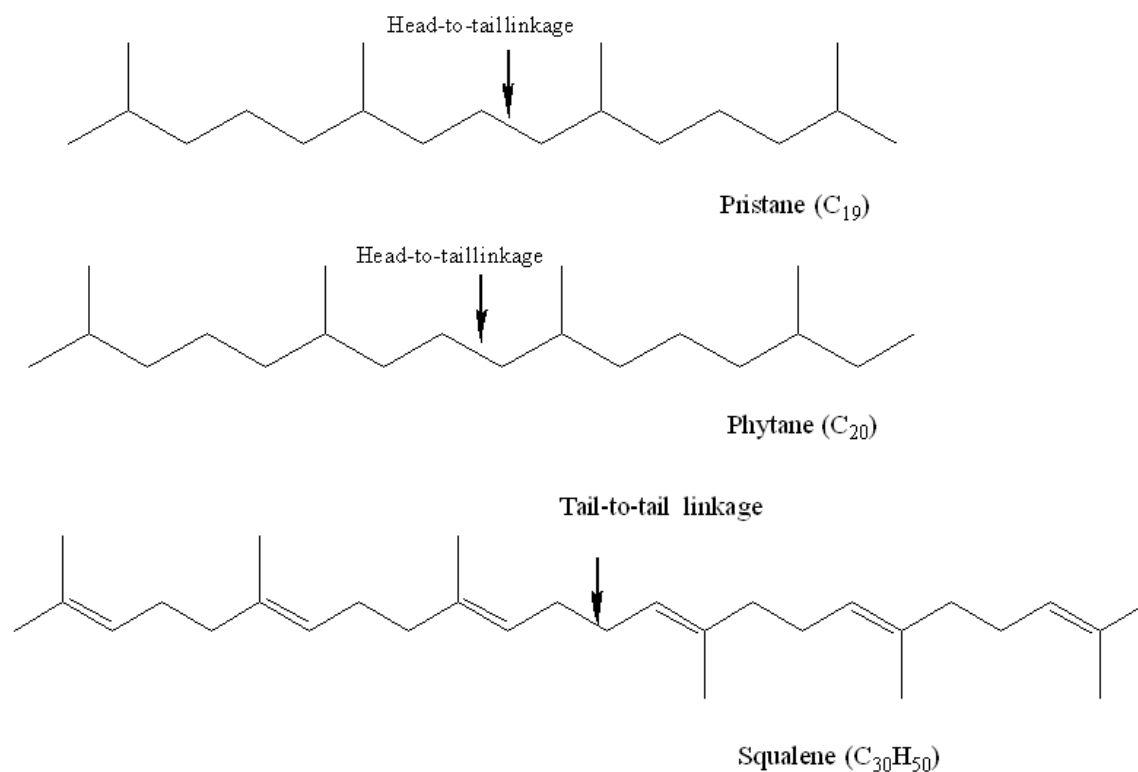


Figure 2.5: Chemical structure of pristane, phytane and squalane. The chemical structures are made in Chemdraw.

modynamic stability [Peters et al., 2005]. They are usually detected using m/z 191 chromatograms [Wang and Stout, 2010].

Gammacerane and Oleanane are other pentacyclic triterpanes that are characteristics in oil spill fingerprinting. Gammacerane is characteristic for marine and lacustrine environments. Oleanane is characteristic for lacustrine environments and are derived from terrestrial plants [Waples, 1985].

Steranes and aromatic steranes

The sterane family is characterized by the 4-cyclic arrangement and 21 to 30 carbon atoms. The homologous series containing 27, 28 and 29 carbon atoms have high source specificity and common steranes. They are detected using m/z 217 [Wang and Stout, 2010].

Among the aromatic steranes, the monoaromatic and triaromatic steranes are

useful in oil spill fingerprinting for differentiation and source identification [Wang and Stout, 2010].

2.6 Chromatography

Gas chromatography-flame ionization detector (GC-FID) and gas-chromatography mass spectrometry (GC-MS) are among the most frequently used instruments for characterizing hydrocarbons in oil [Wang and Fingas, 2003].

In general, chromatography is a separation technique where sample components distribute between a mobile phase and a stationary phase [Ettre, 1993]. Sample components separate and elute through the column at different retention times. This is because they interact differently with the stationary phase. A detector at the end of the column gives a signal of the components as a function of time or volume [Lundanes et al., 2013].

In gas chromatography the mobile phase is an inert gas, which means that the gas does not interact with the components or the stationary phase. The stationary phase can either be a solid adsorbent or a liquid stationary phase, however the latter is more common. Sample components are separated due to different vapor pressure and different interactions with the stationary phase [Grob and Barry, 2004]. The basic components of a GC instrument are the carrier gas tank, flow regulators, the sample injection chamber, the column, detector and a data system as presented in figure 2.6. The most common carrier gases are nitrogen, helium or hydrogen and they are contained in a pressurized cylinder [Lundanes et al., 2013].

The sample is introduced into the column through an injection system. Type of injection system will ultimately depend on the sample, and column type at use [Lundanes et al., 2013]. For capillary columns the most common injection systems are split/splitless injections, on-column and programmed-temperature vaporization [Grob and Barry, 2004]. In a split/splitless injection system the sample is introduced by a syringe needle through a septum into a glass liner. The sample is vaporized and mixed with carrier gas. In split mode the mixture is separated in two parts. One part contains only a fraction of the mixture which enters the column, and the remaining part which contains the largest volume is sent to waste. [Lundanes et al., 2013]. In splitless mode the whole sample volume is directed to the column. Two types of columns exist in gas chromatography; Packed columns and capillary columns. Today the use of capillary columns are frequently used in environmental laboratories [Fingas, 2010]. The column is placed in an oven so that the mixture is maintained in a vaporized state, and temperatures can vary from 40 to 350 °C [Meyers, 2011].

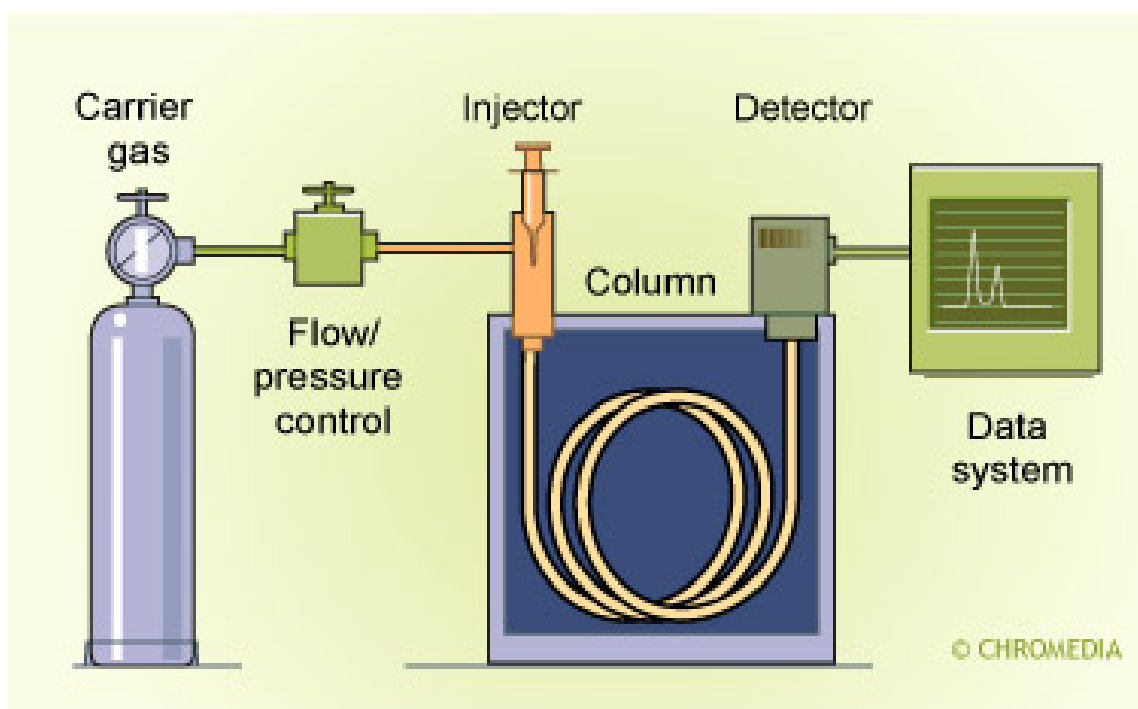


Figure 2.6: Instrument of a Gas Chromatogram. Image courtesy of Chromedia [Chromedia, 2016].

2.6.1 Flame ionization detector

Among the many detectors that exist in gas chromatography, the flame ionization detector is one of the most applied detector [Skoog, 2004]. Some of the reasons for this owes to its response to organic compounds, low detection limits, ease of use, high sensitivity (10^{-13}g/s), and large linear range (10^7g/s) [Grob and Barry, 2004]. The detector responds to ions that are produced in a flame. The carrier gas containing the sample components leaves the column and enters the detection compartment through a jet tip as presented in figure 2.7. Hydrogen gas is added and a flame is started at the end of the jet tip. Air enters the detection compartment through a separate channel. Compounds containing hydrocarbons goes through a set of chemical reactions and form ions. Detection is achieved by having a collector electrode with a potential a few volts higher than the potential of the flame. This generates a current which is proportional with the amount of carbon ions in the flame. It is a mass-sensitive detector and minor changes in temperature, flow or pressure does not affect the response [Skoog, 2004] [Grob and Barry, 2004].

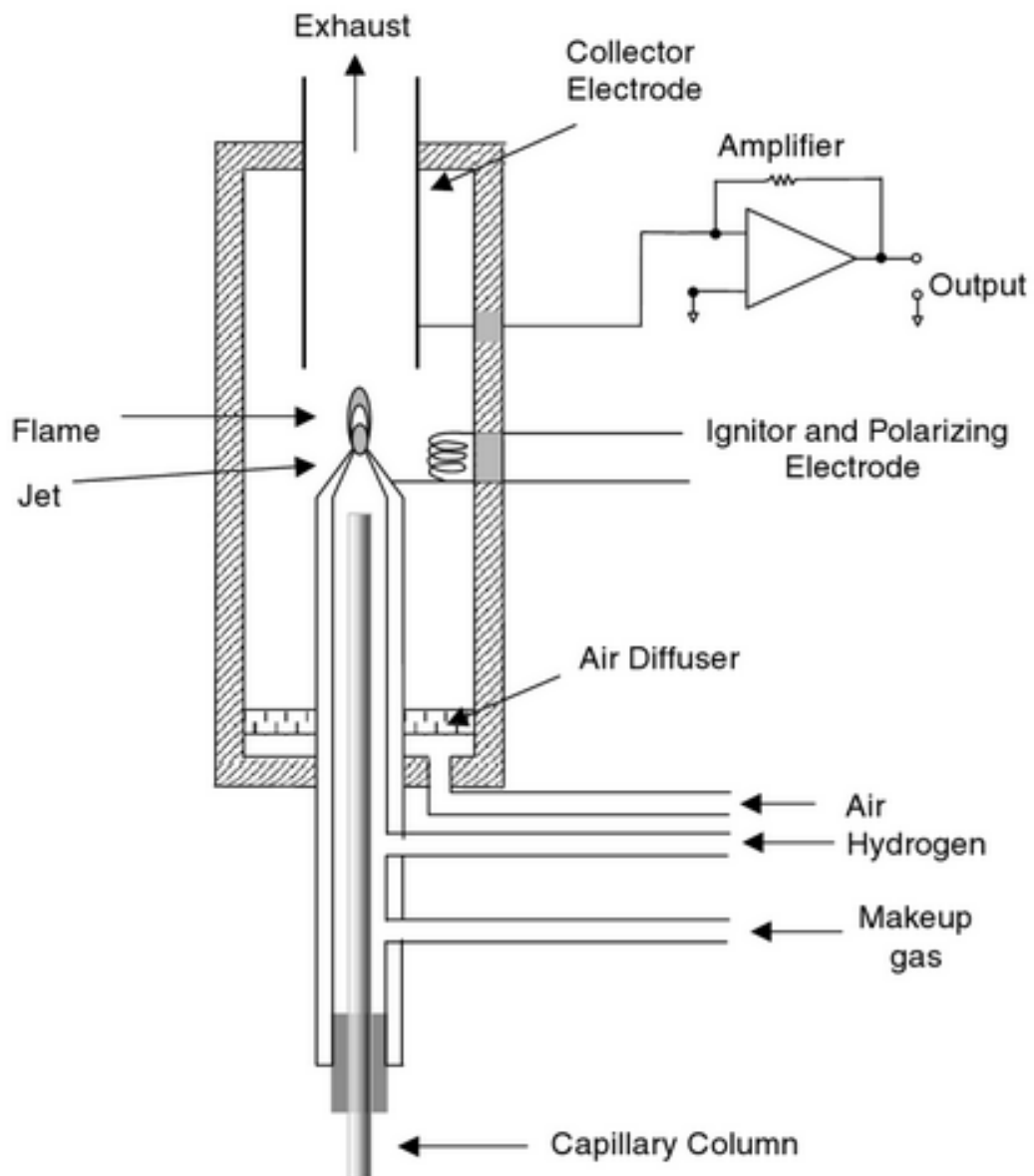


Figure 2.7: Flame ionization detector [Grob and Barry, 2004].

2.6.2 Gas Chromatography Mass Spectrometry

Chromatography separates a sample solution into its individual components and provide information about its peak height or area which makes it possible to gain

quantitative data. However, it cannot provide information about the structure of the components. To achieve this the chromatogram is combined with a mass spectrometer. The results from a mass spectrometer provides information regarding their chemical identity. Moreover GC-MS has a chemical and electron ionization database, which makes identification of unknown compounds possible [Lundanes et al., 2013].

The instrument consist of an ion source, mass analyzer and a detector. The ion source ionizes atoms or molecules to ions. In GC-MS, two types of ionization techniques are commonly used; Electron ionization (EI) and chemical ionization (CI) whereby the most applied technique is EI. In this technique the sample is placed under high vacuum and subjected to a beam of high energy electrons (70 eV). This converts the sample into molecular ions and fragments of molecular ions which makes up the mass spectrum [Poole, 2003] [Grob and Barry, 2004].

The mass analyzer separates the ions based on their mass to charge ratio (m/z). Linear quadrupole is one of the most commonly used mass analyzer in GC-MS, mainly due to fast scanning rate and low cost [Grob and Barry, 2004]. Linear quadrupole is made up of four cylindrical rods alligned pararell to each other. Both radiofrequency (RF) and direct-current (DC) potentials are applied and this create and oscillating electrical fied. Hence ions that have a specific m/z will follow a stable course and reach the detector whereas those with an unstable course will collide with the quadrupole and not be detected. The RF and DC potential can be controlled so, only ions with specific m/z values will reach the detector. Due to its fast scanning rate and low cost the quadrupole is the most commonly used analyzer in GC-MS [Grob and Barry, 2004] [Poole, 2003] . The detector generates an electrical signal and counts the number of ions for each specific mass. This is plotted by the datasystem which creates a mass spectrum [Poole, 2003] [Grob and Barry, 2004].

2.7 CEN methodology for oil spill identification

The existing guideline for oil spill identification is called the CEN methodology (CEN/TR 15522-2:2012). It is a revised and improved technique which is based on the first edition of the CEN guidline from 2006 (CEN/TR 15522-2 Version 1), and the NORDTEST method from 1991 [NORDTEST, 1991] [CEN, 2012] The CEN methodolgy is a forensic tool for "*characterising and identifying the source of waterborne oils resulting from accidental spills or intentional discharges*" [CEN, 2012].It can be applied to waterborne oils, and samples of petrogenic origin that contains significant amounts of hydrocarbons with boiling points exceeding 200°C [Daling et al., 2002]. Examples include crude oils, light refined products and heavy refined products such as bunker oils [Stout and Wang, 2016]. The CEN methodology is divided in two

parts. The first part presents a procedure for sampling, transport and storage of oil, while the second part is concerned with analysis and data processing of results [Stout and Wang, 2016]. The second part is divided into three levels and is illustrated in a flow chart in figure 2.8.

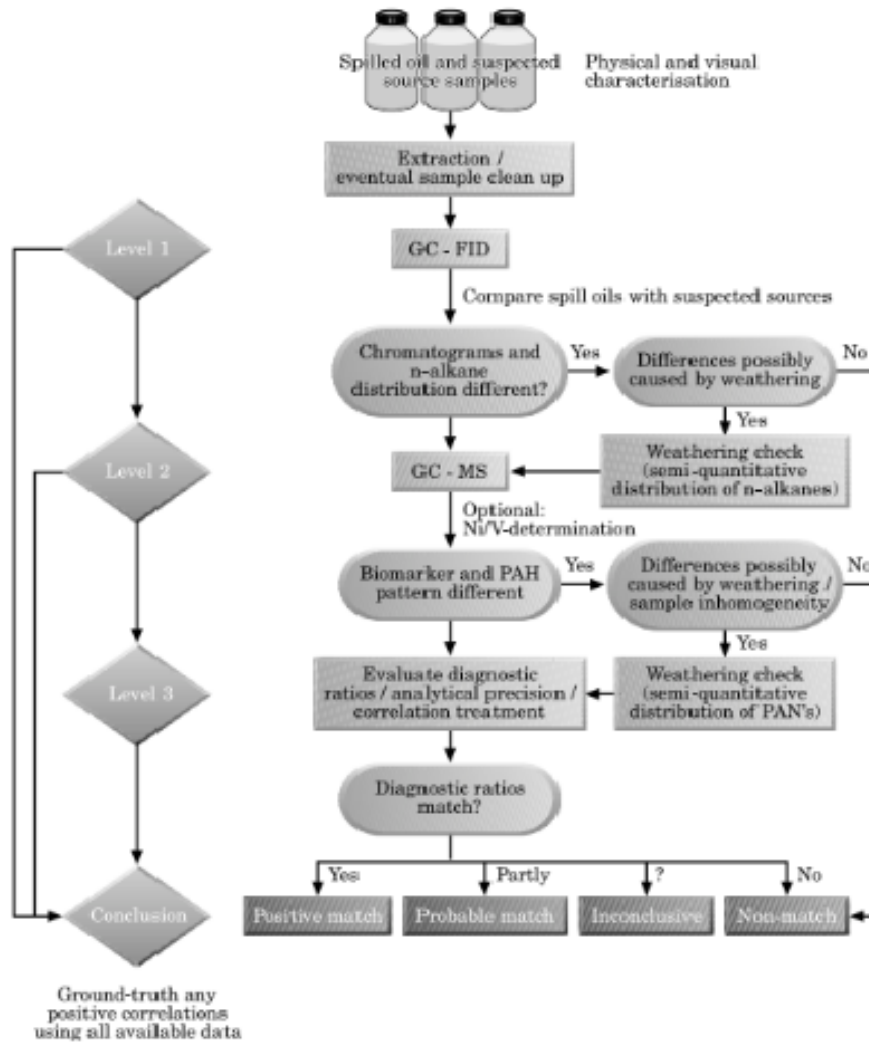


Figure 2.8: Flowchart for oil spill identification [Daling et al., 2002].

2.7.1 Level 1

The first level includes screening all samples by GC-FID. This is mainly to obtain a general overview of the samples, and is helpful for excluding samples that do not match a specific source sample, or exclude samples that turn out not to be oil. The chromatogram provides us with information regarding the distribution of hydrocarbons, type of components in our sample such as n-alkanes and isoprenoids, it gives us an indication about the boiling point range and the degree of weathering. An oil sample that has been analyzed by GC-FID will typically show several peaks of n-alkanes along the X-axis, see figure 2.9. Sometimes a hump under the alkanes is observed and this is described as an unresolved complex mixture (UCM). This means that there is a set of alkane peaks that cannot be separate by the chromatogram. There are two types of isoprenoids that are normally used for measuring the degree of weathering, namely pristane and phytane. Pristane is located close to nC17 and phytane close to nC 18. This is measured by calculating the diagnostic ratio between n-c17/pristane, n-C18/phythane and pristane/phythane [Fingas, 2010]. In cases where you observe that a spill sample and a source sample have similar chromatograms and you suspect that the difference between them is caused by weathering, a “weathering check” can be obtained. This can be done qualitatively by overlaying two chromatograms and compare them, or numerically by normalizing peaks [Daling et al., 2002] [Wang and Stout, 2010].

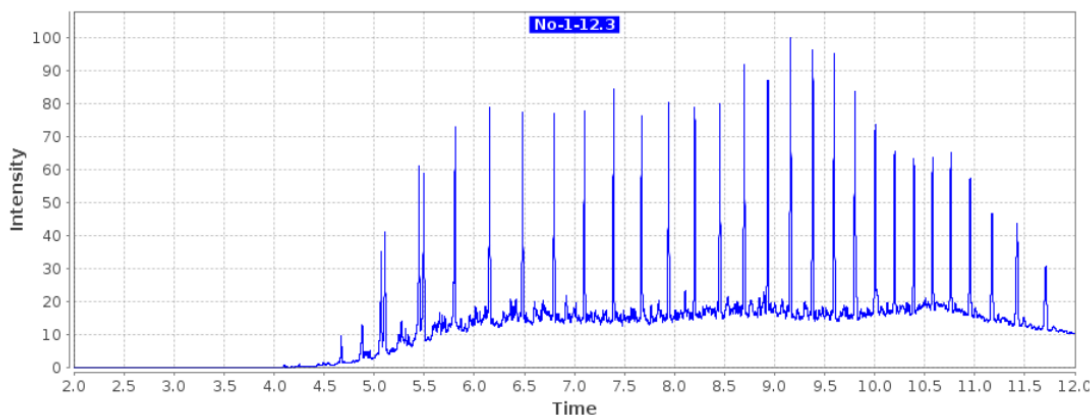


Figure 2.9: Example of a GC-FID chromatogram (sample 2015-0839) displaying n-alkane distribution, and a broad and flat UCM hump.

2.7.2 Level 2

Once the general screening by GC-FID is completed and obvious non-matching samples are discarded, the remaining samples are further analyzed by GC-MS in selected ion-monitoring mode. From a set of selected mass-charge values, we obtain chromatograms that show peaks representing various biomarkers and polycyclic aromatic hydrocarbons (PAH) that can be used in the process of identifying the source of an oil spill. In addition, the chromatograms provides us with helpful information for evaluating which samples are crude oil and bunker oil and if the crude oil is from the north sea or not. This is done by inspecting the chromatogram. [Wang and Stout, 2010] [CEN, 2012].

2.7.3 Level 3

In level three, diagnostic ratios are calculated between two PAHs or biomarkers, preferably within the same m/z value. Diagnostic ratios represents a useful tool for comparing oil samples. They induce a self-normalizing effect on the data. Variations caused by the instrument, operator and matrix effects are minimized [Wang and Stout, 2010]. So when ratios are compared they reflect differences of the biomarker distribution between samples.

For PAHs it is recommended to calculate diagnostic ratios from peak areas, and for biomarkers it is recommended to calculate the diagnostic ratio from peak heights [CEN, 2012]. The diagnostic ratio is calculated by the formula, where A and B represents the peak height and peak area respectively (equation 2.1);

$$DR = A/B \quad (2.1)$$

Diagnostic ratios from a spill sample are typically compared with diagnostic ratios from a reference sample. In order to compare the two samples and evaluate if the diagnostic ratio are identical, univariate statistical analysis is used. The comparison is based upon the repeatability limit (r). A limit, r , is defined as the repeatability limit. The repeatability limit is calculated from a relative standard deviation which is set to 5 %. If you choose a 95 % confidence interval then the repeatability limit is 14 %. This means that when the repeatability limit exceeds 14 % we are 95 % certain that the calculated ratio of the samples are different [CEN, 2012].

A positive match can be concluded if every diagnostic ratio are below the repeatability limit and the only differences between them are due to weathering. A probable match is concluded if only a few diagnostic ratios are above the repeatability limit. The chromatographic pattern are similar to each other and differences are

due to weathering, or contamination [Daling et al., 2002] [Wang and Stout, 2010]. An inconclusive match means that there are some similarities between the samples, but also too many differences and too many DRs above the repeatability limit which may or may not be due to weathering or be due to different samples. A non match is concluded if the chromatographic patterns are different from each other and this cannot be due to contamination or weathering. In addition many DRs are above the repeatability limit.

2.8 Multivariate Statistics

Multivariate statistics is a branch of statistical analysis that enables you to look at multiple variables simultaneously. This makes it possible to identify underlying information in a large dataset which would otherwise be difficult with the conventional univariate statistics [Esbensen et al., 2002]. The components are stored in a matrix X , and consist of n objects and p variables. The variables are listed in the columns and describes the property we want to measure, (e.g. Temperature, pH) and the objects are listed in the rows [Martens and Martens, 2001].

2.8.1 Preprocessing

Preprocessing is performed prior to analysis to improve the data and to obtain a dataset that is reliable for further analysis. The preprocessing methods used in this thesis are mentioned here.

Blank sample

A blank sample can give information about the background noise of a method that has been applied. Different type of blanks can be measured. An example is a method blank which is an analyte-free sample that runs through the same experimental procedure as the actual sample.

A blank enables calculation of a detection limit, which is the lowest concentration or signal that can be detected by a system [Little, 2016]. Limit of quantification (LOQ) is the smallest concentration that can be measured with reasonable reliability by the method [Mocak et al., 1997]. Limit of detection can be calculated by the formula:

$$LOQ = \text{meanblank} + 10 * STDblank \quad (2.2)$$

Centering

Centering is a technique performed to remove an offset in your data, which is the value of the point the variables fluctuate around. The mean of the variables is typically used as the offset. Centering thus changes the origin of your dataset, however the distance between the variables remains the same. This simplifies interpretation because the variation that now can be observed are between the variables only [van den Berg et al., 2006]. Mathematically this is done by subtracting the column average from every data point in the column as presented in equation 2.3. This is done for every column in a matrix:

$$\tilde{x}_{ij} = x_{ij} - \bar{x}_i \quad (2.3)$$

Scaling

Scaling is a technique that can be applied to a dataset in situations where there are variables of different magnitude or different types of variables. This is important to prevent variables with *high values* to completely dominate the model at the cost of variables with *small values*. There are different ways of scaling a dataset, a common type is by standardization, which means that variables are divided by their standard deviation. This gives variables the same variance, thus they influence the model equally [Esbensen et al., 2002] [van den Berg et al., 2006].

Standardization starts by centering data points in a column, before it is divided by the standard deviation of the column, see equation 2.4. This referred to as auto scaling.

$$\tilde{x}_{ij} = \frac{x_{ij}}{\bar{x}_i} \quad (2.4)$$

Drawbacks to this technique can be observed when variables of small values that should only have a minor effect on the model are scaled up. This is because the measurement error, which can be large for small values, will increase and influence the model [van den Berg et al., 2006].

Normality

Many univariate and multivariate techniques require that data are normally distributed, and exceptions from this can lead to results that are difficult to interpret [Hair et al., 2006]. A perfectly normally distributed dataset follows a Gaussian curve (or a bell shape). Unfortunately, this is rarely the case, especially in fields concerning geochemistry and environmental studies. Very often a skewed distribution,

often to the right, can be observed [Reimann and Filzmoser, 2000]. There are various tests to check for univariate normality, some of them being histogram plots, skewness coefficient, kurtosis coefficient, the shapiro-Wilks test and the Kolmogorov-Smirnov test [Hair et al., 2006].

The skewness coefficient displays a numerical value for the shape of the data. A positive value may indicate that data are skewed to the right and a negative value may indicate that data are skewed to the left. When a dataset is symmetric the skewness coefficient equals 0, a number which is not very realistic when dealing with real data. As a rule of thumb the coefficient should be between -1 to 1 in order to describe the data as symmetric or approximately [Remenyi et al., 2011].

The kurtosis coefficient gives a numerical value for the shape of the peak. If the value is larger than 1 this may indicate a sharper peak than the normal distribution and less observations can be found in the tails. If the value is less than 1 it may indicate a less sharp peak than the peak of a normal distribution and more observation can be found in the tails [Remenyi et al., 2011].

Transformations

Transformations can be applied to data that are not normally distributed. Data transformations replace the original values with transformed values to make the distribution more symmetric and suitable for interpretation. Different techniques exist [Bourman, 2009]. One frequently used technique is logtransformation of data, either by taking the natural log or log base 10 (see equation 2.5):

$$\log_{10}x_{ij} \tag{2.5}$$

2.8.2 Principal Component Analysis

In the field of multivariate statics there are various models that enables us to interpret the data we wish to make sense of. Principal component analysis is one of these models, and can be described as one of the most basic tools in the field of multivariate statistics. It was first addressed by Pearson, but also Fisher, Mckenzie, Wold and Hotelling are among the important contributors to this field [Pearson, 1901] [Wold et al., 1987].

PCA aims to express the most relevant information that is contained in a dataset by creating a new set of variables called principal components. These new set of variables are obtained as linear combinations of the original variables. A new coordinate system is formed, the new axes being the principal components. Mathematically the

principal components are calculated from a datamatrix X . The datamatrix is a product of a scores matrix T , a loading matrix P' and a residual matrix E , see equation 2.6.

$$X = TP' + E \quad (2.6)$$

The Scores matrix and loading matrix explain the most important variation in X . The first principal component is always positioned in the direction of maximum variance. The second principal component is orthogonal to the first PC, and describes the next largest variance, and so on. The residual matrix is the unexplained variance and is often noise. The goal is to obtain as much information as possible from a minimum number of principal components. In this way the noise is not a part of the evaluation, however it is important not to miss any important principal components that might be of value [Wold et al., 1987]. The new axes can be visualized in various plots to enhance interpretation. A score plot maps the objects along two score vectors, typically vector 1 and vector 2. This plot shows the distribution of the samples and makes it possible to detect groups of objects and visualize objects that are outliers. Similarly a loading plot is map of the variables along two loading vectors. This plot shows which variables are important for the principal components [Alsberg, 2015]. Variables that are situated in the periphery along a principal component typically has a large contribution to this principal component. Likewise, a variable close to the center do not contribute much to this component. Variables close to each other are positively correlated, and variables located 180 degrees to each other are negatively correlated [Esbensen et al., 2002].

A biplot is a combination of a score plot and a loading plot. The biplot gives an overview over which variables are positively correlated to objects, or groups of objects. Objects that is located nearby a variable in the biplot, typically have high values for that variable. Objects 180 degree to a variable, have a low value for that variable [Esbensen et al., 2002] [Alsberg, 2015].

2.8.3 HCA

Hierarchical cluster analysis is a classification method that combines objects into clusters by various set of rules. This is presented in a dendogram which visualises the formation of clusters. Agglomerative clustering is one type of hierarcial clustering and, this method form clusters based on similarity between objects. The procedure first starts with each object as a separate cluster. In each step the two most similar clusters are grouped into one new cluster. In the next step, the distance between all clusters are updated, and the two most similar clusters are once again grouped into

one new cluster. This process is repeated until all samples are combined into one large cluster. The degree of similarity among objects can be measured by different techniques. The most common measure of similarity are distance measurements. In this category, the Euclidean distance is the most commonly used measurement, and is basically the length (of a straight line) between two objects when visualised graphically [Hair et al., 2006].

The distance between two objects can be calculated by a set of various linkage methods. The most common are single linkage, complete linkage, group average and Ward's method. Single linkage measures the minimum distance between single objects in each of the two clusters A and B. Complete linkage measures maximum distance between two objects in the two clusters A and B. Average linkage measures the average distance from all objects in one cluster A to all objects in another cluster B. Ward's method calculate the average distance from all objects in one cluster. The average distance is subtracted from each object and then squared. The sum of squares are then calculated for each cluster [Hair et al., 2006].

A possible draw back to this type of clustering, is that the procedure will always form clusters, even in cases where there is no structure in the data. This means that samples inside a cluster not necessarily are useful for the analyser. Different methods have been developed for selecting the optimal number of clusters, it is however up to the analyst to make the final conclusion [Hair et al., 2006].

2.8.4 PLS-DA

Partial least square regression is a multivariate calibration technique that consist of two matrices, where X contains the independent variables and Y contains the dependent variables. Typically the X-matrix consist of available measurements, similar to the X-matrix with variables and samples in PCA. The information in the Y-matrix is typically something which is difficult and expensive to measure directly. It is often believed to be a causation between the information of these [Geladi and Kowalski, 1986]. The model explains the maximum variation of the data as well as achieving maximum correlation between the X and Y dataset [Kumar and Mishra, 2015], see equations 2.7 and 2.8.

$$X = TP^T + E_X \quad (2.7)$$

$$Y = UQ^T + E_Y \quad (2.8)$$

T and U are the scores matrix of X and Y, P and Q are the loading matrix of X and Y, E is the residual for X and Y. A common algorithm for solving this is called

NIPALS, which is an iterative algorithm that uses the score matrix of X to update the loading matrix of Y, and the score matrix of Y to update the loading matrix of X. This is performed until the results converge [Alsberg, 2015].

A PLS model is created in a calibration stage where a regression model is build with multivariate analysis techniques from known X- and Y-data. This model describe the relationship between the X- and Y-matrices, but it is unsure whether this will be true for another set of the same data. To test this, it is common to validate the regression model with a new set of data. When the model has been validated it is used to predict new Y-values from new X-measurements [Esbensen et al., 2002].

There are different validation methods available. One method called the Independent validation test works by using one dataset to calibrate the model. Then a completely new dataset are used for validating this calibration model. This is considered to be the best validation technique. However, in many cases there will not be enough samples to create an individual validation test model, especially if it is difficult to acquire data for the Y-matrix [Esbensen et al., 2002].

Often it may be necessary to re-use data. In full cross validation the dataset is split in two subsets. The first subset is used to calibrate the model and the second set is used to test the model. There is only one sample in the test set. The sample which was used in the test set is inserted into the model again, and a new sample is extracted to test the model. This is repeated until all samples have been used to test the model, and all samples have been included in the calibration model [Alsberg, 2015].

Partial least square discriminant analysis arise from PLS. In PLS-DA the Y-variable is a variable with information regarding class membership. It can typically be a binary variable, which means that samples belonging to one class is denoted with a number, for example 1, and samples belonging to the other class is denoted with number -1. The model is then used to classify unknown samples into one of the categories [Esbensen et al., 2002].

2.9 Computerized oil spill identification database

The *computerized oil spill identification* (COSIWeb) is an online database containing raw GC-FID and GC-MS chromatograms of oil samples. The database contain samples from accidental oil spills both of known and unknown origin, crude oils, bunker oils, fresh samples and weathered samples from different countries. It is available through any browser and allows the user to upload GC-FID and GC-MS chromatograms of oil samples in to the database which is then available for other users of the program [Stout and Wang, 2016].

The database automatically integrates components in the chromatogram, compare samples and express sample similarity. The user may then compare the uploaded chromatograms with other chromatograms in the database to investigate possible matches, which is listed in descending order. Similarity is defined by calculating the Pearson correlation coefficient and COSIWeb apply 26 diagnostic ratios (from the CEN methodology) as a basis for calculating correlation coefficients between samples, and provide additional information such as diagnostic ratios, GC overlays and partial weathering plots. This makes it possible to reach a conclusion which is in compliance with the CEN methodology [Stout and Wang, 2016] [COSIWeb, 2016].

If the difference between two diagnostic ratios are below the repeatability limit given by the CEN guideline (14%) these are colored in green. In Cosiweb this corresponds to 7%. If they differ no more than two times the standard deviation (7%-10%) the ratios are colored in yellow. Above this they are colored in red [COSIWeb, 2016].

Chapter 3

Materials and Methods

This section provides information about sample material, sample location, experimental methods as well as data treatment and statistical methods.

3.1 Sample material

Sample material was supplied by SINTEF Sealab in Trondheim, in addition to samples collected by the author. Approximately 400 samples have been collected by students through the course *KJ3050-Marine Organic Environmental Chemistry* and by master students. 112 of these samples have been analyzed by GC-FID and GC-MS, during a time period of 2011-2016, and are included in this thesis. Samples from 2010 have not been a part of this project.

3.2 Sample material from SINTEF

Students have collected samples from 18 different islands along the period 2010-2015. All samples have been inspected properly in the laboratory by students, with the help of experienced staff members. Only a subset of the total sample (approximately 400) size have been further analyzed by GC-FID and GC-MS. The sample subset have been selected based on different properties such as sample location, size of oil sample, smell, stickiness, tar balls etc. In some instances, both the center of a sample and the periphery of the sample have been analyzed, for example if students wanted to investigate if this was a mixture of different oil samples. It should also be mentioned that some samples collected from the field trips, turned out not to be oil when inspected in the laboratory. Appendix A provides an overview over sample

Table 3.1: Sample overview, showing the number of samples included in this thesis

Year	Total size	Collected by	Project size
2011	82	KJ3050	17
2012	65	KJ3050	17
2013	77	KJ3050	22
2014	58	KJ3050	21
2015	67	KJ3050	14
2012	23	Master student (Stine Henriksen)	6
2015	25	Master student (Marie Myrstad)	15
Total	403		112

number, location and description of samples. This information is written by former students, as a part of their field report [KJ3050 students, 2011], [KJ3050 students, 2012] [KJ3050 students, 2013] [KJ3050 students, 2014] [KJ3050 students, 2015].

Table 3.1 gives an overview over the number of samples collected each year, and number of samples analyzed by GC-FID and GC-MS. In 2014, samples from 2011 and 2012 were reanalyzed by a former master student [Vike, 2014], and the raw data from this analysis have been used for these two years. Of the samples size in 2012, 6 of these samples are not from the course KJ3050, but are a subset from samples collected at Sula by a former master student in 2012 [Henriksen, 2012].

3.2.1 Sample location

Oil samples have been collected from 18 different islands. These islands are Vesterkalven (2011), Storkalven (2011), Kunna (2011), Kya (2012), Frøya (2012), Blåskykøya (2013), Olabussøya (2013), Burøya (2013), Storfosna (2013), Vingleia (2014), Gård-søya (2014), Bordholmen (2014), Kråkvåg (2014, 2015), Vassøya (2015), Geitungan (2015), Likøya (2015), Gildklakken (2015) and Sula (2012,2015). Figure 3.1 indicate the approximate area of these islands.

3.3 Sample material from Sula 2015

Material for this thesis have also been collected by the author, and consist of 15 weathered oil samples from Sula. Sula is an island located at the coast of Trøndelag and is a part of Frøya municipality. It is one of the westernmost island on the coast of Trøndelag, thus it is exposed to waves and wind from the North Sea. The island



Figure 3.1: Approximate area of sample location.

have previously been studied by a master student in 2012. During this fieldtrip 23 samples of tarballs and oil were collected. These were mainly found on the west side of the island. Southeast of Sula was also inspected without any findings [Henriksen, 2012]. Based on this, the author was curious to inspect the island once more to search for new samples and investigate other sampling sites on the island.

The field trip took place 13th of may, 2015. Samples were collected on the north-east and north-west side of the island and both the inter-tidal and upper-tidal zones were investigated during low tide. The weather was sunny but windy, with a temperature of 7°C. In total, 25 samples were collected. Of these, 4 were described by the author as tarballs. These were found in typical “wreck bays”, which are areas where marine debris accumulates. The remaining samples were oil stranded on rock surfaces. All samples were described by their smell, size, viscosity and stickiness. Additional notes were taken if samples contained biological material (such as feathers, algae or moss) or other substances (such as plastic). A photo was taken on every

sample station and GPS coordinates recorded. All samples taken from Sula were characterized as either solid or semi-solid thus indicating that samples were quite weathered. See appendix A for a complete description of collected sample material from Sula. Pictures of oil two samples from Sula are shown in figure 3.2 and 3.3.



Figure 3.2: Sample 16 (2016-163) collected on the North-West side of the island

Samples were collected by a spoon and knife and stored in either aluminum containers with paper lids or small glass bottles (40 ml) with plastic lids. The spoon and knife were cleaned with paper towel between every sampling to avoid contamination. The samples were stored in a refrigerator (4 °C) at SINTEF until laboratory work started in February 2016.

3.3.1 Inspection of samples

The laboratory work was carried out in February 2016. All samples were inspected visually and properties that was evaluated during the excursion were evaluated in the laboratory once more. In general, samples that was classified as semi-solid in the



Figure 3.3: Sample 19 (2016-165) collected on the North-West side of the island

field were more viscous, sticky and shiny compared to notes taken at Sula. This did not come as a surprise since the temperature inside the laboratory was much higher compared to the temperature out in the field.

Samples that upon visual inspection in the laboratory were found to be something else than oil were eliminated from further analysis. This eliminated 8 samples. A majority of these samples were recognized to be burned plastic and some turned out to be soil covered in burned material. The plastic samples were extremely solid and impossible to cut through. An interesting observation is that all samples described as tarballs by the author during fieldwork turned out to be plastic. A tarball is a solid or semi-solid globule consisting of oil that forms when they drift at sea for a long period of time. They are highly weathered and are often found at shore together with other marine debris that reaches land.

Sample nr 21 and 22 were omitted from further analysis since they turned out to be extremely weathered. In total, 10 samples from the field trip are excluded from further laboratory work. The remaining samples were prepared for further analysis.

3.3.2 Experimental procedure

In this section, description of the experimental procedure applies to samples collected at Sula in 2015, but the same procedure was followed by students in KJ3050 and by master students Kristine Vike and Stine Henriksen. This is because all samples in this thesis are analyzed with application of the CEN procedure [CEN, 2012].

3.3.3 Experimental procedure of samples from Sula

Laboratory work was done at SINTEF Sealab in Trondheim. Sample preparation and solid phase extraction was carried out by the author, whereas GC-FID and GC-MS analyses were performed by Senior Engineer Marianne Unaas Rønsberg and Senior Engineer Inger Kjersti Almås.

Glass ware and glass wool were baked, meaning that the equipment had been placed in a baking oven with increasing temperature of 225 °C per hour until it reached 450 °C for three hours. All solvents were of analytical grade. A total of 15 samples were analysed.

3.3.4 Preparation of oil samples

A small portion of oil (approximately 0.1 g) was taken out of the original oil sample with a small knife and transferred to a clear glass vial (40 mL). Care was taken to include only a portion from the center of the sample to avoid contamination from other substances. Dichloromethane (DCM, 10 mL) of analytical grade was added to the glass vial containing the sample and stirred firmly. The knife was cleaned with DCM between every new sample and the tip of the knife was changed for every fourth sample. To ensure that the samples dissolved, they were left in room temperature for approximately 12 hours.

Dissolved sample (4 mL) was filtered through Pasteur pipettes into new marked glass vials (20 mL). The pipettes contained Bilsom cotton (glass wool) and 3-4 cm of sodium sulfate powder (Na_2SO_4) in order to eliminate water and unwanted particles from the sample. Samples were further diluted with DCM in order to achieve the right concentration for GC-FID analysis. The samples were clarified for analysis when the color of the samples resembled a dark brown color. Samples were then transferred to marked GC-vials (4 mL) and stored in a refrigerator, awaiting GC-FID analyses.

3.3.5 GC-FID conditions

The GC-FID analyses were carried out using an Agilent 6890N Gas Chromatograph equipped with an Agilent 7683B Series autosampler. Injection was done by splitless injection at 330 °C, with injection volume of 1 uL. The compounds were separated on a Zebron ZB-5 column (30 m long, 0.25 mm ID and 0.25 um film thickness). The oven temperature was initially 40 °C for 1 min then increased to 330 °C, at 6 °C/min and held at 330 °C for 15 min. Carrier gas consisted of helium at a flow rate of 2.5 mL/min. Detector temperature was 330 °C.

3.3.6 Solid phase extraction

In order to remove interference compounds prior to GC-MS analyses, a solid phase extraction purification step (SPE) was conducted.

Prior to applying SPE, a solvent substitution from DCM to hexane was necessary. This was done by transferring the liquid samples from the glass vials (approximately 1 mL) to collection tubes. Hexane was added, and the tubes placed in a heating block of max 37 °C in a tmhosphere of nitrogen (N₂, 0.5 bar) to let DCM evaporate until approximately 1 mL was left. Hexane was added to the test tubes once more and the process repeated, ensuring that the remaining DCM evaporated.

The SPE device and collection tubes were rinsed with DCM prior to extraction. Silica columns were placed on top of the SPE device and conditioned with 3 mL of hexane. Since the columns never were to be left dry, a few drops of hexane was always above the column. Collection test tubes were marked with ID and placed in the SPE device. Sample (0.5 mL in hexane) was applied to the column by help of a pasteur pipette and the collection tubes were rinsed with hexane (0.5 mL) tree times making sure that the entire sample was transferred from the tube to the column. Samples eluted through the column by adding hexane (3 x 2 mL) and by vacuum pressure.

The solution was dried under N₂ stream (0.5 bar) until 1 mL was left in the collection tube. The fraction was transferred from collection tubes to GC-vials. To make sure that the whole sample would transfer to the GC-vials the tubes were rinsed with hexane and transferred to the vials until the GC-vials were almost full (4 mL). The volume were reduced to approximately 1mL under N₂ stream.

3.3.7 GC-MS conditions

The GC-MS analyses were performed using a an Agilent 6890N Gas Chromatograph equipped with an Agilent 5975B quadrupole mass-selective detector (MSD; ionsource: 230 °C). Injection was done by splitless injection at 330 °C, with injection volume of

1 μL . The mass spectrometer was employed in electron ionization (EI) mode with an ionization energy of 70 eV. The compounds were separated on a Zebtron z-5ms db-5 column (30 m long, 0.25mm ID and 0.25 μm film thickness). The carrier gas was maintained at a constant helium flow of 1.1 ml/min. The oven temperature was set to 42 $^{\circ}\text{C}$ for 2.30 min and then increased to 5.5 $^{\circ}\text{C}/\text{min}$ to 330 $^{\circ}\text{C}$ and held at 330 $^{\circ}\text{C}$ for 10 min. Selected ion monitoring mode was used during analysis, targeting 72 biomarkers and PAH compounds. A list of these compounds together with their m/z values is presented in appendix B.

3.4 Data treatment

This section explains how data have been prepared prior to statistical analysis and gives an overview of the statistical methods that have been applied to the data.

3.4.1 Integration

Raw data from GC-MS analyses was uploaded and integrated in Chemstation. Integration was done manually by the author for each sample which comprise 112 samples. For each sample, 72 target biomarkers and PAH components were integrated if they were present in the chromatogram. Peak height was integrated for biomarkers and peak area for the PAH components as recommended by the CEN methodology [CEN, 2012]. Peaks in the chromatogram were identified by comparing them to a *SINTEF oil mixture*. This mixture contains significant levels of all peaks recommended by the CEN guideline. To control that the biomarker and PAH components were integrated correctly, the retention times for all samples were uploaded in Excel and compared with each other. If any deviation was observed, these components were inspected once more and corrected if necessary.

3.4.2 Noise

Noise in the data that was caused by sample preparation and instrument analysis was estimated by preparing a method blank. Three method blanks containing hexane were prepared, by running them through Pasteur pipettes with cotton wool and solid phase extraction prior to GC-MS analysis. Limit of quantification (LOQ) was calculated and subtracted from the dataset prior to calculating diagnostic ratios, by equation 2.2. This resulted in very low values, and eliminated few peaks. According to the CEN methodology, peaks with signal to noise ratio (S/N) > 3 to 5 need to be eliminated for comparing diagnostic ratios. Signal is the peak height, and noise

is the value from peak-to-peak around the signal. Due to much compound noise in the data, which was not detected by the method blank, the CEN procedure was also applied for noise determination.

3.4.3 Inspection of chromatograms

Prior to the use of more advanced statistical methods GC-FID chromatograms and selected ion mass chromatograms were assessed, using a method inspired by level 1 and level 2 in the CEN methodology. This was done to evaluate weathering and if possible, categorize the samples as a crude or bunker oil. GC-FID was used to get an indication of the degree of weathering. GC-FID chromatograms were evaluated by inspecting the boiling point area, analyzing the n-alkane pattern and unresolved complex mixture hump, see appendix D .

Selected ion mass chromatograms were used to identify biomarkers and PAH components that are characteristic for differentiating between crude oil and bunker oil, and are able to separate a crude oil that do not originate from the north sea, and crude oil from the north sea. These biomarkers and PAH components are also known to be resistant to weathering, which is crucial when doing a visual inspection on weathered sample material. The ion fragmentogram displaying methyl-phenanthrenes and methylanthracene (m/z 192) was used to characterize the oil as either crude or bunker. The oil was characterized as a bunker oil if the chromatogram displayed a distinct methylanthracene peak (MA) and if first pair of peaks (3-methyl-and-2-methylphenanthrenes), were clearly more abundant than the second pair of peaks (9-/4 and 1-methylphenanthrenes). In the reverse case, namely that the second pair of peak was more abundant than the second pair of peaks and no MA could be identified, it was characterized as a crude oil. In doubtful cases, for example if the abundance of doublets were more or less the same, and no MA peak could be identified, the sample was characterized as unknown. Figure 3.4 shows the ion chromatogram (m/z 192) belonging to an oil sample characterized as crude oil. Figure 3.5 shows the ion chromatogram (m/z 192) belonging to an oil sample characterized as a bunker oil.

The Retene peak (m/z 234) is another characteristic PAH for crude oil, since this is typically lost in the refinery process [CEN, 2012]. Figure 3.6 displays the ion chromatogram of an oil sample with a retene peak. The ion fragmentogram displaying biomarkers Oleanane(30 O) and Gammacerance (30G) (m/z 191) were used to characterize crude oils that did not originate from the north sea, since these usually are absent in north sea crude oil [Liv-Guri Faksness, 2002]. Figure 3.7 display the ion chromatogram of an oil sample with these two peaks.

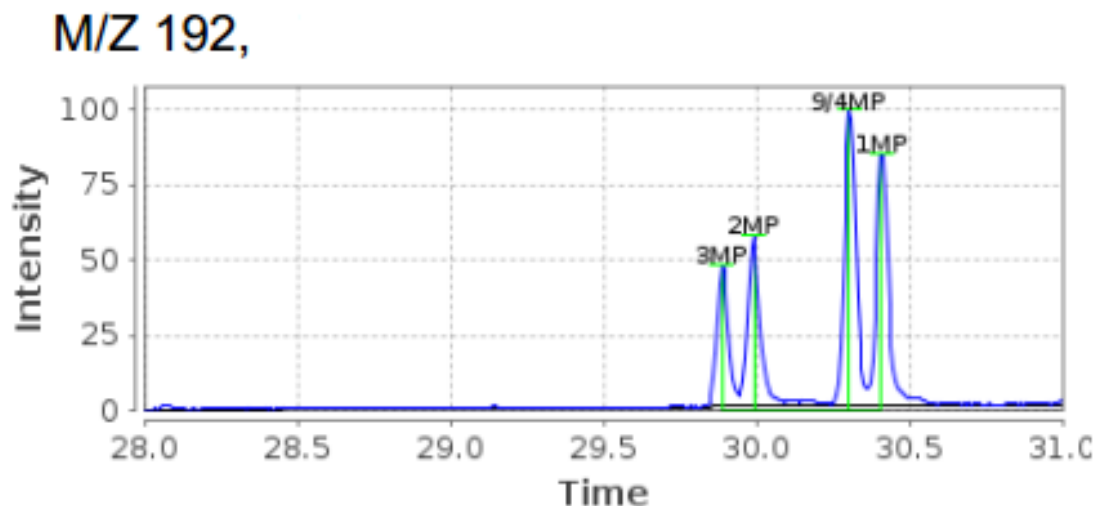


Figure 3.4: Ion chromatogram (m/z 192). This is an example of a sample characterized as a crude oil due to high abundance of 9-/4 and 1-methylphenanthrenes compared to 3-methyl-and-2-methylphenanthrenes

Samples were characterized into one of the following categories "Bunker oil", "Crude oil", "Non north sea crude oil" or "Unknown".

3.4.4 Diagnostic Ratios

Selected diagnostic ratios are listed in appendix C and follows recommendations provided by the CEN methodology. The CEN methodology includes ratios that are valuable for most oil spill identifications. The sesquiterpanes, however were not included since these are known to be affected by weathering [CEN, 2012].

3.4.5 Statistical analysis

Statistical analysis were performed using Excel (2016), Unscrambler (version 10.3) and R studio (version 3.3.1). Raw data from Chemstation was imported to excel and different spreadsheets were created to calculate diagnostic ratios and to calculate univariate statistics for each variable. *NA* was represented for empty data cells. Histogram plots, Kolmogorov-Smirnov tests and results from univariate statistics were evaluated for each variable prior to performing multivariate statistics. This

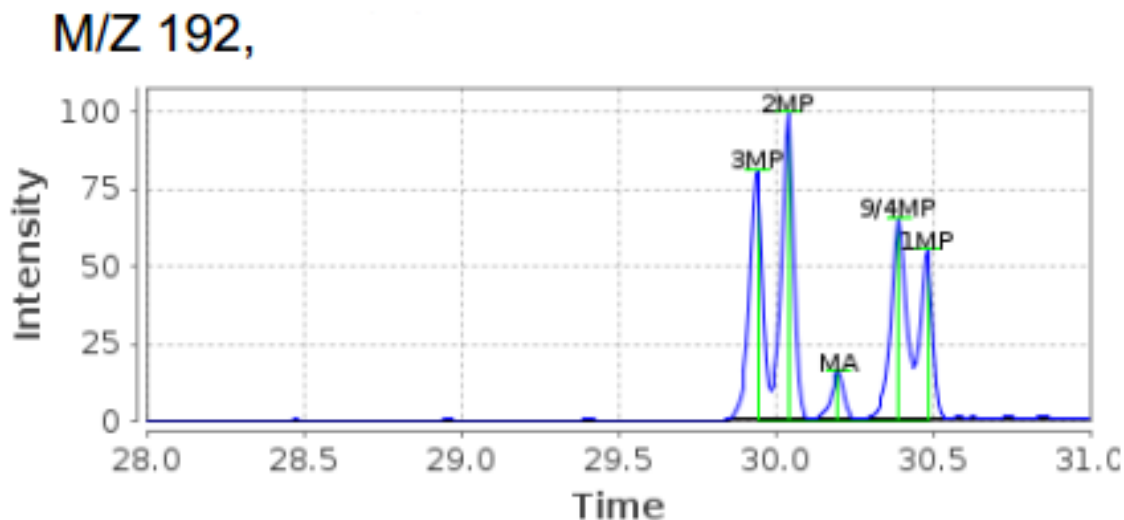


Figure 3.5: Ion chromatogram (m/z 192). This is an example of a sample characterized as a bunker oil due to high abundance of 3-methyl-and-2-methylphenanthrenes compared to 9-/4 and 1-methylphenanthrenes and a distinct methylanthtracene peak

is because both PCA and PLS require normally distributed data [Esbensen et al., 2002].

Diagnostic ratios were exported from Excel to Unscrambler for Principal Component Analysis (PCA) and Partial Least Square Discriminant Analysis (PLS-DA). m or $-0.9973E + 24$ are automatically represented for empty data cells in unscrambler, so missing values do not influence the results [Esbensen et al., 2002]. Diagnostic ratios were exported from Excel to R studio for Hierarchical Cluster Analysis (HCA). "NA" was represented for empty data cells. The statistical package *cluster* was used for plotting HCA, and the agglomerative Nesting (AGNES) function was applied since this function handles empty data cells.

3.4.6 COSIWeb

Selected samples were imported into COSIWeb with means of identifying samples. Samples from 2011, 2012 and selected samples from 2015 (2015-839, 2015-846, 2015-881, 2015-844) was already imported into the database. Samples from 2014 and 2016 (Sula) were imported by the author.

M/Z 234,

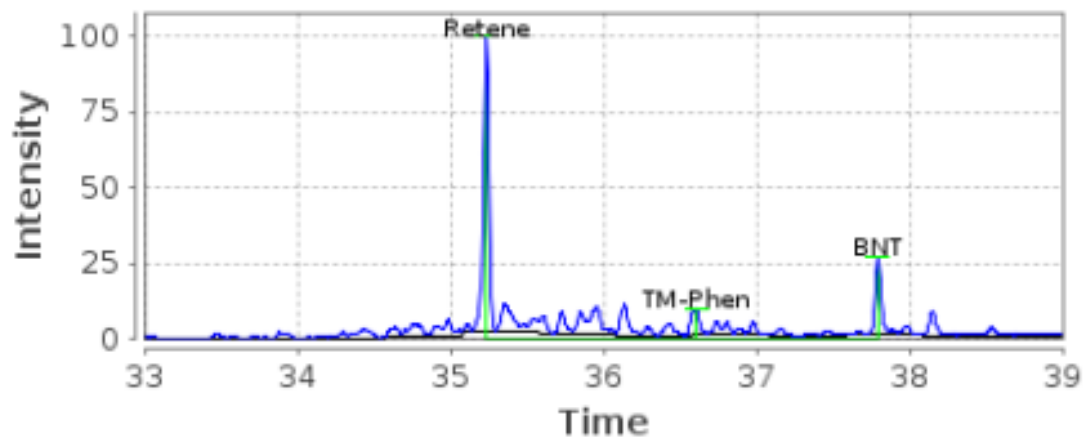


Figure 3.6: Ion chromatogram (m/z 234) displaying a retene peak.

M/Z 191, No-1-4.7

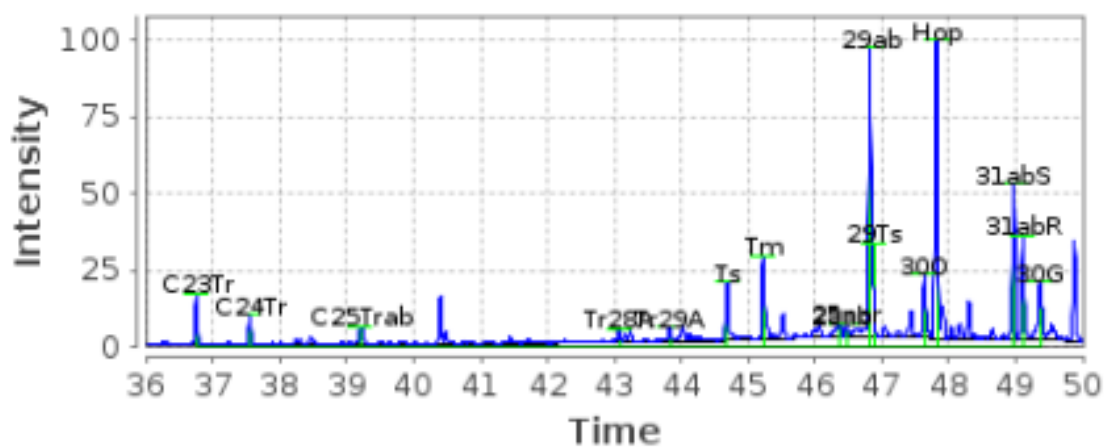


Figure 3.7: Ion chromatogram (m/z 191) displaying the biomarkers Oleanane (30 O) and Gammacerane (30 G)

Raw GC-FID and GC-MS chromatograms were converted to custom format, and

renamed according to the COSIWeb procedure. The files were imported into COSIWeb along with retention times for pristane and phytane. For each sample a short comment about sampling location, oil type and sampling date were provided. Integration of GC-MS chromatograms were provided by COSIWeb, however these were reevaluated by the author in case any doubtful integrations was observed. In these cases, the errors were corrected by using the manual function in COSIWeb.

When a sample is compared to existing samples in the database, COSIWeb automatically identifies samples who returns the highest correlation coefficient to that sample. Only samples with a correlation of 0.98 was inspected further. This number is based on experience by previous master student Kristine Vike [Vike, 2014], whom experienced that correlation should be at least 0.98 for similar samples. A cut of value of 0.98 was thus chosen to get a reasonable amount of potential matches. However it is still possible to have a match for lower values.

Chapter 4

Results

This section presents results from inspection of chromatograms, multivariate methods and results from the COSIWeb database.

4.1 CEN

Inspection of GC-FID and GC-MS chromatograms (level 1 and level 2) resulted in 51 samples being identified as crude oils and 20 samples identified as bunker oils. 40 samples were not identified based on the inspection of chromatograms. Among the samples identified as crude oils, 10 samples were identified as "non north sea crude oils". The results are displayed in a map with location, sample number and colors according to the classification. Samples colored in blue are characterized as "bunker oil", samples colored in green are characterized as "crude oil", samples colored in yellow are characterized as "non north sea crude oil" and samples in red are characterized as "unknown". See figure 4.1, figure 4.2, figure 4.3, figure 4.4, figure 4.5, figure 4.6, figure 4.7, figure 4.8, figure 4.9 and figure 4.10.

Results from the inspection of GC-FID chromatograms are not presented in this section, but can be found in appendix D with chromatograms and a description of the results from the chromatograms. The results from the GC-FID chromatograms gave an indication of the degree of weathering.

4.2 Pretreatment of raw data

Variables with more than 25% missing values were removed from the dataset, and this eliminated 15 variables from further analysis and leaves 29 variables. The removed

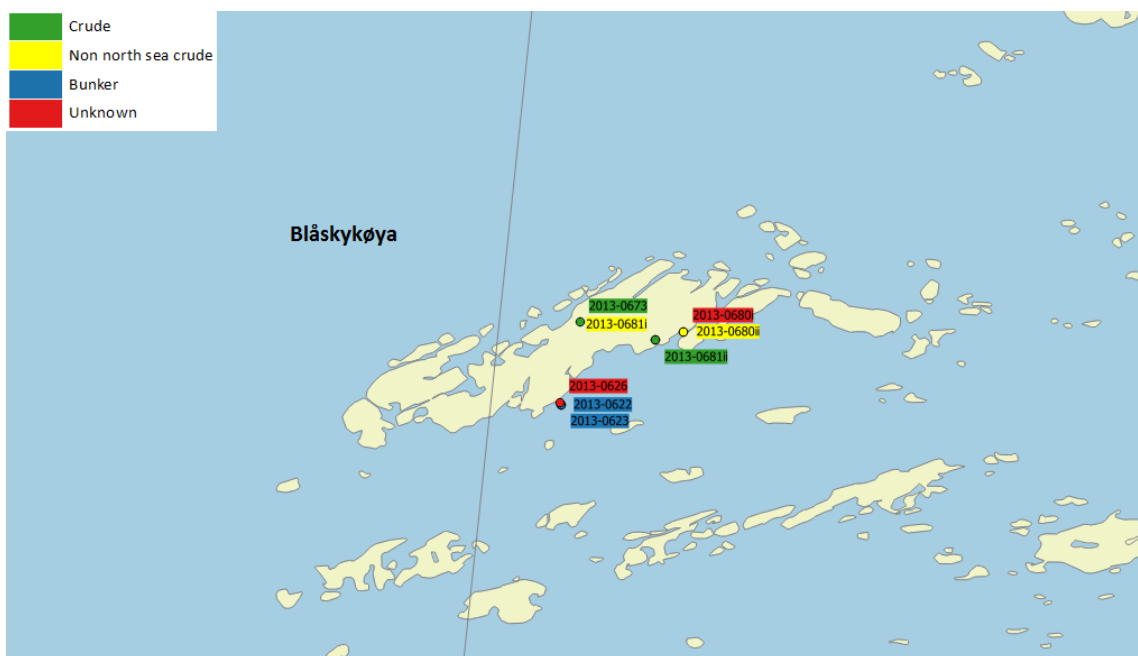


Figure 4.1: Blåskykøya

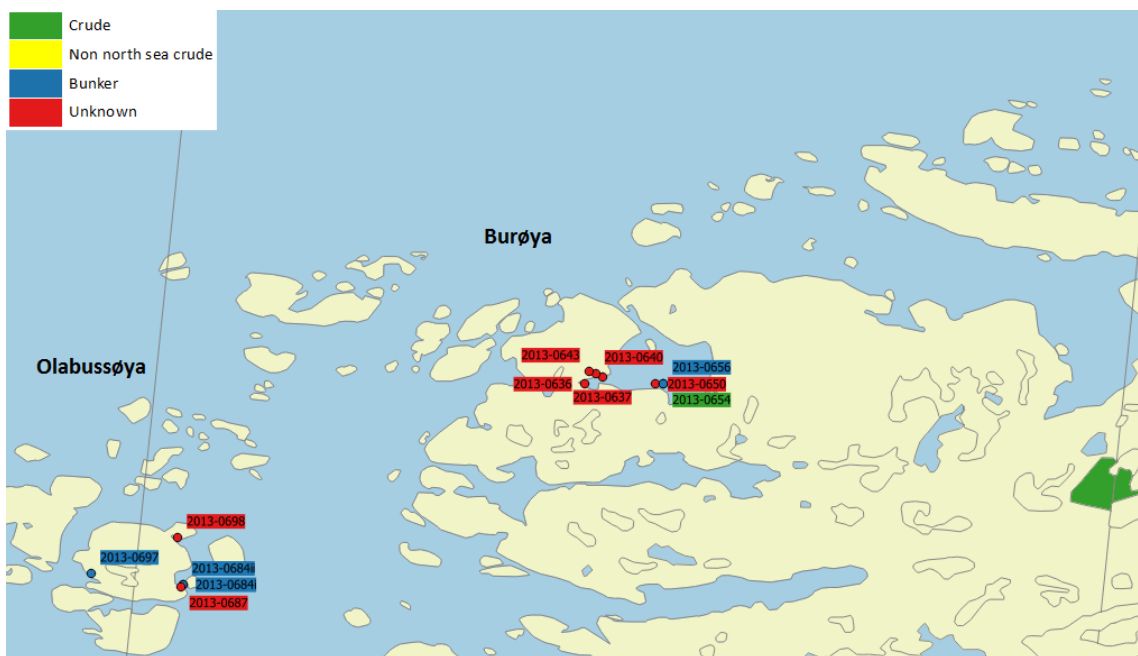


Figure 4.2: Olabussøya, Burøya

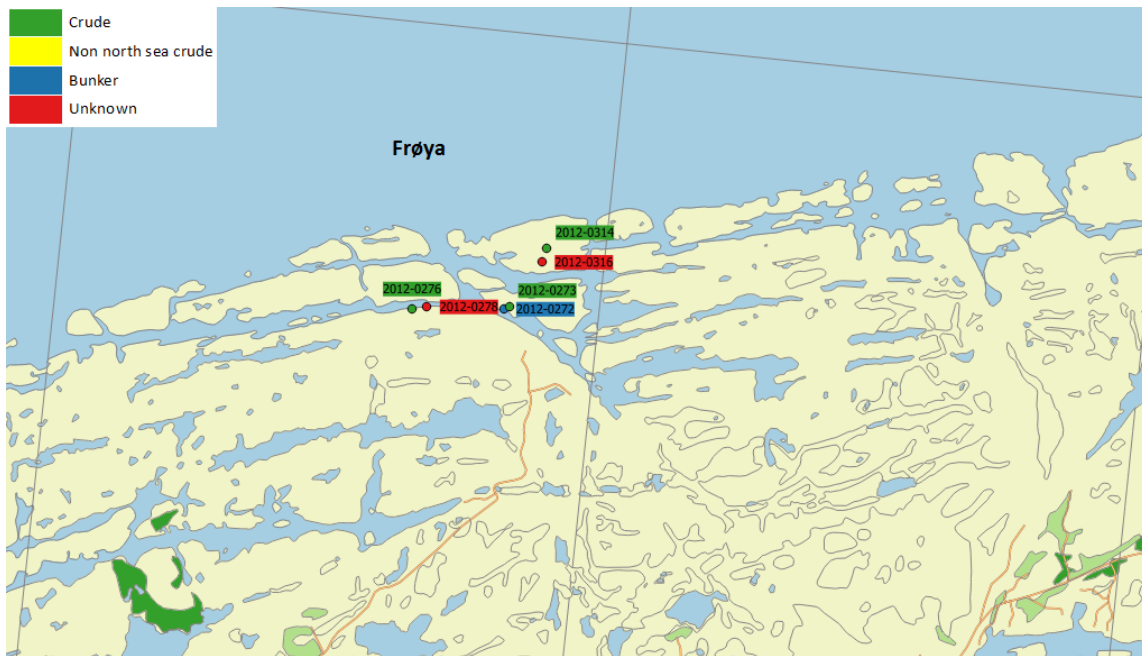


Figure 4.3: Frøya



Figure 4.4: Kya

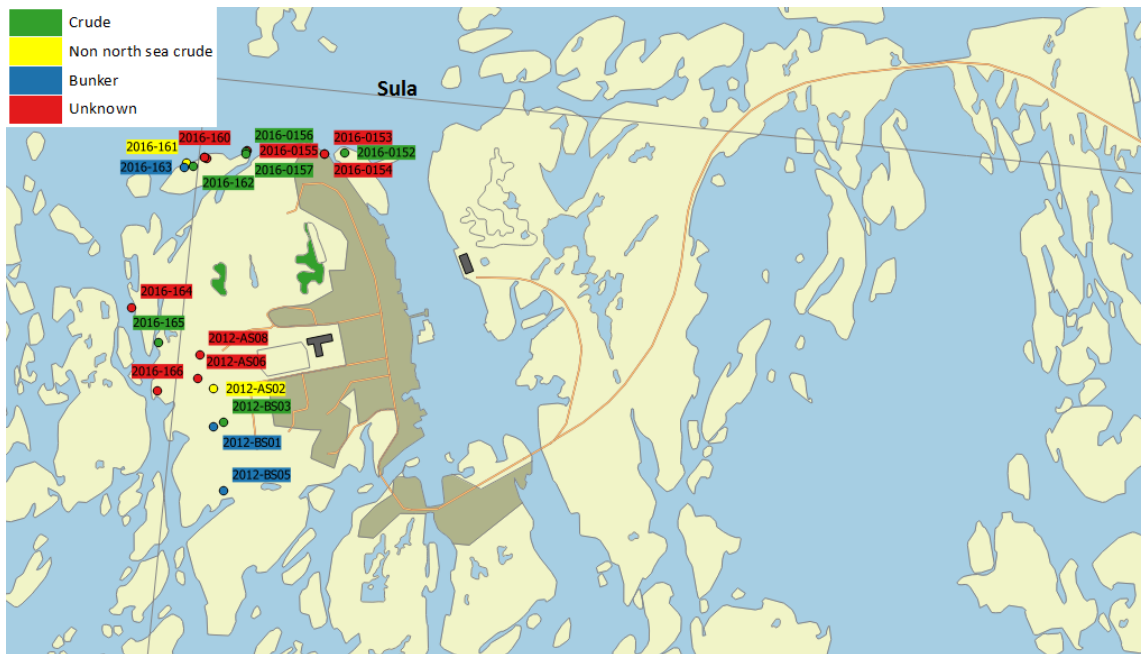


Figure 4.5: Sula

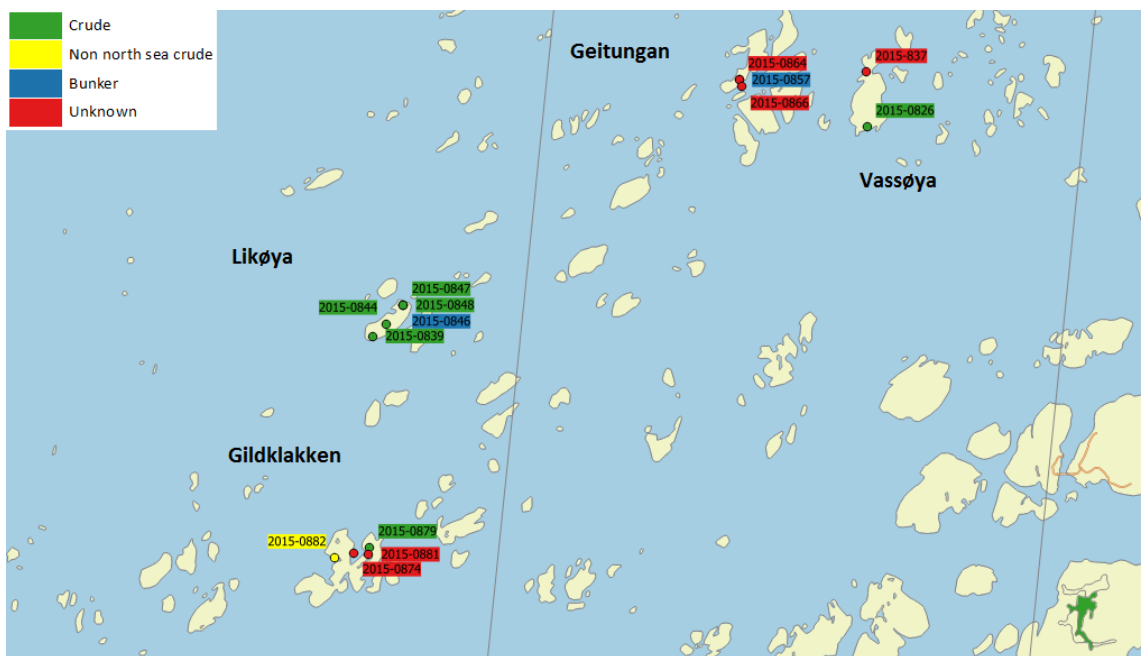


Figure 4.6: Gildklakken, Likøya, Geitungan, Vassøya

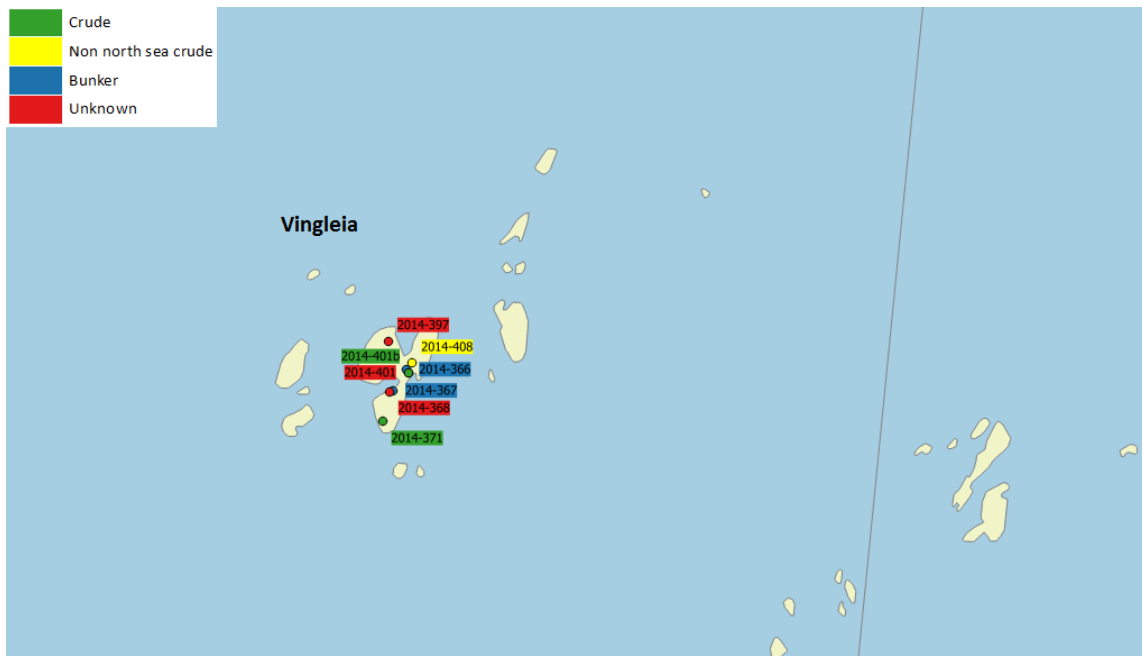


Figure 4.7: Vingleia

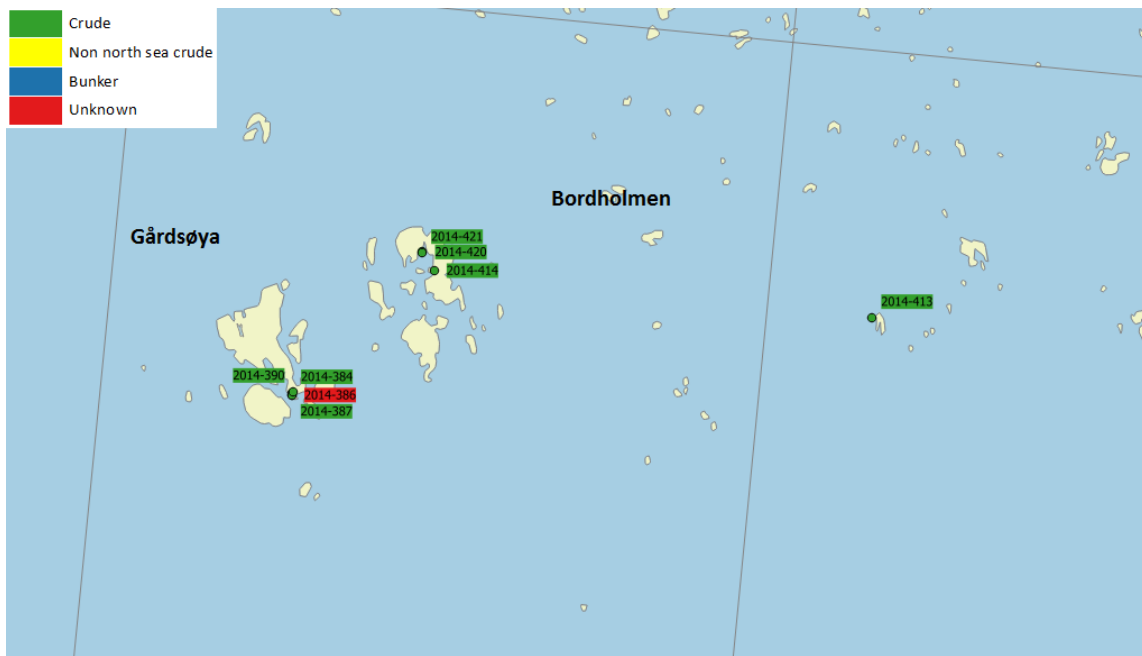


Figure 4.8: Bordholmen, Gårdsøya

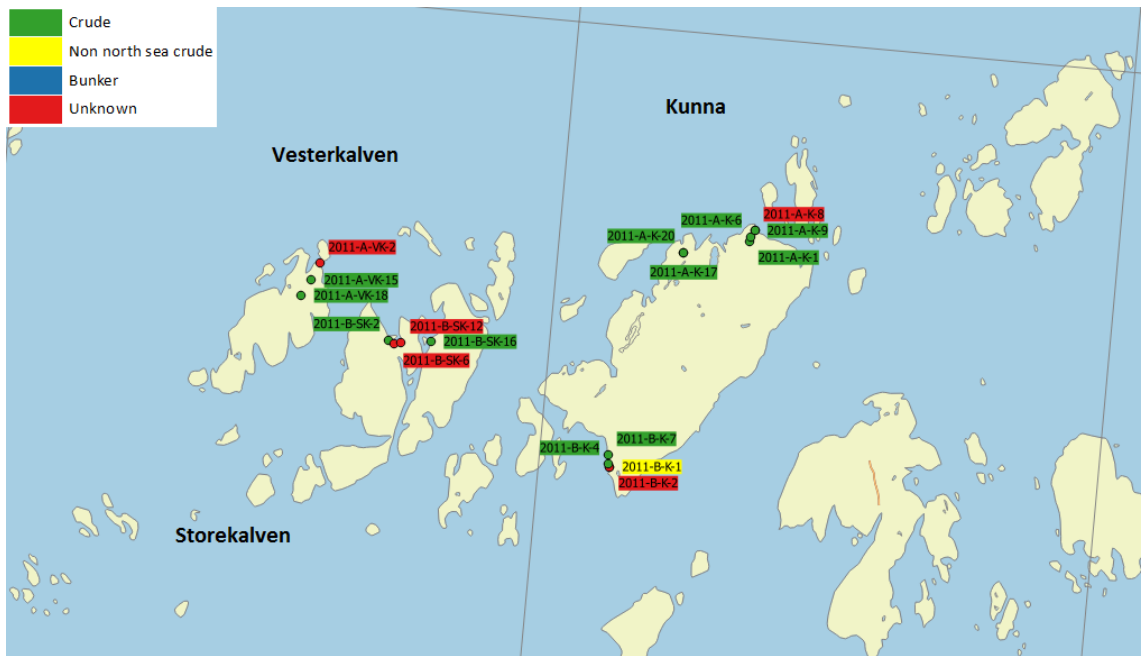


Figure 4.9: Kunna, Storekalven, Vesterkalven

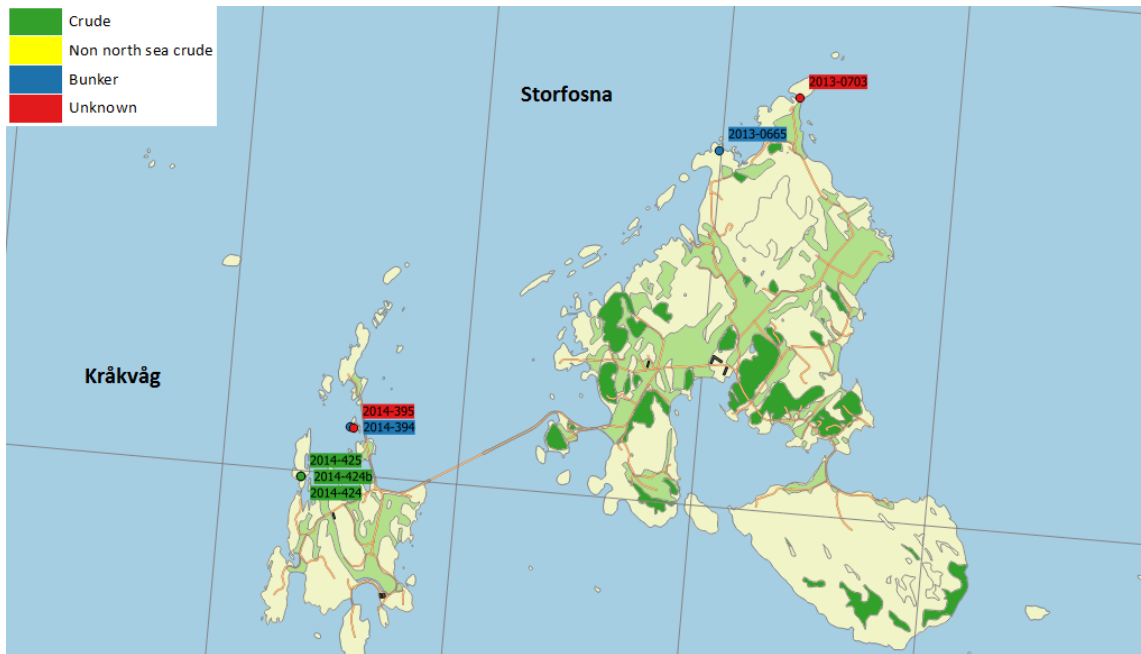


Figure 4.10: Kråkvåg, Storfosna

Table 4.1: Variables with more than 25% missing values. These are removed from further analysis

Ratio name		
DR-C17/pris	DR-B(a)F/4-Mpy	DR-30O/30ab
DR-C18/phy	DR-B(b+c)F/4-Mpy	DR-30G/30ab
DR- pris/phy	DR-Retene/ T-M-phen	DR-27dbR/27dbS
DR-2MFL/4-Mpy	DR-28ab/30ab	DR-MA/1-MP
DR-C28(22S)/30ab	DR-Retene/C4-Phe	DR-29TS/30ab

variables are presented in table 4.1. The decision of removing variables with more than 25% missing values was based on an article by Reimann and Filzmoser, and suggests that a variable can neither be normal or lognormal distributed if more than 25% of its values are missing [Reimann and Filzmoser, 2000].

Samples with more than 70 % missing values were not included in the dataset. This applies to three samples (2014-397, 2013-703 and 2013-698). The remaining 109 samples consist of 55 samples with no missing values, 40 samples with missing values between 10-20%, 10 samples with missing values between 20-40% and four samples with missing values in the range between 40-70%.

4.2.1 Descriptive Statistics

A summary of the descriptive statistics for each variable is presented in table 4.2. This includes the mean, median, skewness, kurtosis, standard deviation (SD) and relative standard deviation (RSD).

The table shows that the mean is slightly higher than the median which is an indicator of skewed data [Remenyi et al., 2011]. The results of the RSD reveals varying values, ranging from 15 % to 240 % between variables. This is as expected, since the diagnostic ratios are specifically used to explain differences between samples.

Table 1.2 shows that 13 variables have skewness coefficients between -1 and 1. 14 variables have skewness coefficients above 1, and two variables have skewness coefficients below -1. The kurtosis coefficient show 11 variables with kurtosis coefficient between -1 and 1. 19 variables show positive kurtosis coefficients above 1. In addition, histogram plots and Kolmogorov-Smirnov tests were assessed for each variable. The histogram plots and Kolmogorov-Smirnov tests for each variable are not presented, but the plots revealed that a large majority of the variables are not normally distributed. These variables are listed in table 4.3. The combined results

Table 4.2: Mean, median, skewness (before/after log transformation), kurtosis (before/after logtransformation), STD and RSD for each variable.

Variable	Mean	Median	STD	RSD	Skewness	Kurtosis
4-MD/1-MD	1,69	1,67	0,72	43%	1,04/-0,23	1,66/-0,08
2-MP/1-MP	1,08	1,03	0,48	44%	0,51/-0,14	-0,78/-0,93
2Mpy/4-Mpy	0,55	0,44	0,26	47%	1,46/ 0,48	3,50/-0,73
1Mpy/4-Mpy	0,46	0,43	0,16	34%	0,71/0,06	-0,17/-0,63
BNT/TM-phen	1,75	1,58	0,98	56%	1,22/-0,31	1,74/0,10
27Ts/30ab	0,16	0,15	0,05	30%	0,70	2,01
27Tm/30ab	0,26	0,27	0,10	39%	0,22/-0,68	-0,48/0,28
29ab/30ab	0,72	0,72	0,23	32%	0,37/-0,23	-0,27/-0,89
31abS/30ab	0,51	0,52	0,08	15%	0,10	-0,11
27bb/29bb	0,95	0,99	0,20	21%	-1,06/-2,63	2,71/9,34
SC26/RC26+SC27	0,27	0,25	0,11	39%	1,70/0,50	3,70/0,46
SC28/RC26+SC27	0,63	0,61	0,14	22%	1,03	1,81
RC27/RC26+SC27	0,64	0,63	0,10	15%	4,23/2,26	28,06/13,58
RC28/RC26+SC27	0,68	0,66	0,22	32%	4,48/1,57	29,99/6,71
C2-dbt/C2-phe	1,01	0,77	0,67	67%	1,27/-0,02	1,58/-0,62
C3-dbt/C3-phe	1,14	1,01	0,74	65%	1,29/-0,02	1,43/0,60
C23Tr/C2-PA	0,10	0,03	0,24	240%	5,18/0,39	31,04/0,36
29aaS/29aaR	0,92	0,91	0,14	15%	-0,10/-1,22	2,21/5,17
C20TA/C21TA	1,04	1,00	0,25	24%	0,79/-0,09	0,96/0,87
C21TA/RC26+SC27	0,29	0,26	0,19	66%	2,93/-1,42	11,59/6,02
Ts/Tm	0,70	0,67	0,33	47%	0,68/-0,09	-0,20/-0,96
30ba/30ab	0,11	0,10	0,03	29%	0,91/0,45	0,18/-0,69
C21TA/RC28TA	0,46	0,39	0,30	66%	2,78/-2,41	9,91/10,78
SC26TA/SC28TA	0,42	0,40	0,15	35%	1,17/0,06	2,08/0,20
RC27TA/RC28TA	0,99	0,94	0,20	20%	0,56/0,04	0,00/-0,20
C27BBSTER	0,52	0,53	0,10	19%	-1,49/-3,03	4,12/11,95
C28BBSTER	0,42	0,41	0,07	17%	0,60	0,96
C29BBSTER	0,59	0,56	0,12	21%	2,00/1,03	6,62/2,15
29bb/29aa	0,64	0,65	0,11	18%	-0,07	0,71

Table 4.3: Variables considered to be normally distributed.

Ratio name	
27Ts/30ab	SC28/RC26+SC27
31abS/30ab	C28BBSTER
29bb/29aa	

from the descriptive statistics, histogram plots and Kolmogorov-Smirnov tests reveals that only four variables can be considered to be normally distributed. Table 4.2 also display skewness and kurtosis coefficient after log transformation (equation 2.5) of variables not considered to be normally distributed. A majority of these coefficients decrease after transformation has been applied, but there are some that increases as well.

4.3 Multivariate data analysis

Multivariate data analysis (MVA) was applied the dataset to identify trends and structures in the data. Three different methods were tested; Principal component analysis (PCA), Partial least square discriminant analysis (PLS-DA) and Hierarchical cluster analysis (HCA). This section aims to present the most relevant findings for each multivariate technique.

Prior to all multivariate data analysis, data was mean-centered and scaled. Samples that was not considered to be normally distributed were log transformed. This includes all variables except those listed in table 4.3.

4.3.1 Principal component analysis

Principal component analysis was first applied to the data. This was done to observe groups in the dataset, identify samples with the same chemical composition and remove possible outliers.

Outliers

PCA applied to mean-centered, scaled and log transformed data resulted in the following score plot, see figure 4.11.

The score plot in 4.11 shows that PC1 and PC2 explains 38% of the variance. It is quite evident that sample 2013-0684ii and sample 2013-0684i are located away

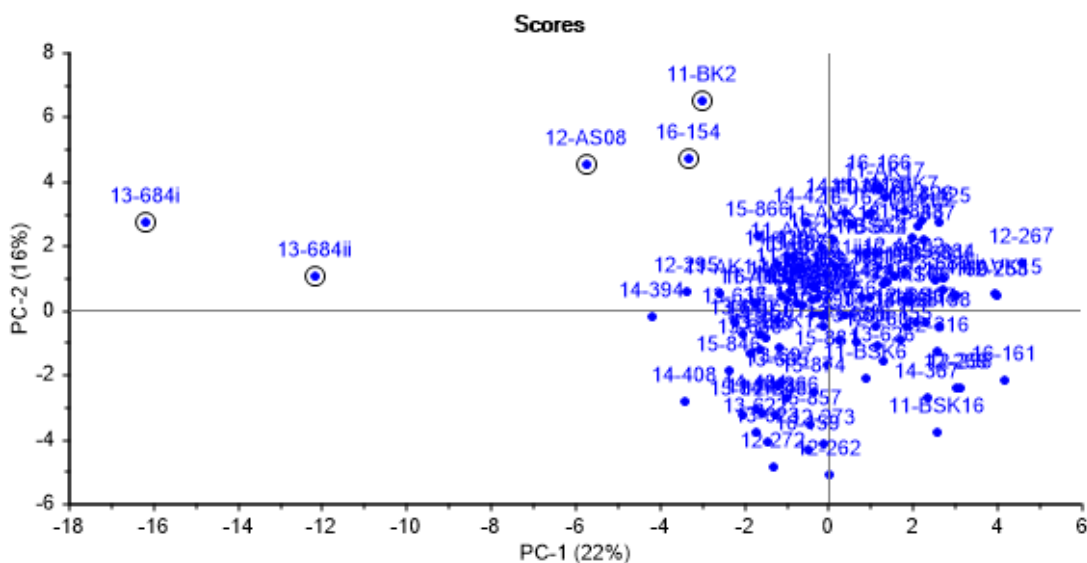


Figure 4.11: First score plot of mean-centered, scaled and logtransformed oil samples.

from the other samples. PC1 are almost solely used to explain the difference between the two samples and the main group of samples. Along the positive axis of PC2 there are three samples (12-AS08, 16-154 and 11-BK2) separated from the main group.

An influence plot was applied to the dataset to measure how much each sample affected the model, see figure Figure 4.12. An influence plot display leverage values on the X-axis and residual values on the Y-axis. A hotelling's T^2 test versus a F-residual test was applied with a 95 % confidence line. What can be observed from the model is that there is no samples exceeding the F-residual confidence line, but nine samples are exceeding the Hotelling's T^2 test confidence line. Five of these are the same samples observed as possible outliers in the score plot, and these are marked with black circles. Especially sample 13-0684ii, 13-0684i and 11-BK2 have increased leverage values and therefore increased influence on the model. This does not necessarily mean that they are outliers. On the contrary, they may be interesting samples with valuable information. However, they should be investigated further to rule out possible outliers [Esbensen et al., 2002]. Samples outside the confidence line were studied by inspecting the loading plot, raw data table, GC-FID and GC-MS chromatograms and students field report.

A table describing which samples were removed with a brief explanation are shown in table 4.4. In total 5 samples were removed from further analysis, mainly due to extreme weathering.

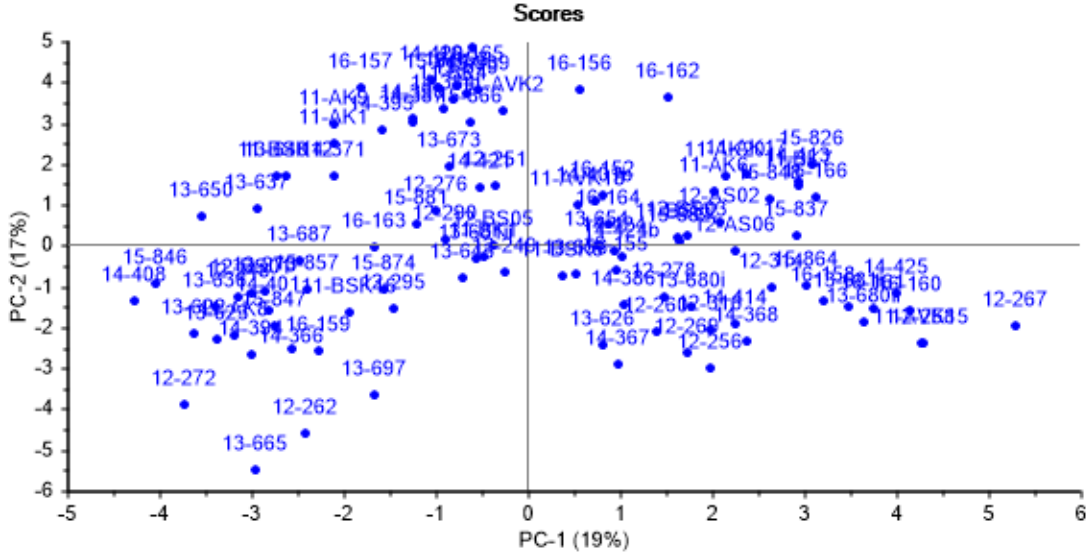


Figure 4.13: Score plot when outliers are removed.

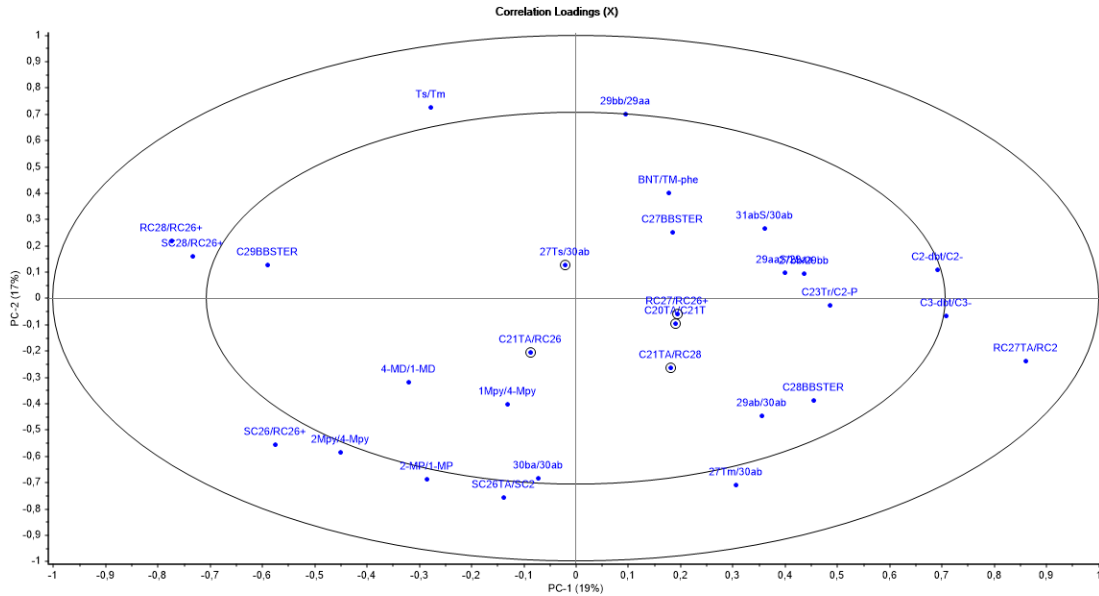


Figure 4.14: Correlation loading plot. Variables between the inner and outer ellipse indicate between 50% to 100% explained variance. Variables close to the center explain very little variance.

Table 4.5: Variables removed from after inspecting them in the correlation loading plot.

Removed variables	
C21TA/RC26+SC27	C20TA/C21TA
RC27/RC26+SC27	C27Ts/30ab
C21TA/RC28TA	

The loading plot in figure 4.14 gives a descriptive picture over variables that are important for each principal component (Esbensen, 2012 KILDE). The Correlation loading plot describe the importance of each variable. It shows two ellipses, the inner ellipse indicate 50% explained variance, whereas the outer ellipse indicate 100% explained variance [CAMO, 1998]. This means that variables between the two ellipses indicate that they have large loadings and are important for explaining the variance. Variables inside the inner ellipse are less important for explaining the variance. Variables with low contributions were removed from the correlation loadings plot to simplify interpretation. These are marked with black circles in the figure. The variables are listed in table 4.5

PCA when outliers are removed

The final PCA score plot and loading plot are presented in figure 4.15 and 4.16. The score plot for PC1 versus PC2 in figure 4.15 explains 42 % of the variance. The samples are evenly distributed along both principal components. Despite this there is no obvious separation of groups in the scores plot, except perhaps for small groups.

The loading plot in figure 4.16, indicate that there are still variables inside the first ellipse, which means that they are not equally important as the variables outside the first ellipse. It was decided to keep these because most of them are located close to the first ellipse and therefore contribute to a greater extent to the model, compared to those variables that were removed.

Hotelling T^2 -test in Figure 4.17 reveals six samples outside the confidence interval. Especially three of these (2013-687, 2014-394 and 2011-BSK6) have high leverage values. These were further inspected, but the author found no reason for removing these samples from further analysis.

The explained variance plot in figure 4.18 shows how much variance are explained by each principal component. The blue line represents calibration variance and the red line represents the validation variance. All seven principal components explain approximately 85 % of the variance, but the idea of PCA is to explain as much as

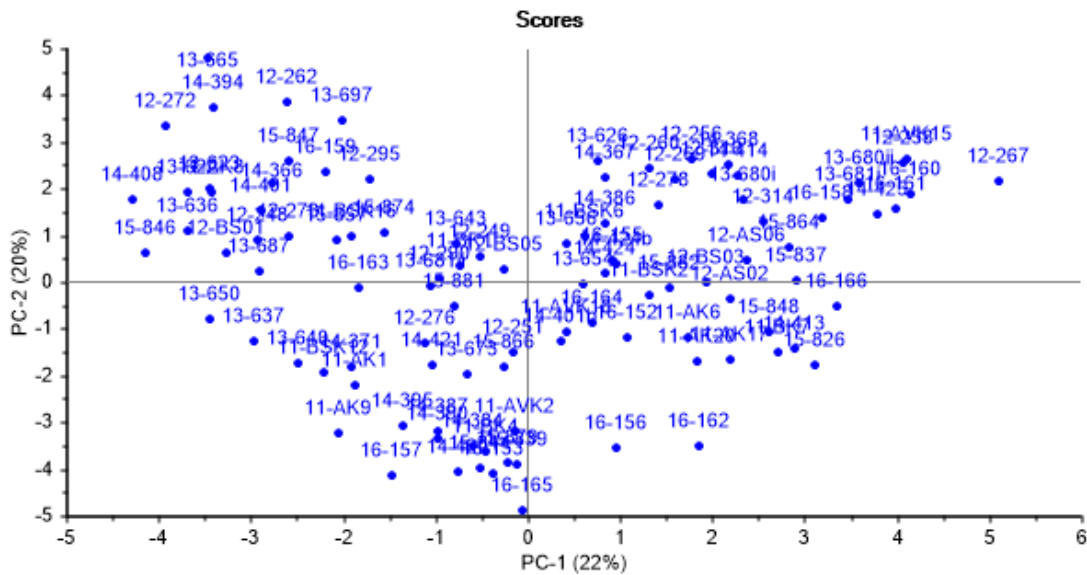


Figure 4.15: Final score plot after removing five outliers and five variables with little explained variance.

possible with as few principal components as possible, since the last PCs typically contain a lot of noise. PC1 and PC2 describes 22% and 20% of the variance respectively and PC3 and PC4 describes 14% and 10% of the variance. Since the score plot for PC1 and PC2 in figure 4.15 did not reveal any obvious distinct groups between the samples, it was desirable to include three principal components in the score plot (PC1, PC2 and PC3) to see if this gave any interesting groups in the score plot.

Identification of groups in the score plot was hence inspected by observing samples in a 3D score plot for PC1, PC2 and PC3 (see figure 4.19 and figure 4.20). Figure 4.19 and figure 4.20 represents different angles if the 3 D plot. This is to illustrate how some samples seem to belong to certain groups, but are spread when you look at the 3D plot in different angles. This was a challenging task. In the end it was decided on five different groups.

In addition 2D score plot for PC1 vs PC2, PC1 vs PC3 and PC2 vs PC3 were studied to confirm that the groups were present in all 3 PCs (see figure 4.21, figure 4.22 and figure 4.23)

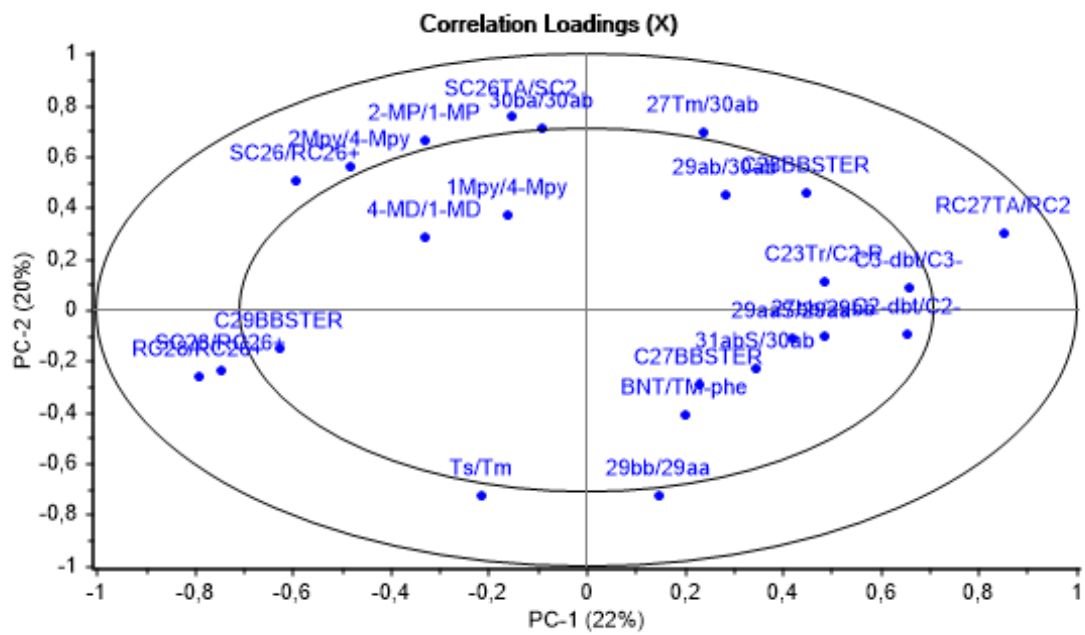


Figure 4.16: Correlation loadings plot after removing five variables with little explained variance.

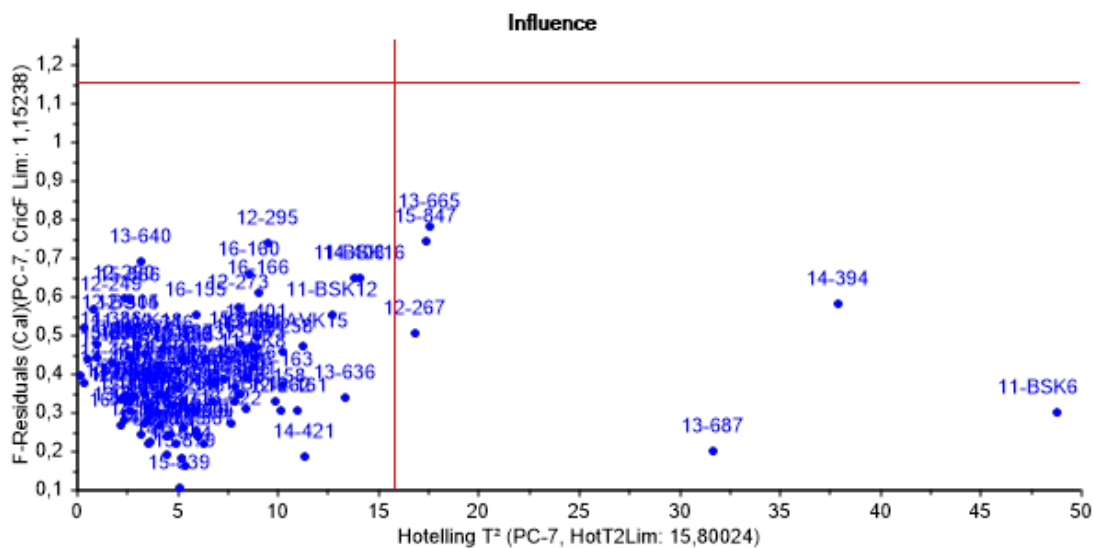


Figure 4.17: Hotelling T^2 test after removing five outliers. There are six samples outside the confidence interval with increased leverage values but, non were removed from the analysis.

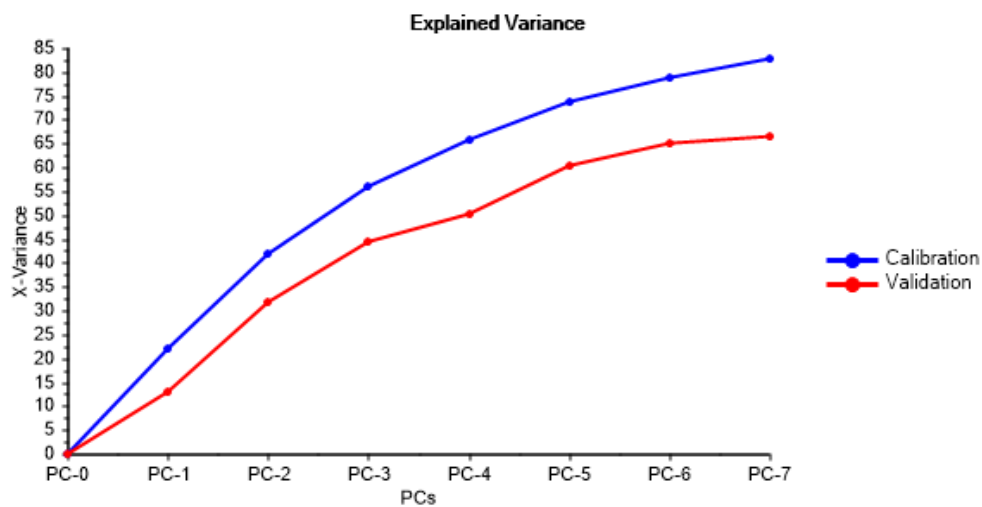


Figure 4.18: Explained variance plot displaying both calibration variance (blue) and validation variance (red).

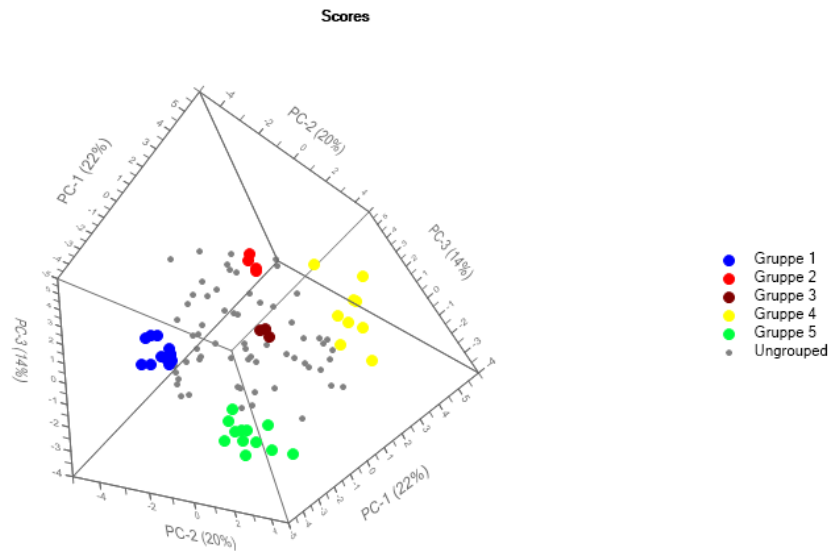


Figure 4.19: Score plot in 3D space with five groups

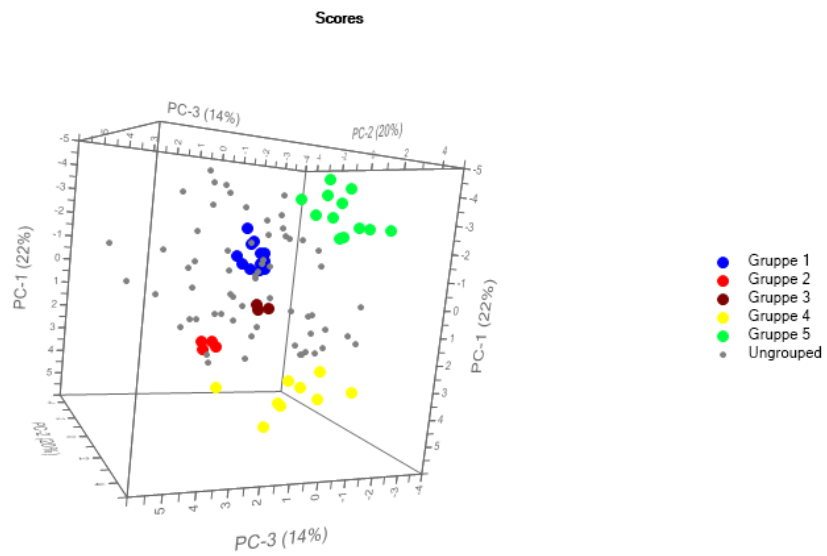


Figure 4.20: Score plot in 3D space with five groups observed from a different angle.

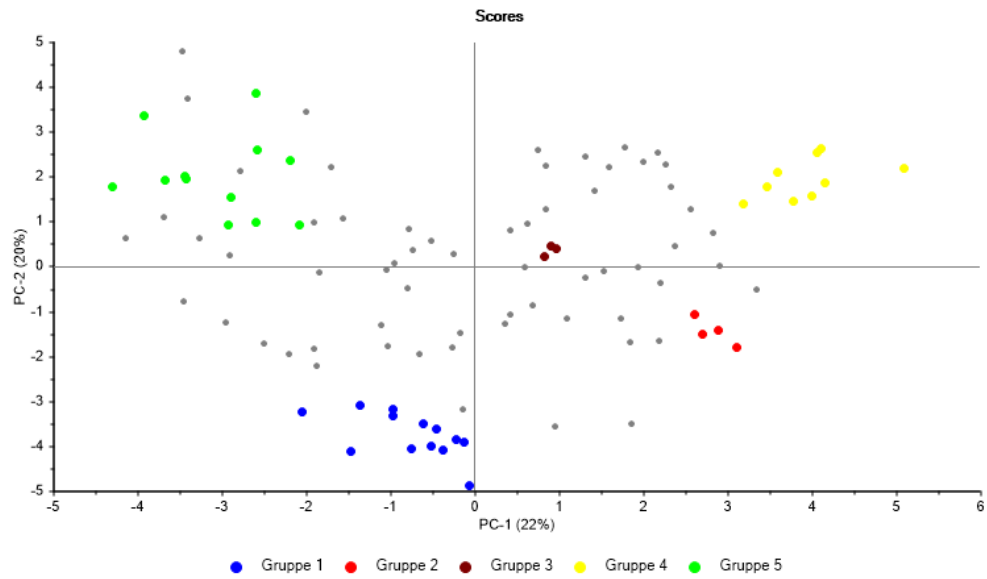


Figure 4.21: Score plot, PC2 vs PC1, with groups made in 3D score space.

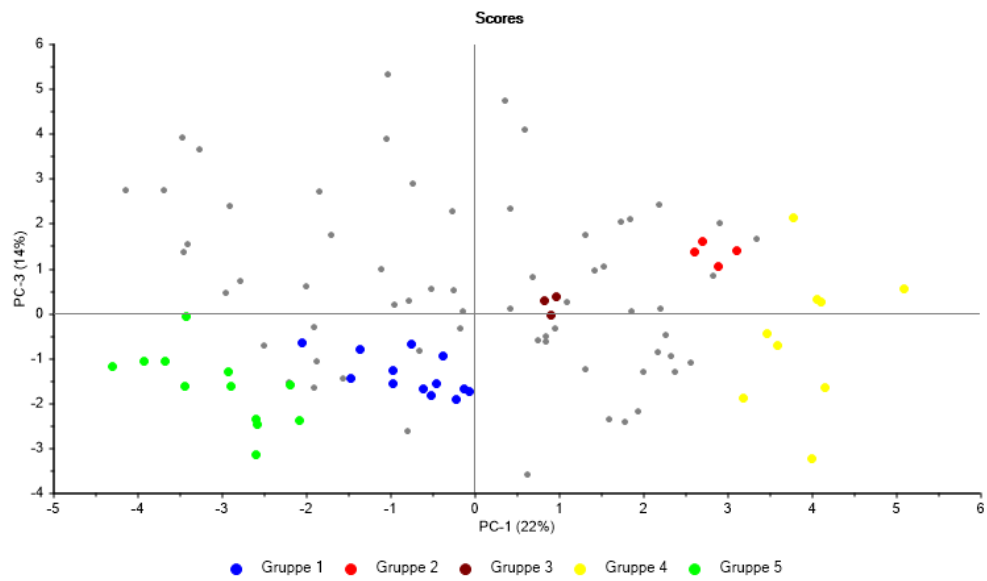


Figure 4.22: Score plot, PC1 vs PC3, with groups made in 3D score space.

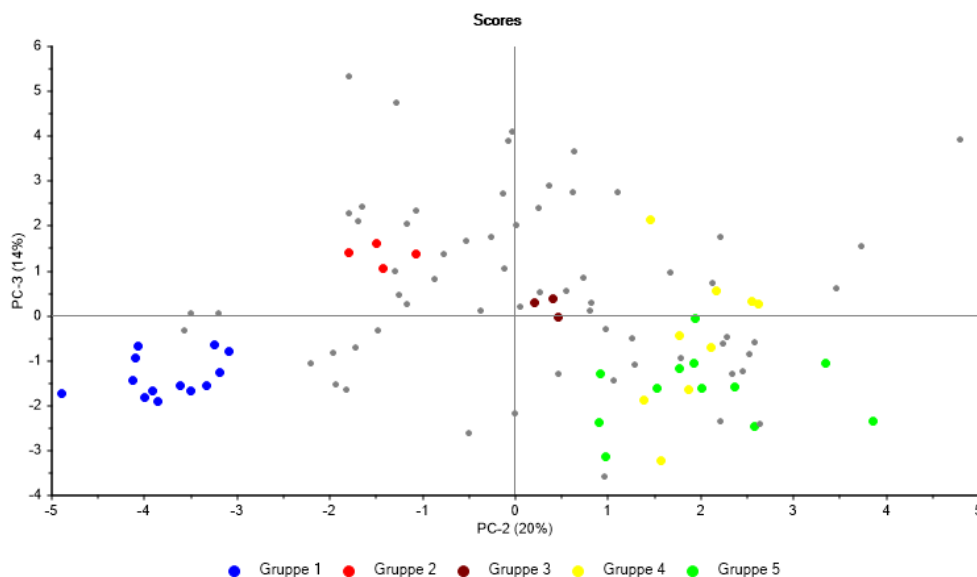


Figure 4.23: Score plot, PC3 vs PC2, with groups made in 3D score space.

When comparing the three score plots, the groups are more distinct and more separated for the score plots displaying PC1 versus PC2 and PC1 versus PC3. In the score plot for PC2 versus PC3, group 1, 4 and 5 merge into each other and are less obvious.

To determine if the different groups actually contain samples that are similar to each other, the relative standard deviation (%) was calculated for variables within each of the five groups. They were then compared with the RSD containing all samples, see table 4.6.

Figure 4.6: For groups with similar samples, one would expect the RSD to decrease. The table shows that the RSD especially decreases for group 1, group 2, group 3 and group 4, but much less for group 5. The RSD decrease for group 5 as well, but not to the same degree.

Color coding of PCA according to classes made by inspection of chromatograms

Samples were color coded according to the classification made by inspection of chromatograms, to look for interesting trends or patterns in the score plot, see figure 4.24. The blue circles are samples classified as bunker oils, the dark green circles are samples classified as crude oils, the yellow circles samples are classified as non north

Table 4.6: Relative standard deviation (%) for variables within the five groups identified in 3D score space and for all samples in the dataset.

	All samples	Group 1	Group 2	Group 3	Group 4	Group 5
Variable	N=109	n= 13	n=4	n=3	n=9	n=12
4-MD/1-MD	43%	21%	17%	21%	19%	43%
2-MP/1-MP	44%	15%	26%	24%	37%	33%
2Mpy/4-Mpy	47%	12%	3%	1%	18%	26%
1Mpy/4-Mpy	34%	24%	7%	30%	37%	48%
BNT/TM-phen	56%	27%	19%	66%	47%	36%
27Ts/30ab	30%	9%	12%	1%	29%	19%
27Tm/30ab	39%	5%	6%	21%	11%	36%
29ab/30ab	32%	5%	4%	2%	12%	12%
31abS/30ab	15%	5%	7%	3%	23%	7%
27bb/29bb	21%	6%	5%	2%	15%	11%
SC26/RC26+SC27	39%	12%	13%	8%	13%	27%
SC28/RC26+SC27	22%	8%	7%	7%	10%	17%
RC27/RC26+SC27	15%	6%	9%	4%	4%	9%
RC28/RC26+SC27	32%	9%	8%	11%	7%	15%
C2-dbt/C2-phe	67%	12%	9%	16%	44%	26%
C3-dbt/C3-phe	65%	13%	5%	3%	44%	24%
C23Tr/C2-PA	240%	79%	33%	24%	110%	104%
29aaS/29aaR	15%	7%	4%	3%	11%	24%
C20TA/C21TA	24%	11%	9%	9%	15%	23%
C21TA/RC26+SC27	66%	14%	21%	20%	21%	72%
Ts/Tm	47%	13%	9%	18%	21%	54%
30ba/30ab	29%	7%	2%	3%	30%	19%
C21TA/RC28TA	66%	18%	13%	32%	25%	81%
SC26TA/SC28TA	35%	8%	13%	2%	23%	27%
RC27TA/RC28TA	20%	4%	2%	9%	7%	10%
C27BBSTER	19%	6%	6%	9%	15%	13%
C28BBSTER	17%	6%	8%	18%	10%	12%
C29BBSTER	21%	6%	5%	9%	9%	8%
29bb/29aa	18%	12%	4%	4%	26%	19%

Table 4.7: Samples selected for use in the two training sets.

Bunker	14-366, 14-394, 15-857, 15-846, 13-665, 12-256, 12-262, 12-276
Crude	14-408, 14-420, 14-421, 14-425, 15-839, 15-844, 15-848, 13-673, 12-AS02, 11-BK4, 11-BK7, 11-BSK16, 11-AK20, 11-AK1, 11-AK9, 11-AK17, 11-AVK18, 12-251, 12-273, 12-273, 12-276, 12-314
Crude oil	14-384, 14-387, 14-420, 15-839, 15-844, 13-673, 11-BK4, 11-AK1, 12-276
NNCrude oil	14-408, 15-882, 13-680ii, 13-681i, 11-BK1, 12-AS02, 12-249, 12-314

group of crude oils. 33 % of the X-variance is used to explain 80 % of the Y variance. Factor 2 is 12 % of the X-variance to explain 8 % of the Y-variance.

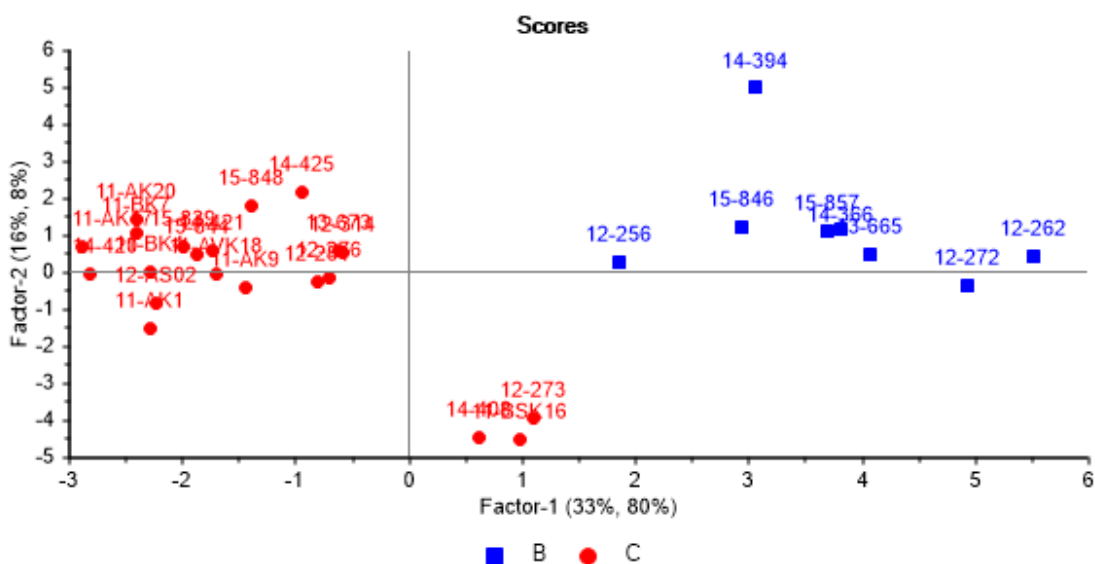


Figure 4.26: PLS-DA Score plot. Samples with red markings are classified as crude oils, samples with blue markings are classified as bunker oils.

Root Mean Square of Calibration (RMSEC) is 0.30 and Root Mean Square of validation (RMSEP) is 0.43. The R^2 value for the calibration model is 0.88 and 0.78 for the validation model. The optimal number of factors are 2.

The regression coefficient plot in fig 4.27 indicate important variables for describing each class. Variables 4-MD/1-MD, 2-MP/1-MP, 2-MPy/4-Mpy, 1-Mpy/4-Mpy are most important for describing class B, and TS/Tm, C23Tr/C2-PA, 27TS/30ab, BNT/TM-phen are most important for describing class C.

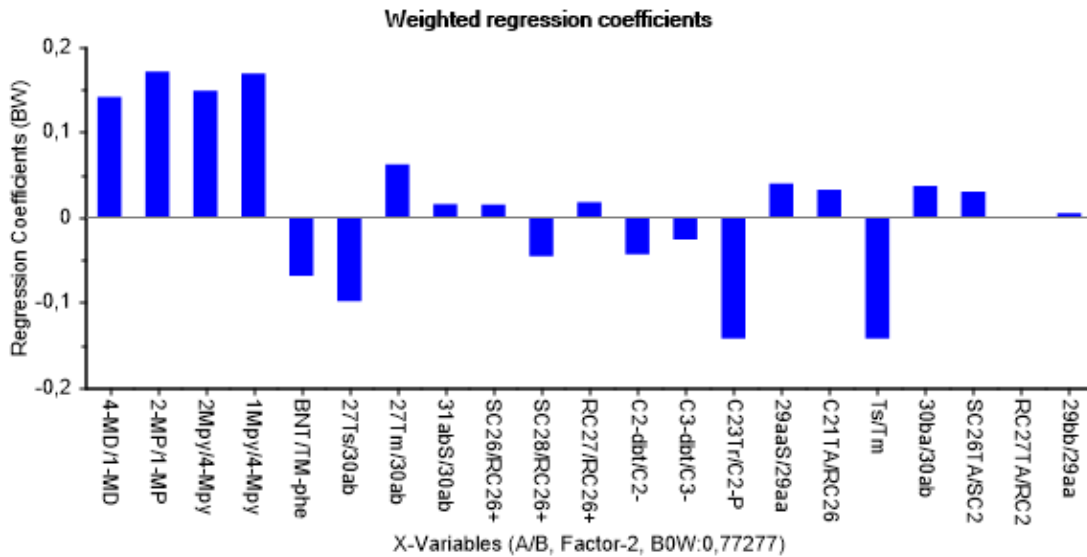


Figure 4.27: Regression coefficient plot, showing important variables for class B and class C.

PLS-DA model are then used to predict samples in the dataset. Prediction was done for all samples in the dataset except for samples that was included in the training set. Predicted samples are shown with their deviation in fig4.28. Samples with values >0.5 and with deviations that do not exceed the 0.5 limit belong the group Bunker oils, however no samples were predicted to be bunker oil from inspecting fig4.28. Samples with values <0.5 and with deviations that do not exceed the -0.5 limit belong to the group crude oil and are colored in green. 22 samples were predicted to be crude oils. These samples are listed in figure 4.8.

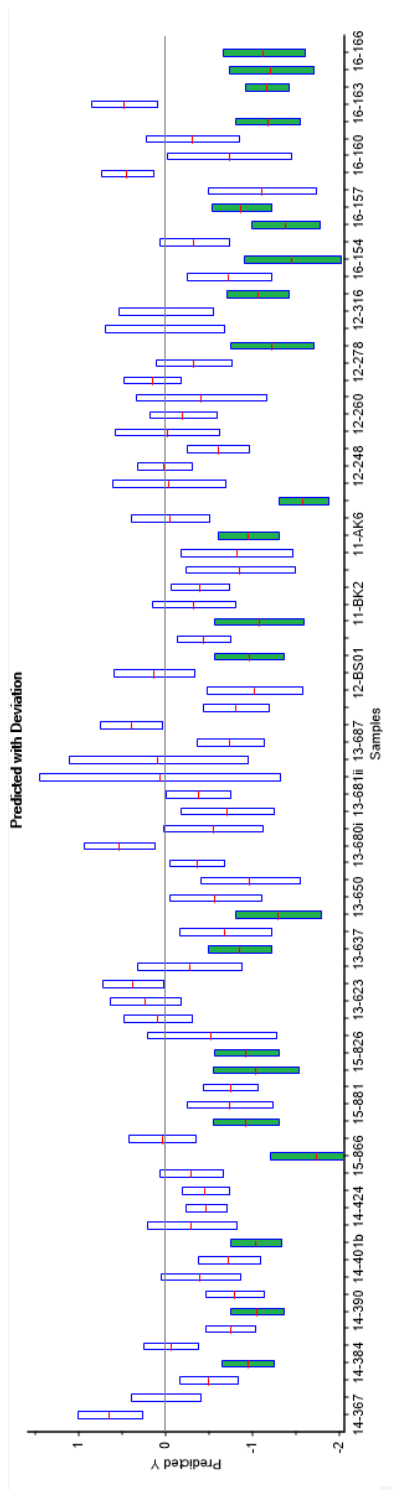


Figure 4.28: Predicted samples with deviations. Samples colored in green are classified in class C, which means that they are characterized as crude oils.

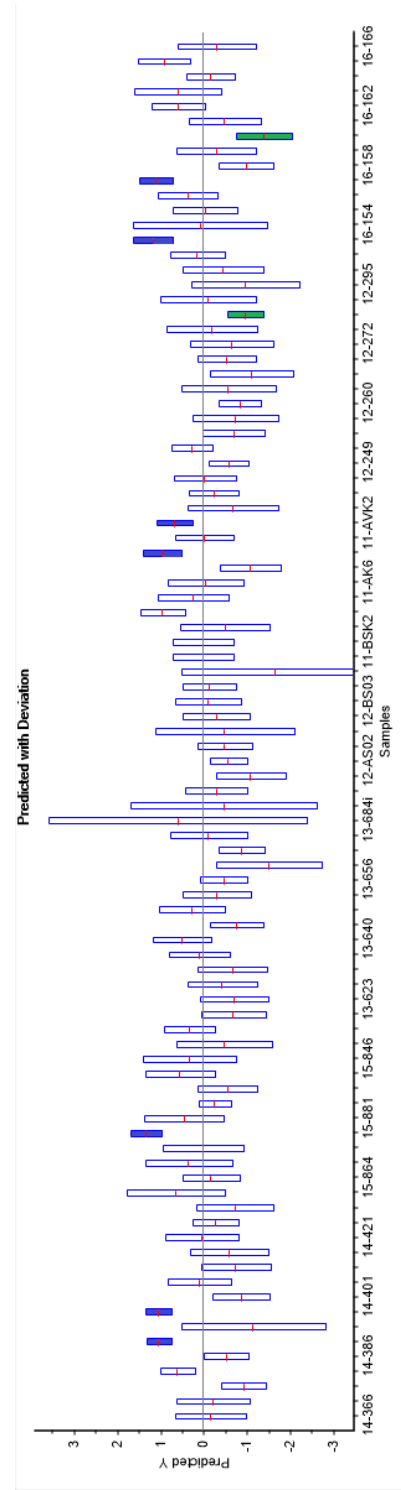


Figure 4.29: Samples predicted to class Crude oil and samples predicted to class non north sea crude oil

Table 4.8: Samples predicted to class C, crude oil. Non samples were predicted to class B, bunker oil.

Crude oil	14-384, 14-390, 14-413, 15-866, 15-879, 15-826, 15-837, 13-637, 13-643, 12-BS03, 11-BK1, 11-AK6, 11-AVK2, 12-290, 16-152, 16-154, 16-156, 16-157, 16-162, 16-164, 16-165, 16-166
Bunker oil	None

A Inlier vs Hotelling T^2 test shows that all samples fall inside the limit line, and indicates that all predictions can be trusted.

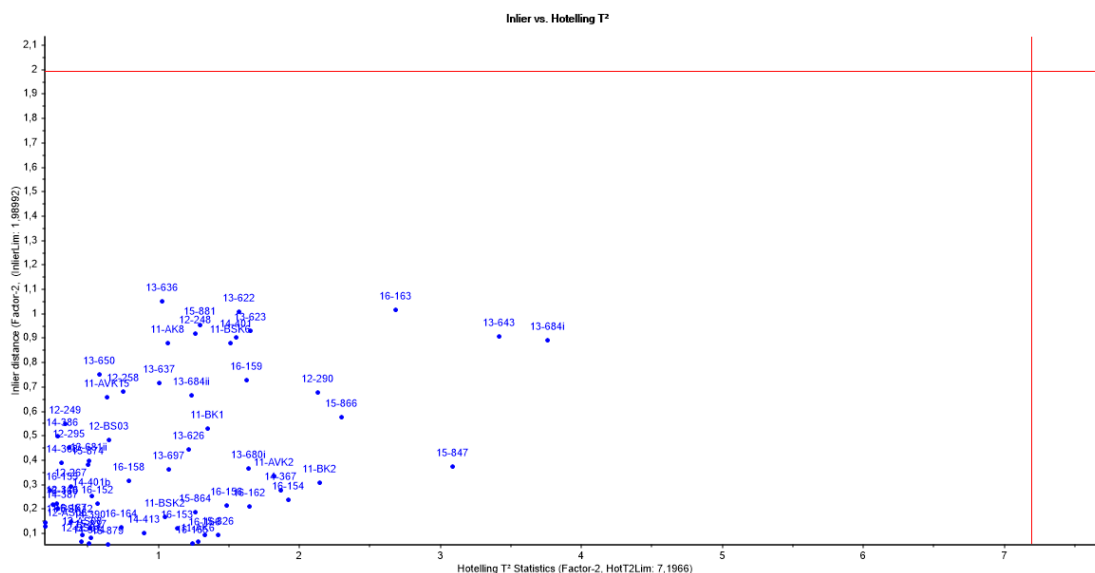


Figure 4.30: Inlier vs Hotelling T^2 shows that samples fall inside the limit, which means that all predictions can be trusted.

4.3.3 PLS-DA: Crude oils

PLS-DA model was created to differentiate between the origin of crude oils, that is, differentiate between "crude oil" (possible from the north sea) and "non north sea crude oil".

The training set involves 17 crude oils, whereby nine samples are classified as crude oils and 8 samples are classified as "non north sea crude oils". The Y-variable is a binary variable with value 1 for member of class C (crude oil) and value -1 for member of class C_NN (non north sea crude oil").

Samples used for classification is presented in table 4.7. The resulting PLS-DA model is presented in figure 4.31. From the scores plot it is mainly factor 1 that discriminates between the two groups.

It takes 53% of the X-variance to explain 83% of the Y variance. For factor 2 it takes 10% of the X-variance to explain 9% of the Y-variance. The RMSEC (root mean square error calibration) and the RMSEP (root mean square error prediction) are 0.41 and 0.49 respectively. The R^2 are 0.83 for the calibration model and 0.79 for the validation model. The optimal number of factor is 1.

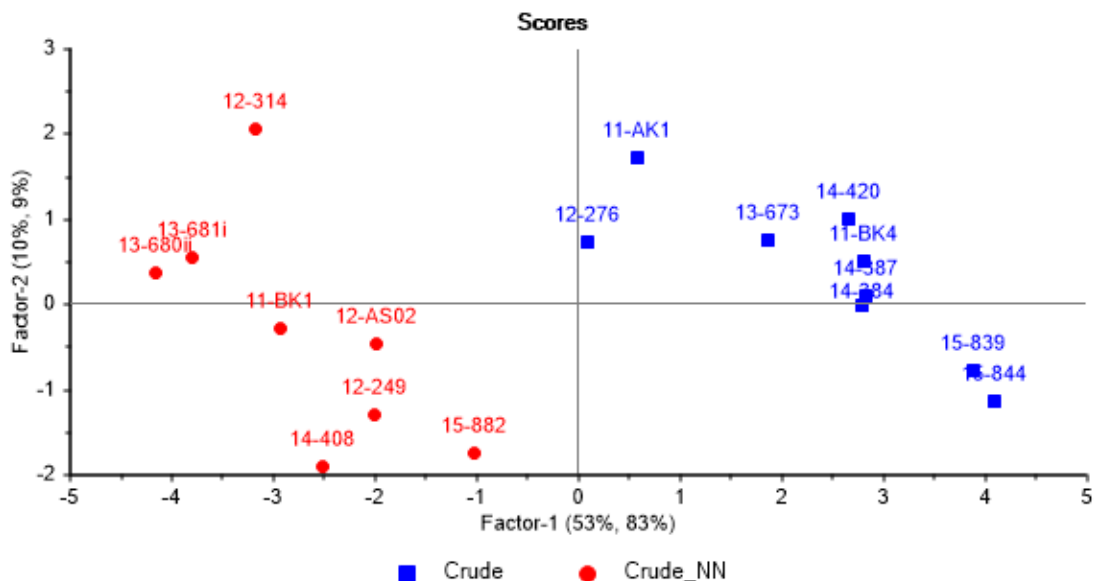


Figure 4.31: PLS-DA Score plot. Samples with red markings are classified as crude oils, samples with blue markings are classified as non north sea crude oils

The regression coefficients plot in fig 4.27 shows important variables for describing each class. Variables BNT/TM-phen, Ts/Tm, 29bb/29aa and RC28/(26+SC27) are important variables for the Y-variable Crude. Variables 27Tm/30ab, 29ab/30ab, C2-dbt/C2-phe, C3-dbt/C3-phe, C23Tr/C2-PA and 30ba/30ab are important variables for the Y-variable Crude not north sea.

The prediction was done for all samples in the dataset except samples in the

Table 4.9: Samples predicted to class Crude oil and samples predicted to class NN_Crude oil

Crude oil	14-390, 14-395, 15-879, 11-AK9, 16-153, 16-157
NN_Crude oil	12-278, 16-160

training set model. The prediction and can be seen in figure 4.32. Samples with values >0.5 and with deviations that do not exceed the 0.5 limit, belong the group "crude oils", and these are colored in blue. Samples with values <0.5 and with deviations that do not exceed the -0.5 limit belong to the group "non north sea crude oil" and are colored in green.

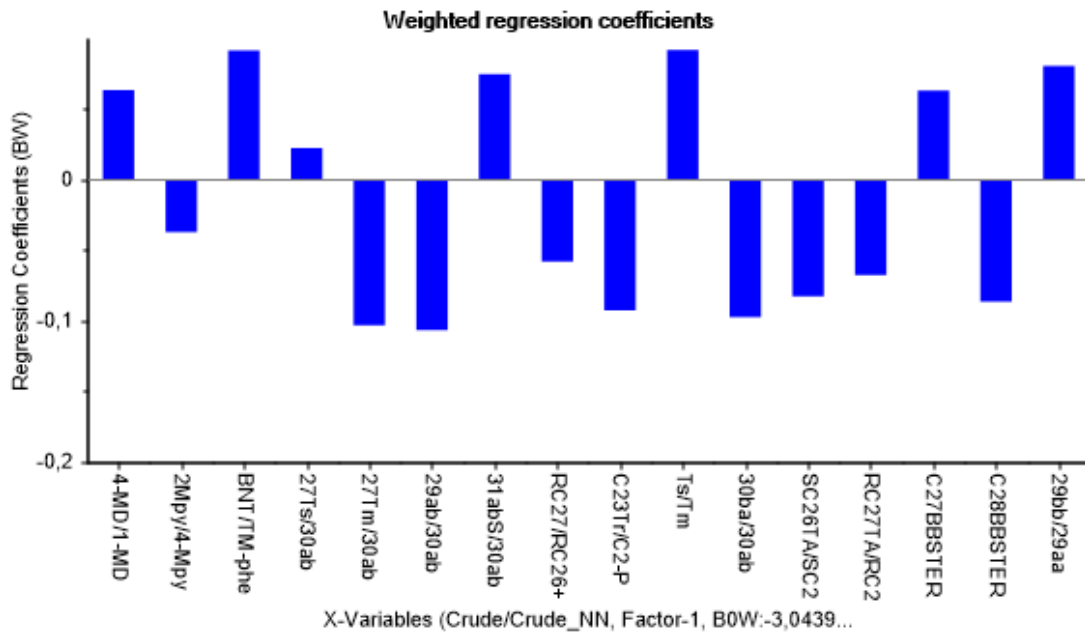


Figure 4.32: Regression coefficient plot, showing important variables for class crude oil and class non north sea crude oil.

Samples belonging to each of these groups are presented in table 4.9

Inlier versus hotelling test indicate that there is one sample where this prediction cannot be trusted (2011-BK2). This sample was suspected to be an outlier in PCA and this strengthens the suspicion further.

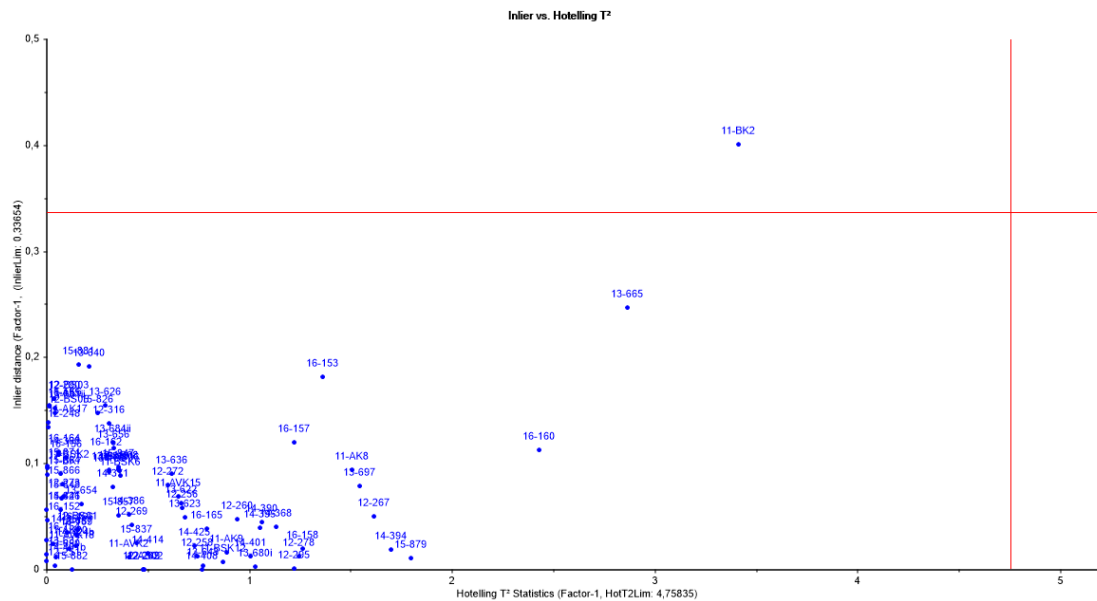


Figure 4.33: Inlier vs Hotelling T^2 shows sample 2011-BK2 fall outside the limit, which means that this prediction cannot be trusted.

4.3.4 Hierarchical Cluster Analysis

HCA was computed with Euclidian distance and was tested with four different linkage methods; Single linkage, Complete linkage, Average linkage and Ward’s method. These results were displayed in a dendrogram. Each dendrogram was validated by calculating the copenetic correlation coefficient (CPCC). The dendrogram displaying the highest correlation was further assessed. Average linkage returned the highest CPCC (0.79), and Ward’s method returned the lowest CPCC (0.50).The dendrogram with average linkage are displayed in figure 4.34.

Optimal number of clusters

Due to missing values in the dataset, the traditional "elbow method" for evaluating the optimal number of clusters were not possible to compute. Hence,the number of clusters were evaluated by manually computing different number of clusters in the dendrogram. Based on these results it was decided that 16 clusters gave the most reasonable separation between samples.This resulted in two large clusters and many small clusters. The two largest clusters hold 37 and 24 samples in each cluster. Table 4.10 provides a simplified overview over samples belonging to each cluster.

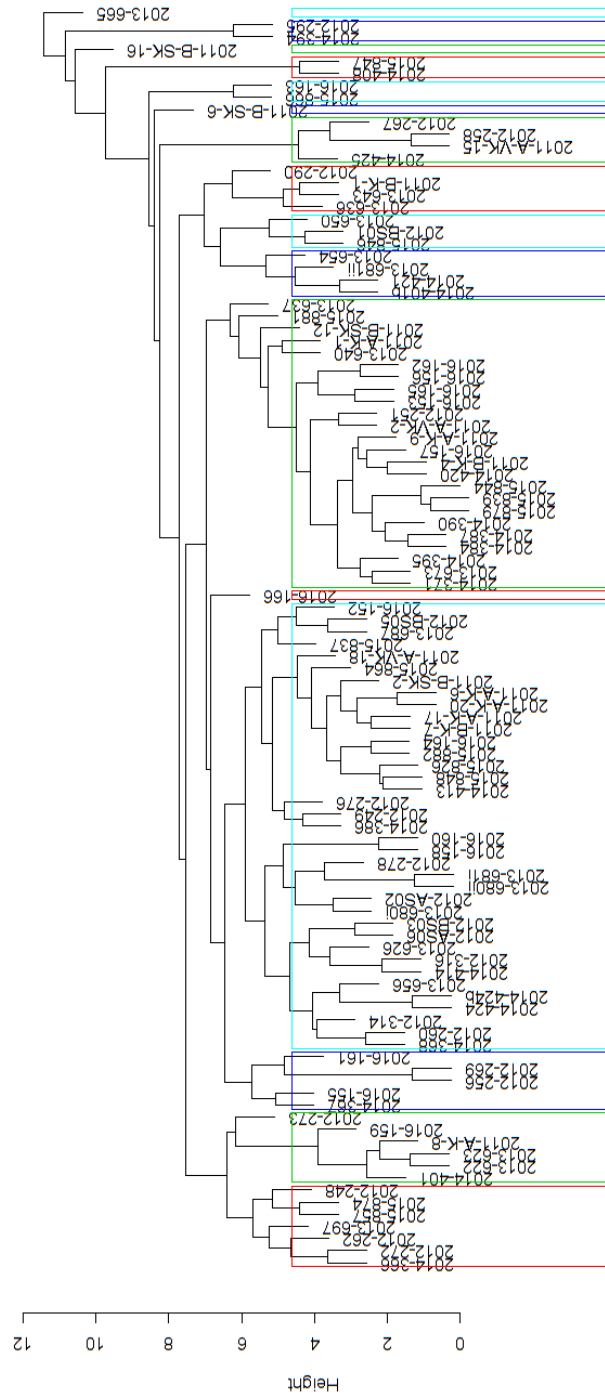


Figure 4.34: Dendrogram of hierarchical clustering analysis with 16 clusters

Table 4.10: Samples belonging to each of the 16 clusters

Cluster 1	2013-665
Cluster 2	2012-295, 2014-394
Cluster 3	2015-847, 2014-408
Cluster 4	2011-B-SK-16
Cluster 5	2016-163, 2015-866
Cluster 6	2011-B-SK-6
Cluster 7	2012-267, 2012-258, 2011-A-VK-15, 2014-425
Cluster 8	2012-290, 2011-B-K-1, 2013-643, 2013-636
Cluster 9	2013-650, 2012-BS01, 2015-846
Cluster 10	2013-654, 2013-681ii, 2014-421, 2014-401b
Cluster 11	2013-637, 2015-881, 2011-B-SK-12, 2011-A-K-1, 2013-640, 2016-162, 2016-156, 2016-165, 2016-153, 2012-251, 2011-A-VK-2, 2011-A-K-9, 2016-157, 2011-B-K-4, 2014-420, 2015-844, 2015-839, 2015-879, 2014-390, 2014-387, 2014-384, 2014-395, 2013-673, 2014-371
Cluster 12	2016-166
Cluster 13	2016-152, 2012-BS05, 2013-687, 2015-837, 2011-A-VK-18, 2015-864, 2011-B-SK-2, 2011-A-K-6, 2011-A-K-20, 2011-A-K-17, 2011-B-K-7, 2016-164, 2015-882, 2015-826, 2015-848, 2014-413, 2012-276, 2012-249, 2014-386, 2016-160, 2016-158, 2012-278, 2013-681i, 2013-680ii, 2012-AS02, 2013-680i, 2012-BS03, 2012-AS06, 2013-626, 2012-316, 2014-414, 2013-656, 2014-424b, 2014-424, 2012-314, 2012-260, 2014-368
Cluster 14	2016-161, 2012-269, 2012-256, 2016-155, 2014-367
Cluster 15	2012-273, 2016-159, 2011-A-K-8, 2013-623, 2013-622, 2014-401
Cluster 16	2012-248, 2015-874, 2015-857, 2013-697, 2012-262, 2012-272, 2014-366

Table 4.11: Groups identified by 3D plotting are color coded after classification made by visual inspection. Samples identified as dark green belong to the group "crude oil", samples identified as yellow belong to the group "non north sea crude oil", samples identified as blue belong to the group "bunker oil" and samples identified as red belong to the group "unkown".

PCA groups	Color coded samples in groups
1	2014-384, 2014-387, 2014-390, 2014-395, 2014-420, 2015-879, 2015-839, 2015-844, 2011-BK4, 2011-AK9, 2016-153, 2016-157, 2016-165
2	2014-413, 2015-826, 2015-848, 2011-BK7
3	2016-155, 2014-424, 2014-424b
4	2014-425, 2013-680ii, 2013-681i, 2011-AVK15, 2012-258, 2012-267, 2016-158, 2016-160, 2016-161
5	2014-401, 2014-408, 2015-857, 2015-847, 2013-622, 2013-623, 2011-AK8, 2012-248, 2012-262, 2012-272, 2012-273, 2016-159

4.3.5 Color coding of groups in HCA and PCA based on results inspection of chromatograms

Groups created in PCA and the eight largest clusters in HCA were color coded according to the classification of samples from inspection of chromatograms. The results are presented in table 4.11 and 4.12. Samples identified as dark green belong to the group "crude oil", samples identified as yellow belong to the group "Non north sea crude oil", samples identified as blue belong to the group "bunker oil" and samples identified as red belong to the group "unkown".

4.3.6 COSIWeb

In total, 63 samples were imported and compared in the COSIWeb database. Only samples with at least one correlation coefficient (to an existing sample) above 0.98 were further evaluated to conclude or deny a match, and only the five highest correlations for each sample was inspected (given that the correlation coefficient was above 0.98). 16 samples returned correlation coefficient above 0.98 and these are listed in table 4.33.

The rows in 4.33 represents the 16 samples, and the column represents the five

Table 4.12: Groups identified by 3D plotting are color coded after classification made by visual inspection. Samples identified as dark green belong to the group "crude oil", samples identified as yellow belong to the group "non north sea crude oil", samples identified as blue belong to the group "bunker oil" and samples identified as red belong to the group "unknown".

HCA groups	Color coded samples in groups
1	2011-A-VK-15, 2012-258, 2012-267, 2014-425
2	2012-290, 2011-B-K-1, 2013-643, 2013-636
3	2013-654, 2013-681ii, 2014-421, 2014-401b
4	2013-637, 2015-881, 2011-B-SK-12, 2011-A-K-1, 2013-640, 2016-162, 2016-156, 2016-165, 2016-153, 2012-251, 2011-A-VK-2, 2011-A-K-9, 2016-157, 2011-B-K-4, 2014-420, 2015-844, 2015-839, 2015-879, 2014-390, 2014-387, 2014-384, 2014-395, 2013-673, 2014-371
5	2016-152, 2012-BS05, 2013-687, 2015-837, 2011-A-VK-18, 2015-864, 2011-B-SK-2, 2011-A-K-6, 2011-A-K-20, 2011-A-K-17, 2011-B-K-7, 2016-164, 2015-882, 2015-826, 2015-848, 2014-413, 2012-276, 2012-249, 2014-386, 2016-160, 2016-158, 2012-278, 2013-681i, 2013-680ii, 2012-AS02, 2013-680i, 2012-BS03, 2012-AS06, 2013-626, 2012-316, 2014-414, 2013-656, 2014-424b, 2014-424, 2012-314, 2012-260, 2014-368
6	2016-161, 2012-269, 2012-256, 2016-155, 2014-367
7	2012-273, 2016-159, 2011-A-K-8, 2013-623, 2013-622, 2014-401
8	2012-248, 2015-874, 2015-857, 2013-697, 2012-262, 2012-272, 2014-366

best matches along with correlation coefficients. A code describing the country that uploaded the sample is denoted. UK is short for United Kingdom, Fi is short for Finland, NI is short for Northern Ireland, DE is short for Deutschland and BE is short for Belgium. The remaining samples are imported from SINTEF. An example is that sample 2014-384 has best match with sample 2014-384 and their correlation coefficient is 0.9988.

A positive, probable, inconclusive or non-match between two samples were given after manual inspection of the GC-FID chromatograms, GC-MS chromatograms and diagnostic ratios in COSIWeb. Samples that were concluded to be a positive match are colored in green. Samples with a probable match is colored in yellow, Samples with an inconclusive match is colored in blue and samples with non match are colored in red (see figure 4.33). In addition, number of DRs exceeding the analytical limit between two samples are shown nearby the correlation coefficients.

COSIWeb was used to identify samples in this dataset with samples from other projects and laboratories. The results show that 6 samples from the dataset were considered a probable match to four samples from the Shetland Islands. These are chromatograms which have been imported from a lab in the United Kingdom, and are collections of oils from bird feathers from the Shetland Islands. The type of oil in these samples are unknown from their description.

Sample UK-1-231.4 was imported into COSIWeb in 2012 and has best match to sample 2014-384 (0.9946). The GC-MS chromatograms look very similar although 13 of the 26 peaks have DR exceeding the 7 % limit which corresponds to 14 % in the CEN. Two of these ratios (30G and MA) display 100 % difference and should probably not have been integrated in the first place since they display low signals in the chromatogram. Of the remaining 11 ratios, there are mostly PAHs and acyclic isoprenoids that are above the limit, and most of these differences are most likely due to weathering (such as C17/pr, C18/ph and pr/ph). Sample UK-1-231.4 has next best match to sample 2014-387 (0.9935) and are considered as a probable match. The GC-MS chromatograms look very similar. 15 of the 26 peaks exceed the limit, mostly PAHs. As with sample 2014-384 there are two ratios (30G and MA) that display 100 % difference and should not have been included. It has third best match to sample 2011-BK-4 (0.9923). There are 17 DR outside the limit, mostly PAHs and acyclic isoprenoids.

Sample UK-1-341.3 were imported into COSIWeb in 2014 and has best match to sample 2014-384 (0.9943). The GC-MS chromatogram display many similarities. 19 of 26 peaks have DR exceeding the limit. Two of these ratios display (30 O and 30G) 100 % difference and should probably not have been integrated in the first place since they display low signals. The remaining 17 ratios above the limit are a mixture of

both hopanes, PAHs, steranes, and acyclic isoprenoids and most of these differences are most likely due to weathering. Despite the many differences in DR they were still considered to be a probable match. Next best match is to sample 2014-387 (0.9929). 22 diagnostic ratios are above the limit and as with sample 2014-384, it is a mixture of different hopanes, PAHs, steranes and acyclic isoprenoids. Other matches are to samples 2011-BK-4 (0.9919), 2014-420 (0.9911). 2011-A-VK-2 (0.9869) and 2016-165 (0.9826). These were also considered a probable match, although between 20-23 DR for these samples are above the limit.

Sample UK-1-332.2 probable match to 2015-846 (0.9856). 11 Diagnostic ratios are above the limit. Again, a majority of these are PAHs and Isoprenoids and most differences are due to weathering.

Sample UK-1-231-5 was imported into COSIWeb in 2012 and has best match to sample 2014-387 (0.9889) and is considered a probable match. 13 DR are outside the limit, mostly PAHs and Isoprenoids.

Table 4.13: Samples with coefficient correlation above 0.98. Samples evaluated as positive match are colored in green, samples displaying a probable match are colored in yellow, samples with no match are colored in red and inconclusive matches are colored in blue. Matches with external samples are shown in bold.

Sample	Match 1	Match 2	Match 3	Match 4	Match 5
2014-384	2014-0387 0.9988/5	UK-1-231-4 0.9946/13	2011-B-K-4 0.9943/8	UK-1-341.3 0.9943/19	2014-420 0.9932/11
2014-0387	2014-0384 0.9988/5	UK-1-231.4 0.9935/15	UK-1-341.3 0.9929/22	2011-B-K-4 0.9894/9	UK-1-231.5 0.9889/13
2014-0420	2011-B-K-4 0.9981/8	2016-165 0.9945/7	2014-0384 0.9932/11	2011-A-VK-2 0.9919/11	Uk-1-341.3 0.9911/20
2014-394	NI-1-3063.1 0.9856/16	Fi-1-13.8 0.9854/20	Ni-1-3313-4 0.9842/22	FI-1-13.17 0.9835/22	De-1-1186.10 0.9833/21
2015-879	2015-839 0.9992/3	2015-844 0.9877/4			
2015-839	2015-879 0.9992/3	2015-844 0.9883/4	De-1-1042.1 0.9833/20		
2015-846	2007-0064* 0.9856/13	Uk-1-332.2 0.9856/11	BE-1-6042.142 0.9856/14		
2015-881	BE-1-6042.45 0.9845/14	BE-1-6042.25 0.9799/13			

2015-844	2015-839 0.9883/4	2015-879 0.9877/4	De-1-042.1 0.9818/19		
2016-165	2014-420 0.9945/7	2011-B-K-4 0.9938/13	2011-A-VK-2 0.9888/13	UK-1-341-3 0.9826/23	2014-0384 0.9815/13
2012-BS01	2012-OBS05 0.9824/19				
2012-BS05	2012-OBS01 0.9824/19	2012-0115 0.982/17			
2011-B-K-4	2014-420 0.9981/8	2014-384 0.9943/8	2016-165 0.9938/13	Uk-1-231.4 0.9923/17	Uk-1-341.3 0.9919/23
2011-A-VK-2	2014-420 0.9919/11	2011-B-K-4 0.9899/11	2016-165 0.988/13	uk-1-341.3 0.9869/23	
2012-0262	DE-1-178230 0.9817/19				

Chapter 5

Discussion

5.1 Sources of errors

Samples in this study have been collected and analyzed at different years and by different students. They have all been analyzed according to the CEN procedure [CEN, 2012] and this minimizes errors arising from different methods of analysis.

Sample material in this project are of different size, thickness and degree of weathering. Even though efforts have been made to select the center of the sample that is least affected by weathering, this can be difficult for small and thin oil samples. This has led to some samples being more weathered than others and may have affected the GC-FID and GC-MS results. In addition contamination may have affected some of the samples. Contamination can for example occur from multiple oil spills or background chemicals. These are factors that can affect identification and make it more difficult to characterize samples.

Although the author has not taken part in fieldtrips and laboratory work for a majority of the samples in this project, manual integration of samples have solely been performed by the author. This was done to measure the peak height for biomarkers, and peak area for PAHs and was assessed for approximately 8000 components. Manual integration is a subjective evaluation, and depends on the experience of the analyzer, so it is an advantage that this has been done by one person. This is because the person in charge of the integration must decide on the peak boundaries, how to handle co-eluting peaks (overlapping), fronting and tailing. The integration will to a large degree depend on the experience of the analyzer. The author had little knowledge of manual integration prior to this thesis, but gained knowledge of this through training under guidance of an experienced employee at SINTEF. Nevertheless, samples analyzed in the beginning may have been analysed and handled differently from

samples towards the end of the integration, although the author attempted to make it as equal as possible.

Samples with biomarkers and PAH components below the detection limit were not included in the analysis, which lead to some diagnostic ratios not being included in the thesis. This could be unfortunate for some of the samples that actually have valid measurements for some of the DRs. An alternative approach would be to fill in missing values for samples below the detection limit, but this was not considered for this work.

5.2 Inspection of chromatograms

Results from the inspection of chromatograms, based on recommendations from the CEN methodology [CEN, 2012] have been used as a starting point for other analyses. This can be a vulnerable technique to classify oil samples. Firstly, the visual inspection will be affected by how weathered the sample material is. This may lead to faulty assumptions or no conclusion at all. Faulty assumptions can be prominent in cases where abundance of a peak or abundance between peaks are used for characterization. Separation between crude oil and bunker oil was based on this, since oils are characterized as crude oils if the first doublet (3-methyl-and-2-methylphenanthrenes) is more abundant than the second doublet (9-/4 and 1-methylphenanthrenes) and vice versa for bunker oils (see figure 3.4 and figure 3.5). However for bunker oils, the characterization is safer in the sense that there also have to be a MA peak present for it to be classified as a bunker oil. Secondly this method also depend largely on the analyst, and the experience of the analyst.

Many of the GC-FID chromatograms were affected by weathering. The best information that could be acquired from these, was basically to separate the most weathered samples from less weathered samples. The results from the GC-MS chromatograms were easier to interpret, which is as expected since this evaluation was based on biomarkers and PAH components that are resistant to degradation. Approximately 36% of the samples could not be identified solely on the inspection of chromatograms. In some samples, even the most resistant biomarkers were highly weathered, which made it hard to characterize them. This was especially evident in ion chromatogram m/z 192, which was used to compare doublets to decide if the oil was crude oil or bunker oil. In some samples the abundance of the doublets were equal which made it difficult to classify them based on the height of the doublets. In some cases a very small MA peak could be identified but due to much noise, in the chromatogram it was difficult to conclude if it was the noise or the biomarker. In these cases they were identified as unknown. The abundance of oleanane og gam-

macerane (see figure 3.6 and 3.7) is to some degree a subjective evaluation, by which the person decides to what the degree of abundance must be to classify them into certain groups.

The results from the visual inspection were compared to students classification in the course KJ3050 - Marine Organic Environmental Chemistry, and previous master students. Of 97 samples, 27 samples were classified differently by the author compared to the students. Of 27 samples, 25 of these were classified as either crude or bunker by other students, while the author classified these as unknown. One sample was classified as bunker by the author, and this was classified as crude oil by the students. This concerns sample 2014-394. One sample was classified as crude oil by the author, and this was classified as bunker oil by the students. This is sample 2011-A-VK-15. Moreover the author characterized less crude samples into the group "Non north sea crude oil", compared to the students. Since this classification was intended to be used in the rest of the work, it was important to be prudent and to avoid any false identifications. To classify something wrong, would be worse for the analysis than having one more unknown.

For islands with many sampling points such as Sula, there were a mixture of both crude oils, bunker oils and unknown samples. For islands with small sample sizes there are some islands with samples from only one classification, in addition to unknown samples. Examples are Olabussøya where there are five sampling points, three samples are characterized as bunker oils and two samples are characterized as unknown. At Gårdsøya and Borholdmen, two very close islands, almost all samples are characterized as bunker oils, except one which is characterized as unknown. However, no conclusions about identification of unknown samples can be based solely by looking at the map since there is no guarantee, with this limited sample size, that there is only one type of oil on the island.

Interestingly, many of the biomarkers that were used to identify oil samples based on the CEN methodology, were not included in the multivariate analysis. This includes Oleanane, Gammacerane, Methylantracene and Retene (see figure 3.4, figure 3.5, figure 3.6 and figure 3.7). This is because their diagnostic ratios were missing for more than 25% of the samples. Methylantracene is only visible in bunker oils, and since only a small part of the sample material was characterized as bunker, and some probably were not characterized due to weathering, it is not surprising that this diagnostic ratio is missing many values. This also applies to the biomarkers oleanane and gammacerane. These only shows high abundance in few samples. When it comes to retene, this biomarker is only source specific for some samples.

5.3 Principal component analysis

In the analysis, three principal components were considered, which in total explained 55% of the variance. PC1 and PC2 explained 42 % of the variance, thus PC3 was also considered. Since the focus of this work was to study the larger trends in the data, it was not desirable to include too many principal components. If a more detailed view is desired, it is possible that more components should be considered, although this increase the risk of including noise.

The identification of groups in 3D score space resulted in 5 groups selected visually (figure 4.19). But there might be several groups that the author did not identify or other samples belonging to one of the groups. The RSD, which describe how similar the results within a group are, show that especially group 1-4 consist of internally similar samples. Group 5 had a higher RSD, indicating that it is a less tightly knit group (figure 4.6).

5.3.1 Combining results from inspection of chromatograms with results from PCA-2D score plot

It was not suggested groups of samples based on the 2D score plot, only in terms of the 3D plot. Instead, samples in the 2D score plot were assigned colors according to the category it belonged to from inspecting chromatograms. The score plot for PC1 and PC2 (see figure 4.24) shows that samples classified as bunker oils are distributed in the upper side of the plot, especially in the upper left corner. The majority of samples classified as crude oils are located on the lower side of the plot, although there are samples classified as crude oils in the upper right side of the plot as well, especially non north sea crude oil. The unknown samples are spread around the entire score plot. From this it seems like PC1 is related to separation of oil type.

5.3.2 Combining results from inspection of chromatograms with results from PCA-3D score plot

Table 4.11 shows that group 1-4 which was created in the 3D plot are dominated by oils that have been classified as crude oils, and most bunker oils are in group 5. The difference between north sea and non north sea crude oil is less clear. However, there is some indication of this as well, mainly because there are no “non north sea” crude oils in group 1, 2 and 3. There are two samples in group one that are categorized as unknown, namely 2014-395 and 2016-153. Since the majority of these samples have been characterized as crude oils, it can be reasonable to assume that

these unknown samples are crude oils as well. In group 4 there is a mixture of crude oils, and three unknown samples. From this one may assume that the unknown samples also are crude oil, but it is difficult to say anything about the type of crude oil. Group 5 had the highest RSD which means that there is less similarity between these samples compared to the other groups. This also reflected by looking at the color classification. The majority of samples are classified as bunker oils, but there are some samples classified as crude oil and some unknown as well.

5.3.3 Biplot

The biplot can be used to determine which variables are important for characterizing different samples. In a previous study [Sun et al., 2015] PCA biplot was used on a set of crude oils, light fuel oils, heavy fuel oils, weathered fuel oils and weathered crude oils. They used 43 diagnostic ratios, of which 23 are shared with the analysis in this thesis, and of which 13 are included in this PCA biplot.

They concluded that n-alkanes were positively correlated with light fuel oil (LFO), aromatic hydrocarbons were positively correlated with Heavy fuel oil (HFO) and terpanes and steranes were indicative of crude oil.

The biplot in figure 4.25 shows that variables located to samples classified as bunker oils are 2-Mp/1-MP, 1Mpy/4-Mpy, 2Mpy/4-Mpy 4-MD/1-MD, SC26TA/SC28TA, 30ba/30ab, SC26/RC26+SC27, 27Tm/30ab and 29ab/30ab. The first three of these were also identified as such by the previous study. Since both studies have identified these, it is an indication that these DRs are useful for identifying bunker oils. But there are also differences. From their study the DRs C2-dbt/C2-phe and C3-dbt/C3-phe are also shown to be characteristics of bunker oil. From our biplot it seems like these two diagnostic ratios are more prominent for samples classified as crude oils, but they are also located close to a group of many unknown samples.

Variables that are closest located to samples classified as crude oil (both from the north sea and not from the north sea) are C29BBSTER, SC28/RC26+SC27, RC28/RC26+SC27 and Ts/Tm, 29bb/29aa, BNT/TM-phe, C27BBSTER, 31abs/30ab, 29aaS/29aa, 27bb/29bb, C2-dbt/C2-phe and C3-dbt/3-phe, 31aBs/30ab, C23Tr/C2-P, RC27TA/RC28TA. In their study, the results shows that Ts/Tm and RC27TA/RC28TA are DRs indicative of crude oils, but they also show this for DR SC26TA/SC28TA which in our study seems more indicative of bunker oils.

Moreover the article says that the diagnostic ratios C29ab/30ab, C27BBster, C3-D/C3-C were positively correlated with LFO and distinguished LFO from both HFO and crude oils. From our biplot, the ratios 27Tm/30ab and 29ab/30ab are positively

correlated with some samples classified as bunkers, which are located on the upper right side of the plot, but nothing conclusive can be said about these being light fuel oil, since the material is extremely weathered. Moreover diagnostic ratios C27BBster is positively correlated with many crude oils in this plot. The same can be said about the diagnostic ratio C2-D/C3-C although there are more unknown samples close to this ratio.

5.4 PLS-DA

5.4.1 PLS-DA-for crude oil and bunker oil

The PLS-DA model classified 22 samples as crude oils, but no samples as bunker oils. Of 13 samples classified as crude oils, these have been classified as crude oils by inspecting chromatograms (2014-384, 2014-390, 2014-0413, 15-879, 15-826,16-152, 16-156, 16-157, 16-162, 16-165, 12-BS03, 11-BK1 and 11-AK6). The remaining 9 samples have been classified as unknown by inspecting chromatograms (15-866,15-837, 16-637, 16-643, 11-AVK2,12-290, 16-154,16-164, 16-166). Thus the model is able to predict some unknown samples into a category. Since the model predicted only known crude oil samples (from inspection of chromatograms) and no bunker oils into the crude oil category, this makes it more plausible that the unknown samples in fact are crude oil samples.

Important variables for predicting crude oils from the regression coefficient plot in figure 4.27 were TS/TM, C23Tr/C2-PA, 27TS/30ab, BNT/TM-phen. These are also variables which have been correlated to crude oils in PCA, except variable 27TS/30ab which was removed from the loading plot in PCA since it did not contribute much to the model.

The performance of the model was assessed by looking at the RMSEC, RMSEP and R^2 values for the calibration and validation model. RMSEC and RMSEP was 0.30 and 0.43 respectively. The increased prediction value is as expected since this is tested on samples that was not in the calibration model. Ideally these should be equal and as low as possible. The R^2 value for the calibration model was 0.88 and for the validation model 0.78.

In the training set only 8 samples were used in the bunker oil category, whereas 22 samples were used in the Crude oil category. The low prediction for bunker oils indicate that the samples were not representative for the training set. Important variables for predicting bunker oils were 4-MD/1-MD, 2-MP/1-MP, 2-Mpy/4-Mpy, 1-Mpy/4-Mpy. This was also seen in the PCA score plot.

5.4.2 PLS-DA-for crude oil and non north sea crude oil

The PLS-DA model predicted 6 samples as crude oils and 2 samples as non north sea crude oil. Of the 6 samples classified as crude oils, 4 of these samples (14-390, 15-879, 16-157,11-AK9) were classified as crude oils by inspection of chromatograms. The two other samples (16-153,14-395) are classified as unknown when inspecting chromatograms. Of the 2 samples classified into the group non north sea crude oil (12-278 and 16-160) these have been classified as unknown by inspection of chromatograms.

The RMSEC and RMSEP are 0.41 and 0.49. R^2 values are 0.83 and 0.79. These values are lower than the predictive ability for the PLS-DA model between crude oil and bunker oil. Variables that were important for describing samples in the group crude oil was variables BNT/TM-phen, Ts/Tm, 29bb/29aa and RC28/(26+SC27). Variables that were important for predicting the Y-variable "Crude oil not north sea" are 27Tm/30ab, 29ab/30ab, C2-dbt/C2-phe, C3-dbt/C3-phe, C23Tr/C2-PA and 30ba/30ab. These are trends that we also can identify in the PCA biplot, and all of these variables are located in or close to the upper right quadrant in PCA where the majority of samples classified as non north sea crude oil are located.

5.5 HCA

In the analysis it was determined to have 16 clusters, found by manual inspection. Fewer clusters lead to many small clusters, with one or two samples, and one large group with the remaining samples. Due to missing values it was difficult to get an objective suggested number of clusters. There are methods for accomplishing this, with missing values, but this was not prioritized as the decision for number of cluster ultimately will be a subjective evaluation in the end.

With 16 clusters, the two largest clusters were made up of 24 samples and 37 samples. In figure 4.12, eight of the largest clusters are color coded based on the category it was given by inspection of chromatograms. Group 1-4 are dominated by crude oils. Group 8 only consist of samples characterized as bunker oils. Group 5 is a large group with a mix of different types, but almost all non north sea crude oil are in this group. Since the majority of samples are classified as crude oils in group 1-4 (possibly from the north sea) this can indicate that the unknown samples in these groups also are crude oils. This applies to 2012-267 (group 1), 2012-290, 2013-643, 2013-636 (group 2), 2013-637, 2015-881, 2011-B-SK-12, 2013-640, 2016-153, 2011-A-VK-2 and 2014-395 (group 4).

Sample 2014-395 and 2016-153 are grouped together in PCA in group 1, where the majority of samples are classified as crude oils. This strengthens the assumption

that at least these two samples are crude oils. This is also seen for sample 2012-267 which is grouped in PCA in group 4, where the majority of samples are classified as crude oils (from both origins).

5.6 Search in the international database COSIWeb

High correlation coefficient between two samples in COSIWeb does not necessarily mean that the diagnostic ratios are below the limit given by CEN methodology. In fact, all comparison between samples resulted in a number of DRs exceeding the limit set by CEN (see figure 4.33). In some samples this was only evident for DRs that are prone to weathering such as acyclic isoprenoids, but in many samples this was also the case for resistant biomarkers, such as hopanes. Some differences in DRs were just slightly above the limit and some were extremely different. Even though the DRs always were considered when comparing two samples, it was inspection of the different ion chromatograms that eventually was the deciding factor. The inconsistency between high degree of correlation despite DRs exceeding the limit between many samples, may be due to weathering and contamination of samples, and may show that COSIWeb has some issues when analyzing very weathered samples.

COSIWeb found samples with correlation coefficient above 0.98 for 16 of 63 samples imported into cosiweb. Among these samples, many display high correlation coefficients with each other. This did not come as a surprise, since the material is sampled from the same area, and it is reasonable to assume that some sample originate from the same source. Although it is interesting to study samples that are correlated within this dataset, this have already been accomplished the multivariate methods just described. The strength of COSIWeb is that it can provide information about similar samples that have been collected in connection to different projects and by different countries.

Interestingly samples 2014-387, 2014-384, 2014-420, 2011-A-VK-2, 2011-B-K-4 were identified as a probable match to samples from the Shetland Islands, namely UK-1-231-4, UK-1-341-3, UK-1-231-5. From the results in this thesis, it has been established that the matching samples likely are crude oils, possibly from the north sea. Sample 2015-846 were identified as a probable match UK-1-332.2. From the result in this thesis it, the sample was characterized as a bunker oil. Similar matches have also been identified by a previous master student ([Vike, 2014])

The spreading of oil samples between these two locations may be a result of drifting of wind currents and coastal currents, such as the North Atlantic current. These results is also an example of how oil spill may affect seabirds, since the samples from Shetland are oil collected from bird feathers.

5.7 Comparing multivariate methods

PCA and HCA categorized many samples into the same groups. For example, group 1 in PCA contain the same samples as group 4 in HCA, although group 4 contain additional samples. HCA were also able to categorize one cluster only consisting of bunker oil. PLS-DA identified samples into crude oil classes, but no samples as bunker oils.

HCA and PCA were able to locate some of the unknown samples into groups that consist of one type of oil, and so the multivariate techniques seems robust to weathered samples. Of the unknown samples 2011-A-VK-2, 2014-395, 2016-153, 2012-267, 2016-158, 2016-160, 2013-637, 2015-881, 2011-BSK-12, 2013-640 are most likely crude oil samples. The two PLS-DA models categorized together 13 unknown oil samples, most likely as crude oil samples (from both origins).

Since there is little knowledge of the sample material, and many samples are very weathered it is difficult to say anything absolute about the correctness of the results. However, these result can be used as a indication for future controlled experiments. The biplot in PCA and correlation loading plot in PLS-DA made it possible to identify diagnostic ratios that are useful in these type of weathered samples.

PCA have in this project been used as a subjective technuqie when it comes to identifying groups of samples. This can be a disadvantage since it depends on the expertise of the person in charge of the analysis. However this subjective evaluation can be verified by different techniques, such as calculcating the RSD for each group and comparing them to the whole sample set. There are more formal statistical techniques that could have been used for selecting groups in PCA, but this was not included in this thesis.

HCA is a method that collects all samples into different clusters, and may be a less time consuming method. Each cluster should ideally represent homogeneous groups, which depends on the right selection of clusters. However in cases were there are little similarities between samples, it may lead the analyst to erroneous conclusions. PLS-DA is an objective technuqie and is able to objectively select samples into groups, but this will depend on how representative the training set is.

Chapter 6

Conclusion

A majority of the samples in this thesis were categorized as crude oil (possibly from the north sea). This seems logical, since there are many oil fields in the north sea, and transportation routes along the Norwegian coastline. In addition some samples were identified with high abundance of oleanane (30 O) and gammacerane (30 G), thus indicating long distance oil migration.

There was not observed any distinct relationship between oil type and sampling location, different oil types were randomly located on the different islands. This have also been shown by a previous master student who inspected samples from 2011 and 2012 [Henriksen, 2012].

PCA, HCA, PLS-DA have demonstrated their ability to categorize weathered samples into groups and clusters, and categorized some unknown samples. Compared to univariate methods, such as the CEN methodology this provides a fast and (to some degree) objective measure of sample similarity, especially when there are many samples in a dataset.

An interesting finding is that the multivariate approaches are able to group samples without using some of the typical identifying biomarkers that are used for inspection of chromatograms, or some of the typical DR used in univariate oil spill forensics (MA, 30O, 30G and retene). This shows that multivariate techniques, could be a promising method for identifying heavily weathered samples that often have inconclusive or missing measurements for these typically used biomarkers and diagnostic ratios.

By applying COSIWeb, external oil spill samples from Shetland proved a match to six samples in the dataset, which makes it possible to form a picture over the distances these samples have traveled. It also shows that oil stranded on shore not necessarily are related to a local incidents.

Chapter 7

Further work

This thesis have investigated samples from 18 different islands around the Froan area during a time period of 2010-2016. During these periods, new islands have almost always been investigated each year. It could be interesting to go back to some of the islands previously visited and investigate if newer oil samples are observed, which could be helpful to say something about the frequency of oil spills along the shore. This have to some extent been studied Sula, but this could be performed on several islands.

It would be interesting to enhance the PLS-DA model by including weathered oil samples of known origin in the training set. However, this may be difficult to achieve since, it would require known samples that have weathered over a long period. It could also be advantageous to include a trainig set with equal sizes of bunker and crude oil to see if this improved the prediction ability of bunker oils.

Samples with correlation 0.98 was investigated in the database COSIWeb, lower correlations may also be looked in to, in case this results in additional matches with external samples.

It would also be interesting to go deeper into specific correlations between diagnostic ratios and samples in PCA biplot and PLS-DA plot, to evaluate if more could be said about their source.

Bibliography

- [Alsberg, 2015] Alsberg, B. K. (2015). *Chemometrics*, volume 0.40. Alsberg research group. Compendium for the NTNU courses TKJ4175/KJ8175.
- [Bonn Agreement, 2016] Bonn Agreement (2016). About Bonn Agreement. <http://www.bonnagreement.org/>. Last accessed on Oct 14, 2016.
- [Bourman, 2009] Bourman, B. (2009). Guidelines for surveying soil and land resources. *Geographical Research*, 47(2):224–225.
- [Brandvik and Daling, 2014] Brandvik, P. J. and Daling, P. S. (2014). *Crude oil composition, properties and laboratory methods to characterise crude oils*. NTNU. Lecture Compendium for KJ3050 Marine Organic Environmental Chemistry, autumn, 2014.
- [CAMO, 1998] CAMO (1998). The unscrambler user manual. *CAMO ASA Norway*.
- [CEN, 2012] CEN (2012). Oil spill identification-waterborne petroleum and petroleum products-part 2: Analytical methodology and interpretation of results based on gc-fid and gc-ms low resolution analyses. Technical report, Standard Norge.
- [Christensen et al., 2004] Christensen, J. H., Hansen, A. B., Tomasi, G., Mortensen, J., and Andersen, O. (2004). Integrated methodology for forensic oil spill identification. *Environmental science & technology*, 38(10):2912–2918.
- [Christensen and Tomasi, 2007] Christensen, J. H. and Tomasi, G. (2007). Practical aspects of chemometrics for oil spill fingerprinting. *Journal of Chromatography A*, 1169(1):1–22.
- [Chromedia, 2016] Chromedia (2016). Picture of a gas chromatogram. <http://www.chromedia.org/>. Last accessed on Oct 14, 2016.

- [COSIWeb, 2016] COSIWeb (2016). Manual for use of the cosiweb database. <http://cosi.bsh.de:8080/CosiWeb/>. Last accessed on Oct 14, 2016.
- [Daling et al., 2002] Daling, P. S., Faksness, L.-G., Hansen, A. B., and Stout, S. A. (2002). Improved and standardized methodology for oil spill fingerprinting. *Environmental Forensics*, 3(3-4):263–278.
- [Esbensen et al., 2002] Esbensen, K. H., Guyot, D., Westad, F., and Houmoller, L. P. (2002). *Multivariate data analysis-in practice: an introduction to multivariate data analysis and experimental design*. Multivariate Data Analysis.
- [Ettre, 1993] Ettre, L. (1993). Nomenclature for chromatography (IUPAC recommendations 1993). *Pure and Applied Chemistry*, 65(4):819–872.
- [Fingas, 2010] Fingas, M. (2010). *Oil spill science and technology*. Gulf professional publishing.
- [Geladi and Kowalski, 1986] Geladi, P. and Kowalski, B. R. (1986). Partial least-squares regression: a tutorial. *Analytica chimica acta*, 185:1–17.
- [Grob and Barry, 2004] Grob, R. L. and Barry, E. F. (2004). *Modern practice of gas chromatography*. John Wiley & Sons, 4th edition edition.
- [Hair et al., 2006] Hair, J., Black, W., Babin, B., Anderson, R., and Tatham, R. (2006). Multivariate data analysis sixth edition pearson education. *New Jersey*, pages 42–43.
- [Henriksen, 2012] Henriksen, S. (2012). Kartlegging og kjemisk karakterisering av oljeforurensing langs trøndelagskysten. Master’s thesis, NTNU, Trondheim.
- [International Maritime Organization, 2016] International Maritime Organization (2016). International convention for the prevention of pollution from ships. [http://www.imo.org/en/About/Conventions/ListOfConventions/Pages/International-Convention-for-the-Prevention-of-Pollution-from-Ships-\(MARPOL\).aspx](http://www.imo.org/en/About/Conventions/ListOfConventions/Pages/International-Convention-for-the-Prevention-of-Pollution-from-Ships-(MARPOL).aspx). Last accessed on Oct 14, 2016.
- [ITOPF, 2016] ITOPF (2016). Data and statistics. <http://www.itopf.com/knowledge-resources/data-statistics/>. Last accessed on Oct 14, 2016.

- [Kao et al., 2015] Kao, N.-H., Su, M.-C., Fan, J.-R., and Chung, Y.-Y. (2015). Identification and quantification of biomarkers and polycyclic aromatic hydrocarbons (pahs) in an aged mixed contaminated site: from source to soil. *Environmental Science and Pollution Research*, 22(10):7529–7546.
- [KJ3050 students, 2011] KJ3050 students (2011). Field report 2011. Internal SINTEF report.
- [KJ3050 students, 2012] KJ3050 students (2012). Field report 2012. Internal SINTEF report.
- [KJ3050 students, 2013] KJ3050 students (2013). Field report 2013. Internal SINTEF report.
- [KJ3050 students, 2014] KJ3050 students (2014). Field report 2014. Internal SINTEF report.
- [KJ3050 students, 2015] KJ3050 students (2015). Field report 2015. Internal SINTEF report.
- [Kumar and Mishra, 2015] Kumar, K. and Mishra, A. K. (2015). Application of partial least square (pls) analysis on fluorescence data of 8-anilino-naphthalene-1-sulfonic acid, a polarity dye, for monitoring water adulteration in ethanol fuel. *Journal of fluorescence*, 25(4):1055–1061.
- [Little, 2016] Little, T. A. (2016). Method validation essentials, limit of blank, limit of detection, and limit of quantitation. <http://www.biopharminternational.com/>. Last accessed on Oct 14, 2016.
- [Liv-Guri Faksness, 2002] Liv-Guri Faksness, Hermann M. Weiss, P. S. D. (2002). Revision of the nordtest methodology for oil spill identification. Technical report, SINTEF.
- [Lundanes et al., 2013] Lundanes, E., Reubsaet, L., and Greibrokk, T. (2013). *Chromatography: basic principles, sample preparations and related methods*. John Wiley & Sons.
- [Martens and Martens, 2001] Martens, H. and Martens, M. (2001). *Multivariate analysis of quality. An introduction*. IOP Publishing.
- [Meyers, 2011] Meyers, R. A. (2011). *Encyclopedia of analytical chemistry*. John Wiley & Sons.

- [Mocak et al., 1997] Mocak, J., Bond, A., Mitchell, S., and Scollary, G. (1997). A statistical overview of standard (iupac and acs) and new procedures for determining the limits of detection and quantification: application to voltammetric and stripping techniques. *Pure and Applied Chemistry*, 69(2):297–328.
- [Nielsen et al., 2012] Nielsen, N. J., Ballabio, D., Tomasi, G., Todeschini, R., and Christensen, J. H. (2012). Chemometric analysis of gas chromatography with flame ionisation detection chromatograms: a novel method for classification of petroleum products. *Journal of Chromatography a*, 1238:121–127.
- [NORDTEST, 1991] NORDTEST (1991). Oil spill identification. Technical report.
- [Norwegian Environment Agency, 2016] Norwegian Environment Agency (2016). The norwegian sea. <http://www.environment.no/topics/marine-and-coastal-waters/the-norwegian-sea/>. Last accessed on Oct 14, 2016.
- [Norwegian Petroleum, 2016] Norwegian Petroleum (2016). Acute pollution and oil spill preparedness and response. <http://www.norskpetroleum.no/en/environment-and-technology/oil-spill-preparedness-and-response/>. Last accessed on Oct 14, 2016.
- [Pearson, 1901] Pearson, K. (1901). On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572.
- [Peters et al., 2005] Peters, K. E., Walters, C. C., and Moldowan, J. M. (2005). *The biomarker guide*, volume 1. Cambridge University Press.
- [Poole, 2003] Poole, C. F. (2003). *The essence of chromatography*. Elsevier.
- [Reimann and Filzmoser, 2000] Reimann, C. and Filzmoser, P. (2000). Normal and lognormal data distribution in geochemistry: death of a myth. consequences for the statistical treatment of geochemical and environmental data. *Environmental geology*, 39(9):1001–1014.
- [Remenyi et al., 2011] Remenyi, D., Onofrei, G., and English, J. (2011). *An introduction to statistics using Microsoft Excel*. Academic Conferences Limited.
- [S. Boitsov and Dolva, 2013] S. Boitsov, J. K. and Dolva, H. (2013). Experiences from oil spills at the norwegian coast. Technical report, Institute of marine research.

- [Simanzhenkov and Idem, 2003] Simanzhenkov, V. and Idem, R. (2003). *Crude oil chemistry*. CRC Press.
- [Skoog, 2004] Skoog, D. A. (2004). *Fundamentals of analytical chemistry*. Grupo Editorial Norma.
- [Speight, 2014] Speight, J. G. (2014). *The chemistry and technology of petroleum*. CRC press.
- [Statoil, 1998] Statoil (1998). Regional konsekvensutredning for Haltenbanken/Norskehavet. Technical report, Statoil.
- [Stout and Wang, 2016] Stout, S. and Wang, Z. (2016). *Standard Handbook Oil Spill Environmental Forensics: Fingerprinting and Source Identification*. Academic Press.
- [Stout et al., 2001] Stout, S. A., Uhler, A. D., and McCarthy, K. J. (2001). A strategy and methodology for defensibly correlating spilled oil to source candidates. *Environmental Forensics*, 2(1):87–98.
- [Sun et al., 2015] Sun, P., Bao, M., Li, F., Cao, L., Wang, X., Zhou, Q., Li, G., and Tang, H. (2015). Sensitivity and identification indexes for fuel oils and crude oils based on the hydrocarbon components and diagnostic ratios using principal component analysis (pca) biplots. *Energy & Fuels*, 29(5):3032–3040.
- [van den Berg et al., 2006] van den Berg, R. A., Hoefsloot, H. C., Westerhuis, J. A., Smilde, A. K., and van der Werf, M. J. (2006). Centering, scaling, and transformations: improving the biological information content of metabolomics data. *BMC genomics*, 7(1):1.
- [Vike, 2014] Vike, K. (2014). Oil spill forensics: Identification of sources for oil spills by using data generated by gc-ms and icp-ms combined with multivariate statistics and the cosiweb database. Master’s thesis, NTNU.
- [Visit Norway, 2016] Visit Norway (2016). Froan nature reserve. <https://www.visitnorway.com/listings/froan-nature-reserve/91327/>. Last accessed on Oct 14, 2016.
- [Wang et al., 1999] Wang, Z., Fingas, M., and Page, D. S. (1999). Oil spill identification. *Journal of Chromatography A*, 843(1):369–411.

- [Wang and Fingas, 2003] Wang, Z. and Fingas, M. F. (2003). Development of oil hydrocarbon fingerprinting and identification techniques. *Marine pollution bulletin*, 47(9):423–452.
- [Wang and Stout, 2010] Wang, Z. and Stout, S. (2010). *Oil spill environmental forensics: fingerprinting and source identification*. Academic Press.
- [Wang et al., 2006] Wang, Z., Stout, S. A., and Fingas, M. (2006). Forensic fingerprinting of biomarkers for oil spill characterization and source identification. *Environmental Forensics*, 7(2):105–146.
- [Waples, 1985] Waples, D. (1985). *Geochemistry in petroleum exploration*. IHRDC Press, Boston, MA.
- [Wold et al., 1987] Wold, S., Esbensen, K., and Geladi, P. (1987). Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52.

Appendix A

Description of Oil Samples

Tables and GC-FID chromatograms used in this thesis are presented here. For more detailed information contact: marie.myrstad@gmail.com

Table A.1: Description of oil samples from field trips based on field reports [KJ3050 students, 2011], [KJ3050 students, 2012] [KJ3050 students, 2013] [KJ3050 students, 2014] [KJ3050 students, 2015] [Vike, 2014]

ID	Field ID	GPS N/S	GPS W/E	Island	Sample Size	Odour	Stickyness	Description	Information	Previous class	Classification
2011-B-K-1	BK01	64,032717	9,156367	Kunna		Strong oil smell		Shiny, Sticky, rubbery	Particles/Contaminants in matrix Hard/old outside, Soft inside	Non-NS-crude	Non-NS-Crude
2011-B-K-2	BK02	64,032717	9,156367	Kunna		Not distinct		Shiny,soft		Unknown	Unknown
2011-B-K-4	BK04	64,032817	9,156233	Kunna		Strong oil smell		Shiny, and sticky	Hard	Crude	Crude
2011-B-K-7	BK07	64,033167	9,156117	Kunna		Slight oil smell		Not shiny, slightly Sticky	Soft, refuses to dissolve properly.	Non-NS-crude	Crude
2011-B-SK-2	BSK02	64,036867	9,135633	Storekalven		Yes, not strongly		Soft	Soft, refuses to dissolve properly	Non-NS-crude	Crude
2011-B-SK-6	BSK06	64,036733	9,136167	Storekalven		Weakly		Rubbery Inside	Semi-Solid, Sticky	Unknown	Unknown
2011-B-SK-12	BSK12	64,036817	9,136767	Storekalven		Not distinct		Rubbery, soft, dense	Hard outside, soft inside, does not dissolve properly.	Unknown	Unknown
2011-B-SK-16	BSK16	64,036967	9,139417	Storekalven		Yes, not strongly		Dense, slightly hard inside	Does not dissolve properly. Large samples size	Non-NS crude	Crude
2011-A-K-20	AK20	64,041367	9,161217	Kunna		No oil-smell		Hard and sticky inside	Soft, refuses to dissolve properly, tarball	Non-NS crude	Crude
2011-A-K-1	AK01	64,042033	9,166983	Kunna		Oily smell		A flat sample. Soft and sticky inside	Feathers/contaminants in matrix	Crude	Crude
2011-A-K-6	AK06	64,042217	9,167067	Kunna		Oily smell		Semisolid and sticky inside	Soft, refuses to dissolve properly, tarball	Non-NS crude	Crude
2011-A-K-8	AK08	64,042467	9,167367	Kunna		Oily smell		Flat and small sample. Hard and sticky inside	Soft, refuses to dissolve properly	Bunker	Unknown
2011-A-K-9	AK09	64,042467	9,167367	Kunna		Oily smell		Irregular shape. Sticky and solid	Soft, refuses to dissolve properly	Crude	Crude
2011-A-K-17	AK17	64,041367	9,161217	Kunna		Slight oil smell		A piece of a tarball. Hard and sticky inside	Soft, refuses to dissolve properly.	Non-NS crude	Crude
2011-A-VK-2	AVK02	64,039617	9,128833	Vesterkalven		Slight oil smell		Thin layer of oil. Sticky inside	Sticky, fair bit of matrix	Crude	Unknown
2011-A-VK-15	A-VK-15	64,038933	9,12825	Vesterkalven		Slight oil smell		FLaky, 1x2 cm. Shiny and sticky inside	Hard and shiny	Bunker	Crude
2011-A-VK-18	A-VK-18	64,0383	9,127467	Vesterkalven		No		6x3 cm. Soft and sticky inside	Lots of matrix contaminants	Non-NS-crude	Crude
2012-AS02	AS02	63,84433333	8,451583333	Sula		Old oil		Bendable, slightly soft	Hard/old outside, sticky inside. Some gravel	Non-NS-crude	Non-NS-crude
2012-AS06	AS06	63,8445	8,450883333	Sula		Oil		Soft, bendable	Feathers in matrix	Bunker	Unknown
2012-AS08	AS08	63,84493333	8,450866667	Sula		Oil		Sticky, bendable, slightly soft	Small, hard lump	Bunker	Unknown
2012-BS01	BS01	63,84365	8,451766667	Sula		Oil		Bendable, very sticky	Some matrix contaminant	Bunker	Bunker
2012-BS03	BS03	63,84373333	8,45215	Sula		Oil		Soft, sticky	Hard, shiny	Non-NS-crude	Crude
2012-BS05	BS05	63,84248333	8,452433333	Sula		Forest Oil		Soft, sticky, bendable	Hard, matte, some contaminants	Bunker	Bunker
2012-0248	LA1	63,7728	8,31145	Kya	15x10	2	2		A fair bit of matrix (contaminants)	Bunker	Bunker
2012-0249	LA2	63,772783	8,311517	Kya	15x15	3	3		A fair bit of matrix (contaminants)	Non-NS crude	Non-NS Crude
2012-0251	LA4	63,77265	8,31155	Kya	5x8	2	3		Lots of contaminants	Crude	Crude
2012-0256	LA9	63,774867	8,310433	Kya	20x30	5	5		Semi-liquid, smooth	Bunker	Bunker
2012-0258	LA11	63,773283	8,309133	Kya	x	3	1		Hard/sticky	Crude	Crude
2012-0262	LA13	63,77205	8,307817	Kya	7x4	2	1		Hard/sticky	Unknown	Unknown
2012-0262	LA15	63,771883	8,308617	Kya	13x7	4	3		"Grainy" texture	Bunker	Bunker
2012-0267	LA20	63,771067	8,308117	Kya	10x15	5	0		Big lump, quite solid	Crude	Unknown
2012-0269	LA22	63,770833	8,308483	Kya	70x8	5+	2		Hard, weathered, liquid inside	Bunker	Bunker
2012-0290	A05	63,773717	8,3134	Kya	12x15x1	2	4		Hard, white on top (fungi?)	Unknown	Unknown
2012-0295	A10	63,774	8,31335	Kya	3x3	0	0		Hard and shiny, but sticky. On the small side	Bunker	Bunker
2012-0272	IB3	63,6887	8,341217	Frÿya	10x8	1	2			Bunker	Bunker
2012-0273	IB4	63,688817	8,341717	Frÿya	8x5.5	1	5		Sticky, stubborn, hard	Bunker	Crude

2012-0276	1B7	63,688317	8,33285	Frÿya	19x111	3	5	Bird feathers, "stubborn" soft	Non-NS crude	Crude
2012-0278	1B9	63,688467	8,3342	Frÿya	12x10	2	3	Hard, slightly sticky top, contaminants and "stubborn". Small sample	Unknown	Unknown
2012-0314	B29	63,691333	8,3446	Frÿya	5x5	1	0	Bird feathers, sticky	non-NS crude	Non-NS Crude
2012-0316	B31	63,690767	8,3443	Frÿya	5x10	1	4		Unknown	Unknown
2013-0673	BA3	8,305583	8,305583	Blfskyÿya	1.5x1	3	1	Soft	Crude	Crude
2013-0680i	BA10	8,312883	8,312883	Blfskyÿya	13x15	5	5	Hard, liquid under	Non-NS Crude	Unknown
2013-0680ii	BA11	8,31095	8,31095	Blfskyÿya	12x15	5	3	Hard, liquid under	Non-NS Crude	Non-NS Crude
2013-0681i	BB3	8,352517	8,352517	Olabusÿya	7x5	5	3	Hard, softer under	Non-NS Crude	Non-NS Crude
2013-0684ii	BB3	8,352517	8,352517	Olabusÿya	7x5	5	4	Soft	Bunker	Bunker
2013-0687	BB6	8,3524	8,3524	Olabusÿya	30x0.4	5	2	Soft	Bunker	Bunker
2013-0697	BB16	8,347733	8,347733	Olabusÿya	5x2	6	3	Liquid	Crude	Crude
2013-0698	BB17	8,352	8,352	Olabusÿya	7x11	4	3	Soft	Crude	Unknown
2013-0703	BD4	63,689217	9,415667	Lyngholm	4x6	4	4	Hard	Unknown	Unknown
2013-0622	A2	63,56235	8,304867	Blfskyÿya	Small	1	1	Medium, semi rubber	Bunker	Bunker
2013-0623	A3	63,56235	8,304917	Blfskyÿya	Small	2	3	Medium	Bunker	Unknown
2013-0626	A6	63,562417	8,3048	Blfskyÿya	Small	2	3	Medium	Non-NS Crude	Unknown
2013-0636	B3	63,584717	8,37215	Burÿya	Medium	3	1	Soft	Unknown	Unknown
2013-0637	B4	63,584917	8,373033	Burÿya	Medium	4	2	Soft	Unknown	Unknown
2013-0640	B7	63,584983	8,372667	Burÿya	Small	2	1	Soft	Crude	Unknown
2013-0643	B10	63,585017	8,3723	Burÿya	Medium	4	4	Medium	Non-NS Crude	Unknown
2013-0650	B18	63,584883	8,375783	Burÿya	Large	4	1	Soft	Non-NS Crude	Unknown
2013-0654	B23	63,584917	8,376217	Burÿya	Medium	2	4	Hard	Non-NS Crude	Crude
2013-0656	B25	63,584917	8,376217	Burÿya	Medium	3	1	Medium, soft inside	Bunker	Bunker
2013-0665	C1	63,683633	9,399517	Storfosna	Small	1	2	Medium	Bunker	Bunker
2014-366	Seaotters 6, Vingleia south	63,915783	8,675767	Vingleia	20x30	oil	2	Semi solid	Bunker	Bunker
2014-367	Seaotters 7, Vingleia south	63,915283	8,675183	Vingleia	4x20	Oil + Asphalt	1	Semi Solid	Bunker	Bunker
2014-368	Seaotters 8, Vingleia south	63,91525	8,675017	Vingleia	4x15	Asphalt	1	Solid	Non-NS Crude	Unknown
2014-371	Seaotters 11, Vingleia south	63,914567	8,674817	Vingleia	10x15	Oil	1	Liquid	Crude	Crude
2014-397	Bear 1, Vingleia	63,916367	8,674717	Vingleia	In water	0	In water	In water	Bunker	Unknown
2014-401	Bear 5a, Vingleia	63,915717	8,6759	Vingleia	8x6	NA	NA	Solid	Non-NS Crude	Unknown
2014-401b	Bear 5b, Vingleia	63,915717	8,6759	Vingleia	5x20	NA	1	Semi Solid	Crude	Crude
2014-408	Bear 12, Vingleia	63,915933	8,676033	Vingleia	20x30	Oil/asphalt	2	Semi Solid	Non-NS Crude	Non-NS Crude
2014-413	Bear 5, Borchholmen	63,941067	8,80805	Borchholmen	5x7	Oil	1	Semi Solid	Non-NS Crude	Crude
2014-414	Bear 6, Borchholmen	63,941217	8,774433	Borchholmen	20x60	Oil/asphalt	2	Semi Solid	Non-NS Crude	Crude
2014-420	Bear 12, Borchholmen	63,9418	8,773417	Borchholmen	10x7	NA	0	Solid	Crude	Crude
2014-421	Bear 13, Borchholmen	63,941767	8,773417	Borchholmen	7x3	NA	1	Semi Solid	Non-NS Crude	Crude
2014-384	Seaotters 24, Gfrdsÿya	63,936667	8,764583	Gfrdsÿya	50x20	Asphalt	1	Semi Solid	Crude	Crude
2014-386	Seaotters 26, Gfrdsÿya	63,9366	8,764567	Gfrdsÿya	15x15	Asphalt	0	Solid	bunker	Unknown
2014-387	Seaotters 27, Gfrdsÿya	63,9366	8,764567	Gfrdsÿya	10x25	Asphalt	0	Semi Solid	Crude	Crude
2014-390	Seaotters 30, Gfrdsÿya	63,936717	8,764683	Gfrdsÿya	Tarball 1: 2x2 Tarball 2: 4x2	Oil	0	Solid	Non-NS Crude	Crude
2014-394	Seaotters 34, Krlkvlg NW	63,654483	9,326467	Krlkvlg	200x20	Oil	0	Solid	Crude	Bunker
2014-395	Seaotters 35, Krlkvlg NW	63,654483	9,327133	Krlkvlg	4x3	Oil	0	Semi Solid	Unknown	Unknown

2014-424	Bear 2, Kríkvíg	63,649417	9,316917	Kríkvíg	5x6	NA	1	Semi Solid			Non-NS Crude	Crude
2014-424b	Bear2b, Kríkvíg	63,649417	9,316917	Kríkvíg	5x6	NA	1	Solid			Non-NS Crude	Crude
2014-425	Bear 3, Kríkvíg	63,649417	9,316867	Kríkvíg	30X8	NA	1	Solid			Non-NS Crude	Crude
2015-0826	LA8	63,884444	8,53	Vassfya	7 x 5 (tarball), 10 x 6 (tarball)			Compact, hard outside	Two tarballs, not in splash zone		Non-NS crude	Non-NS Crude
2015-0829	1B2	63,873889	8,486389	Likfya	8 x 6 x 2			Hard outside, soft inside	Tarball		crude	Crude
2015-0844	1B7	63,874444	8,4875	Likfya	5 x 5 (tarball)			Hard outside			crude	Crude
2015-0846	1B9	63,875278	8,488889	Likfya	11 x 5 x 4			Medium			Bunker	Bunker
2015-0847	1B10	63,875278	8,488889	Likfya	35 x 20 x 4			Medium/soft	Dark brown color		Non-NS crude	Crude
2015-0848	1B11	63,875278	8,488889	Likfya	13 x 8 x 0.5			Hard	Grey outside, black inside		Non-NS crude	Crude
2015-0857	2A3	63,885833	8,517778	Geitungan	20 x 10			Hard top, sticky under			Bunker	Bunker
2015-0864	2A10	63,885833	8,517778	Geitungan	15 x 5 x 5			Tarball			Non-NS crude	Unknown
2015-0866	2A12	63,885556	8,518056	Geitungan	10 x 8			Medium/hard			Crude	Unknown
2015-0874	2B6	63,865	8,486667	Gldklakken	20 x 25 x 1.5			Sticky	A piece of glass stuck to the sample		Non-NS crude	Unknown
2015-0879	2B11 (a, b, c)	63,865278	8,488056	Gldklakken	a=6 x 5 x 4, b=5 x 3 x 3			Tarball, Hard outside	One of them are more black and less porous		Non-NS crude	Crude
2015-0881	2B13	63,865	8,488056	Gldklakken	c=4 x 2 x 1 30 x 4 x 1.5			Tarball	Found by soil and grass areas		Crude	Unknown
2015-0882	2B15	63,864722	8,485	Gldklakken	10 x 6 x 2			Hard top, sticky under			Non-NS crude	Non-NS Crude
2015-837	A19	63,886667	8,529444	Vassfya	30 x 15			Very soft	Sticky, looks fresh		Bunker	Unknown
2016-0152	1	63,84889	8,456	Sula	2X2	Asphalt	0	Solid	Gravel		Crude	Crude
2016-0153	3	63,84883	8,4552	Sula	12x7	Asphalt	1	Solid	Moss		Unknown	Unknown
2016-0154	4	63,84883	8,4552	Sula	17x4	Asphalt	1	solid	Moss		Unknown	Unknown
2016-0155	7	63,84874	8,45196	Sula	21x14	Oil and Asphalt	1	semi			Crude	Crude
2016-0156	8	63,84872	8,45192	Sula	12x14	Oil	2	semi			Crude	Crude
2016-0157	9	63,84868	8,45193	Sula	10x10	Oil	2	semi			Crude	Crude
2016-0158	10	63,84852	8,45033	Sula	10x4	Oil	1	semi			Crude	Crude
2016-0159	12	63,84853	8,45024	Sula	21x6	Oil	1	semi			Crude	Crude
2016-160	13	63,84853	8,45024	Sula	10x9	Oil and Asphalt	2	semi			Crude	Crude
2016-161	14	63,8484	8,44953	Sula	6x5	Oil and Asphalt	1	semi			Crude	Crude
2016-162	15	63,84835	8,44982	Sula	14x8	Asphalt	1	semi			Crude	Crude
2016-163	16	63,84832	8,44945	Sula	27x16	Oil	1	semi	Bird feather		Bunker	Bunker
2016-164	18	63,84566	8,44786	Sula	8x5	Oil	1	semi	Moss		Crude	Crude
2016-165	19	63,84567	8,44913	Sula	NA	Oil	1	solid	Plastic rope		Crude	Crude
2016-166	23	63,84418	8,44928	Sula	13x7	Asphalt	1	semi	Grass		Crude	Unknown

Appendix B

PAH and Biomarkers

Table B.1: Target PAH and biomarker compounds

Target PAH and Biomarkers analysed by GC-MS-SIM	Abbreviations
Target PAH (Groups and Single compounds)	
Naphthalene	N0
C1-Naphthalenes	N1
C2-Naphthalenes	N2
C2-Dibenzothiophenes	D2
C3-Dibenzothiophenes	D3
C2-Phenanthrenes/Anthracenes	P2
C3 Phenanthrenes/Anthracenes	P3
C4-Phenanthrenes/Anthracenes	P4
C2-highest peak Phenanthrenes/Anthracenes	
C1-Fluorenes	F1
C2-Fluorenes	F2
C2 benzothiophenes	C2-bt
C2 Fluoranthrenes/Pyrenes	FP2
C1 chrysenes/benzanthracenes	C1
C1-dekalin	DE1
2-Methylanthracene	2MA
4-Methyl Dibenzothiophene	4MD
1-Methyl Dibenzothiophene	1MD
2-methyl phenanthrene	2MP
1-methyl phenanthrene	1MP
Retene	R

Tetramethyl-phenanthrene	T-M-phe
Benzo[b]naphtho(1,2-d)thiophene	BNT
Benzo(b+c)fluorene	B(b+c)
2-methylpyrene	2Mpy
4-methylpyrene	4Mpy
1-methylpyrene	1Mpy
2-Methylfluoranthene	2MFL
Benzo(a)-fluorene	B(a)F
Pentacyclic triterpanes (hopanes)	
18 α (H)-22,29,30-trisnorhopane	27-TS
17 α (H)-22,29,30-trisnorhopane	27-TM
17 α (H),21 β (H)-28,30-bisnorhopane	28ab
17 α (H),21 β (H)-30-norhopane	29ab
18 α (H)-30-norhopane	29Ts
18 α (H)-oleanane	30O
17 α (H),21 β (H)-hopane (hop)	30ab
Gammacerane	30G
17 α (H),21 β (H)-30-homohopane	31abS
17 β (H),21 α (H)-moretane	30ba
17 α (H),21 β (H), 22S-bishomohopane	32abS
Tricyclic Triterpanes	
C23 Tricyclic triterpane	C23 Tr
C24 Tetracyclic triterpane	C24 Tr
C25 Tricyclic triterpane	C25 Tr
C28 Tricyclic triterpanes	C28 (22S)
C29 Tricyclic triterpanes	C29 (22S)
Triaromatic Steroids	
C20-Triaromatic sterane	C20TA
C21-Triaromatic sterane	C21TA
C26,20S-Triaromatic sterane	SC26TA
C28,20S-Triaromatic sterane	SC28TA
C27,20R-Triaromatic sterane	RC27TA
C28,20R-Triaromatic sterane	RC28TA
C26,20R- + C27, 20S-Triaromatic sterane	RC26TA+SC27TA
Steranes	
C27 13 β (H), 17 α (H), 20S-Diacholestane (diasterane)	C27dbS
C27 13 β (H), 17 α (H), 20R-Diacholestane (diasterane)	C27dbR
C28 24-methyl-5 α (H),14 α (H), 17 α (H), 20R-cholestane	C28aaR

C29 24-ethyl-5 α (H), 14 α (H), 17 α , 20S-cholestane	C29aaS
C29 24-ethyl-5 α (H), 14 β (H), 17 β , 20S+R-cholestane	C29bbR+S
C29 24-ethyl-5 α (H), 14 α (H), 17 α , 20R-cholestane	C29aaR
C27 5 α (H), 14 β (H), 17 β (H), 20R+S-cholestane	C27bbR+S
C29 24-ethyl-5 α (H), 14 β (H), 17 β (H), 20R+S-cholestane	C29bbR+S
C29 24-methyl-5 α (H),14 β (H), 17 β (H), 20R+S-cholestane	C28bbR+S
Sesquiterpanes	
C15H28-sesquiterpanes	SES1
C15H28-sesquiterpanes	SES2
C15H28-8 β (H)-drimane	SES3
C15H28-sesquiterpanes	SES4
C16H30-8 β (H)-homodrimane	SES8
Selected n-Alkanes and Acyclic Isoprenoids	
Heptadecane	C17
Pristane	Pri
Octadecane	C18
Phytane	Phy

Appendix C

Diagnostic ratios

Table C.1: Diagnostic ratios

Recommended diagnostic ratios for PAH		
Ratio name	Definition	Ion mass
DR-2-MP/1-MP	2-methylphenanthrene/1-methylphenanthrene	192
DR-MA/1-MP	Methylanthracene/11-methylphenanthrene	192
DR-4-MDBT/1-MDBT	4-methyldibenzothiophene/ 1-methyldibenzothiophene	198
DR-C2-DBT/C2-phe	C2-dibenzothiophenes/C2-phenanthrenes	212/206
DR-C3-DBT/C3-phe	C3-dibenzothiophenes/C3-phenanthrenes	226/220
DR-Retene/C4-phe	Retene/C4-phenanthrenes	234
DR-Retene/T-M-phe	Retene/Tetra-methyl-phenantrene	234
DR-BNT/T-M-phe	Retene/ Tetra-methyl- phenantrene	234
DR-2MFL/4-MPy	2-Methylfluoranthene/4-methylpyrene	216
DR-B(a)F/4-MPy	Benzo(a)fluorene/4-methylpyrene	216
DR-B(b+c)F/4-MPy	Benzo(b+c)fluorene/4-methylpyrene	216
DR-2MPy/4-MPy	2-methylpyrene/4-methylpyrene	216
DR-1MPy/4-Mpy	1-methylpyrene/4-methylpyrene	216
DR-C23Tr/C2-phe-hp	C23 Tricyclic diterpane/Highest peak of the C2-phenantrenes/anthracenes.	191

Table C.2: Diagnostic ratios of biomarkers. Ratios colored in red was not used in this thesis.

Recommended diagnostic ratios of Sesquiterpanes		
Ratio name	Definition	Ion mass
DR-SES 1/3	SES1/SES3	123
DR-SES 2/3	SES2/SES3	123
DR-SES 4/3	SES4/SES3	123
DR- SES 8/3	SES8/SES3	123
Diagnostic ratios for the tricyclic terpanes and hopanes		
Ratio name	Definition	Ion mass
DR-C28/30ab	C28(22S)/30ab	191
DR-27Ts/30ab	27Ts/30ab	191
DR-27Tm/30ab	27Tm/30ab	191
DR-28ab/30ab	28ab/30ab	191
DR-29ab/30ab	29ab/30ab	191
DR-Ts/Tm	Ts/Tm	191
DR-29Ts/30ab	29Ts/30ab	191
DR-30O/30ab	30O/30ab	191
DR-30ba/30ab	30ba/30ab	191
DR-31abS/30ab	31abS/30ab	191
DR-30G/30ab	30G/30ab	191
Diagnostic ratios of the regular steranes and diasteranes		
Ratio name	Definition	Ion mass
DR-27dbR/27dbS	27dbR/27dbS	217
DR-29 α S/29 α R	29 α S/29 α R	217
Diagnostic ratios of Triaromatic Steroids		
Ratio name	Definition	Ion mass
DR-C20TA/C21TA	C20TA/C21TA	231
DR-C21TA/RC26+SC27	C21TA/RC26TA +SC27TA	231
DR-SC26/RC26+SC27	SC26TA/RC26TA +SC27TA	231
DR-SC28/RC26+SC27	SC28TA/RC26TA+SC27TA	231
DR-RC27/RC26+SC27	RC27TA/RC26TA +SC27TA	231
DR-RC28/RC26+SC27	RC28TA/RC26TA+SC27TA	231
DR-C21TA/RC28TA	C21TA/RC28TA	231
DR-SC26TA/SC28TA	SC26TA/SC28TA	231

DR-RC27TA/RC28TA RC27TA/RC28TA 231

Diagnostic ratios of selected n-Alkanes and Acyclic Isoprenoids

Ratio name	Definition	Ion mass
DR- C_{17} /pri	n-heptadecane/pristane	85
DR- C_{18} /phy	n-octadecane/phytane	85
DR-pr/phy	pristane/phytane	85

Appendix D

GC-FID

Table D.1: Description of GC-FID chromatograms.

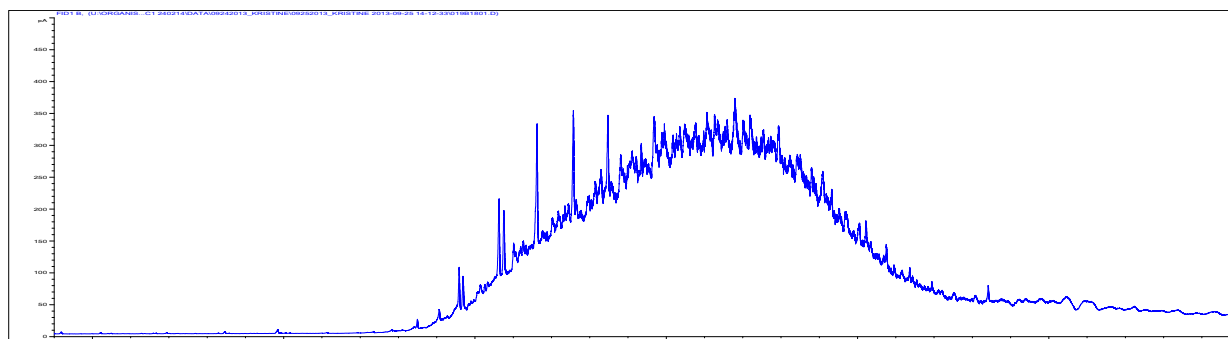
Sample ID	Approximate Carbon range	UCM hump
2014-366	c15-c30	broad,medium
2014-367	c20-c35	broad,medium
2014-368	c20-c35	narrow,high
2014-371	c21-c35	broad, medium
2014-384	c15-c35	broad,flat
2014-0386	c15-c30	broad,medium
2014-0387	c15-c35	broad, flat
2014-390	c15-c40	broad, flat
2014-0394	no carbon range observed even with SPE cleanup	narrow, high
2014-0395	c21-c40	broad, flat
2014-0397	no carbon range observed	
2014-0401	no carbon range observed	narrow, high
2014-0401b	c25-c35	flat
2014-0408	C30-C40	broad, flat
2014-0413	C23-C40	flat
2014-0414	no carbon range observed	broad, medium
2014-0420	no carbon range observed	broad, medium
2014-0421	C23-C40	broad, flat
2014-0424	C22-C34	flat
2014-0424b	C21-C34	flat
2014-0425	no carbon range observed	broad, medium
2015-864	c26-c46	broad, flat

2015-866	c20-c40	Narrow, high
2015-857	c20-c40	broad, high
2015-879	C20-C40	broad,high
2015-881	c24-c40	broad, flat
2015-882	c24-c46	broad, flat
2015-874	c23-c40	narrow, flat
2015-837	c30-c44	broad, medium
2015-839	c17-C44	broad, flat
2015-844	C17-C43	Broad, flat
2015-846	C15-C40	Broad, high
2015-847	C18-C46	broad, flat
2015-848	C24-C46	broad, flat
2013-622	C25-C46	broad, medium
2013-623	C24-C45	broad, medium
2013-626	C18-C40	narrow, high
2013-636	C18-C40	narrow, high
2013-637	C29-C40	broad, flat
2013-640	C18-C40	broad, flat
2013-643	C30-C45	broad, high
2013-650	C25-C45	narrow, high
2013-654	C21-C45	broad, flat
2013-656	C30-C45	narrow, high
2013-665	C21-C45	narrow, high
2013-673	C18-C40	broad, flat
2013-680i	C27-C45	narrow, high
2013-680ii	C24-C40	narrow, medium
2013-681i	C28-C40	narrow, medium
2013-681ii	C30-C45	broad, flat
2013-684i	C23-C45	narrow, high
2013-684ii	C23-C45	narrow, high
2013-687	C22-C45	narrow, high
2013-697	C17-C45	broad, medium
2013-698	c16-C45	broad, flat
2013-703	C25-C45	broad, flat
2016-152	c30-c42	broad, medium
2016-153	c28-c42	broad, medium
2016-154	SPE cleanup is needed	
2016-155	SPE cleanup is needed	

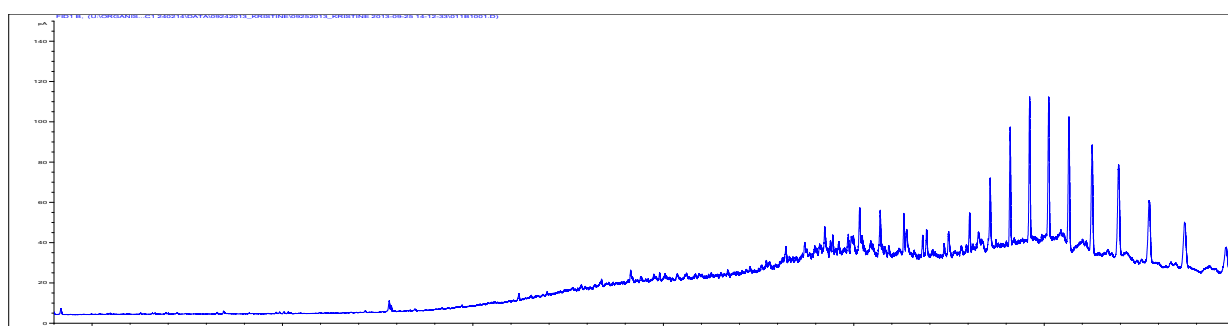
2016-156	C31-C41	broad, medium
2016-157	C29-C41	broad, medium
2016-158	SPE cleanup is needed	
2016-159	C26-C43	broad, medium
2016-160	SPE cleanup is needed	
2016-161	C31-C43	broad, medium
2016-162	C28-C43	broad, flat
2016-163	SPE cleanup is needed	
2016-164	C27-C43	broad, flat
2016-165	C30-C43	
2016-166	SPE cleanup is needed	
2012-AS02	C29-C45	broad, medium
2012-AS06	no carbon range observed	broad, medium
2012-AS08	no carbon range observed	broad, medium
2012-BS01	C21-C45	broad, medium
2012-BS03	C40-C45	broad, flat
2012-BS05	no carbon range observed	
2011-BK01	C26-C40	broad, medium
2012-BK02	C21-C45	broad, medium
2012-BK04	C18-C45	broad, medium
2012-BK07	C26-C45	broad, flat
2012-BSK02	C14-C45	broad, medium
2012-BSK06	C25-C45	broad, medium
2012-BSK12	C28-C45	broad, flat
2012-BSK16	C18-C36	narrow, high
2012-AK20	C28-C40	broad, flat
2012-AK01	C26-C45	broad, medium
2012-AK06	C24-C45	broad, flat
2012-AK08	C22-C45	broad, flat
2012-AK09	C17-C45	broad, flat
2012-AK17	C30-C40	broad, flat
2012-AVK02	C17-C45	broad, medium
2012-AVK15	C30-C45	broad, flat
2012-AVK18	C26-C45	broad, flat
2012-0248	C17-C45	broad, medium
2012-249	C18-C45	broad, flat
2012-0251	C24-C45	broad, flat
2012-0256	C14-C45	broad, medium

2012-0258	C25-C45	broad, flat
2012-0260	C17-C45	broad, medium
2012-0262	C18-C45	broad, medium
2012-0267	C28-C45	broad, medium
2012-0269	C17-C45	broad, medium
2012-0272	C21-C45	broad, medium
2012-0273	C24-C45	broad, medium
2012-0276	C17-C45	broad, medium
2012-0278	C30-C45	broad, flat
2012-0290	C24-C45	broad, flat
2012-0295	no carbon range observed	broad, narrow
2012-314	C19-C45	broad, medium
2012-0316	C16-C45	broad, medium

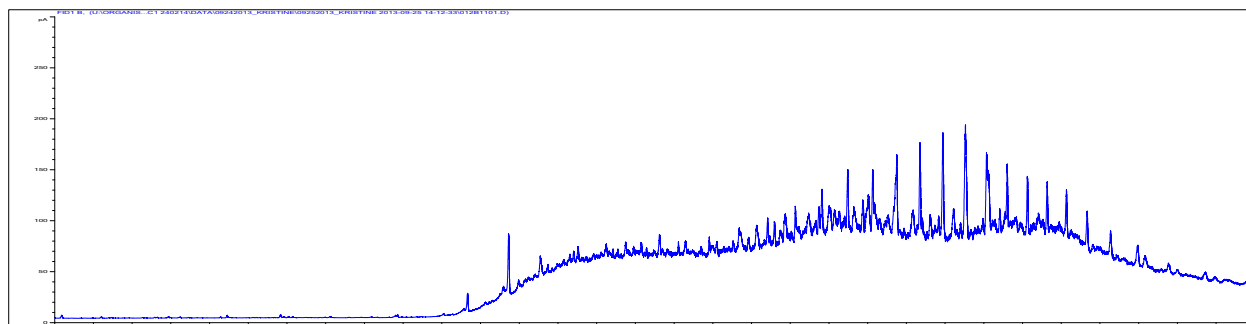
GC-FID 2011



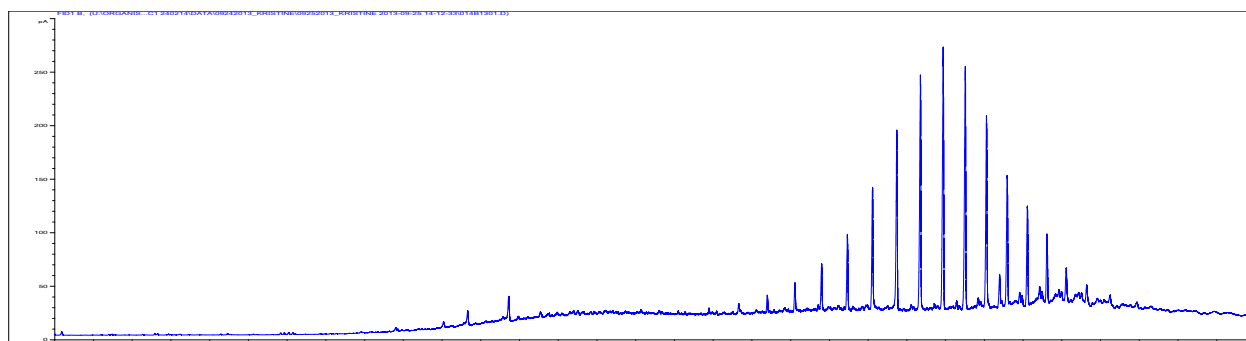
2011-BK01



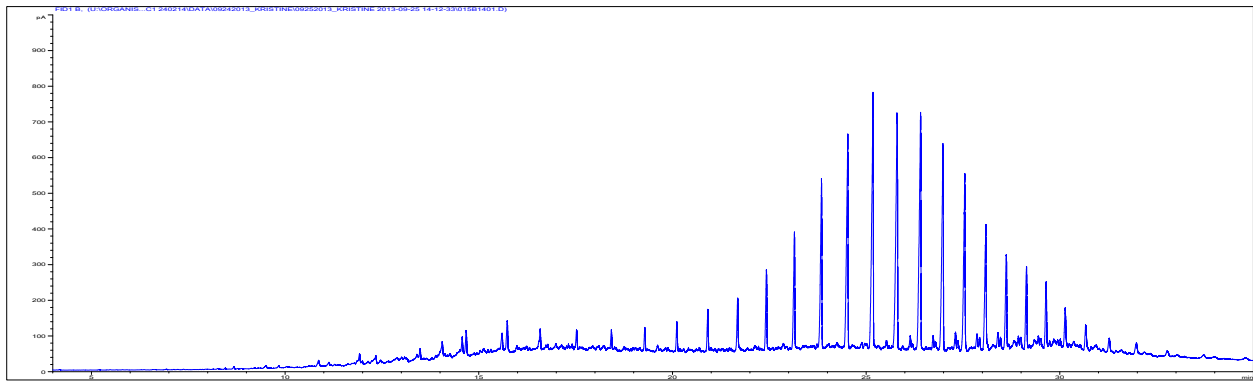
2011-BK02



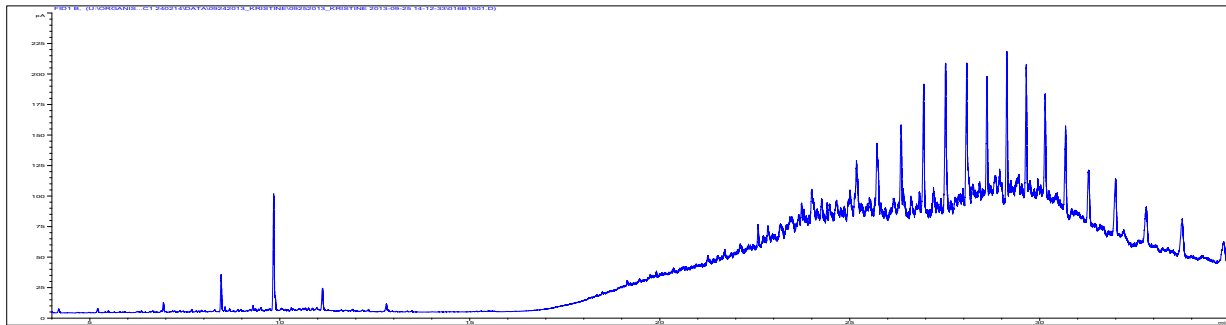
2011-BK04



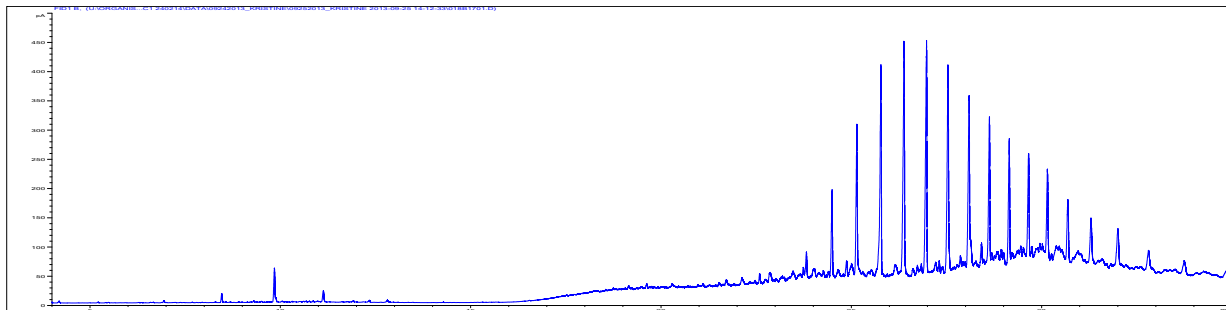
2011-BK07



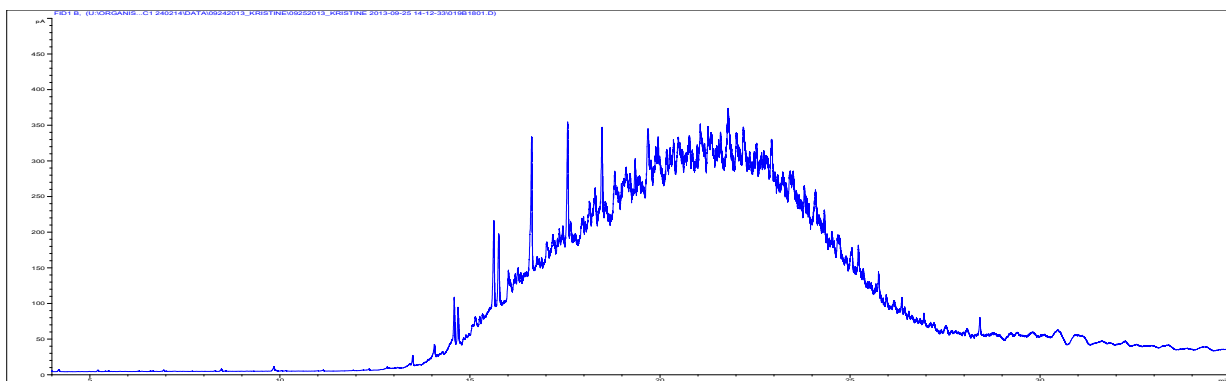
2011-BSK02



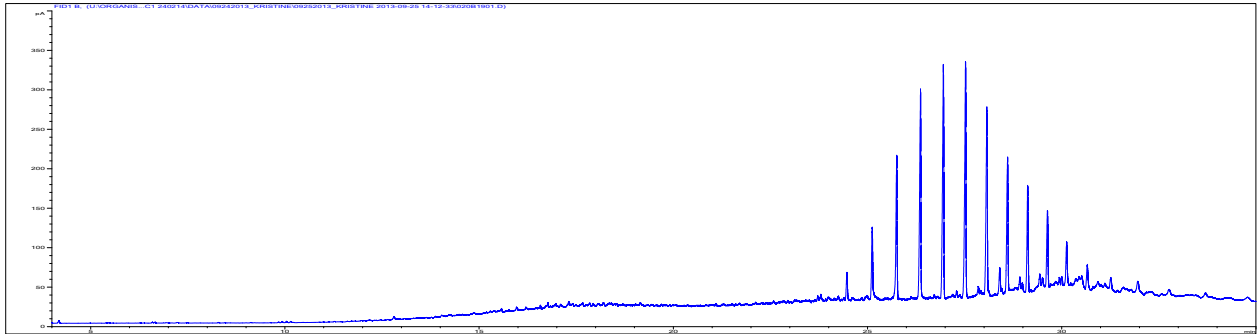
2011-BSK06



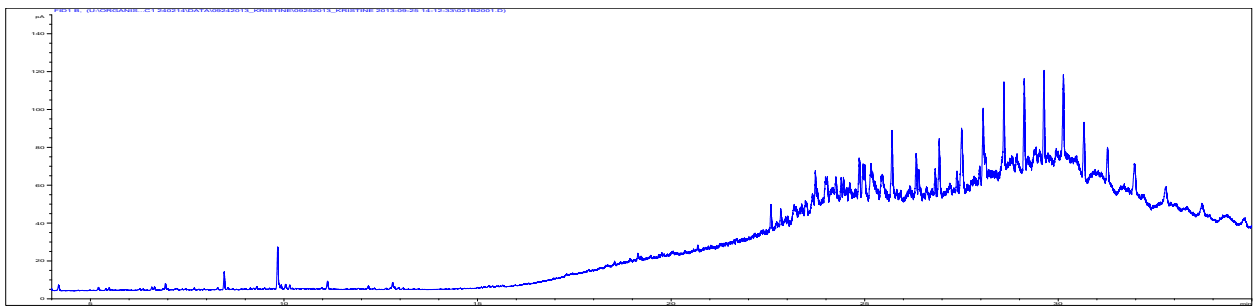
2011-BSK12



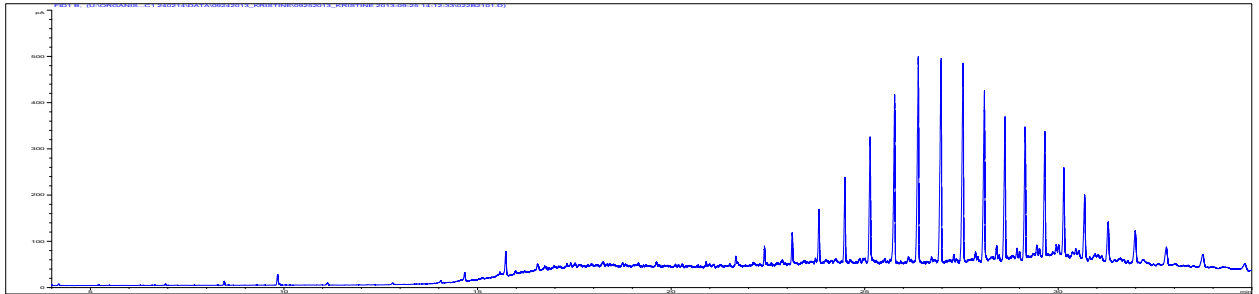
2011-BSK16



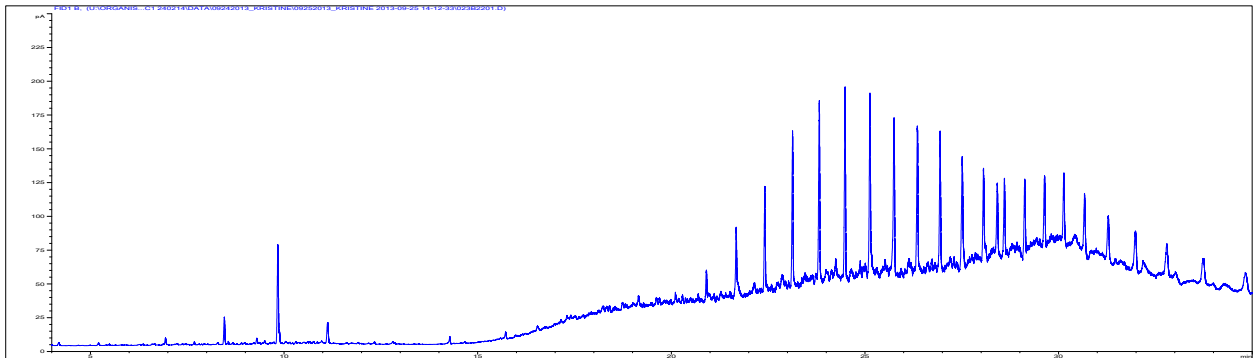
2011-AK20



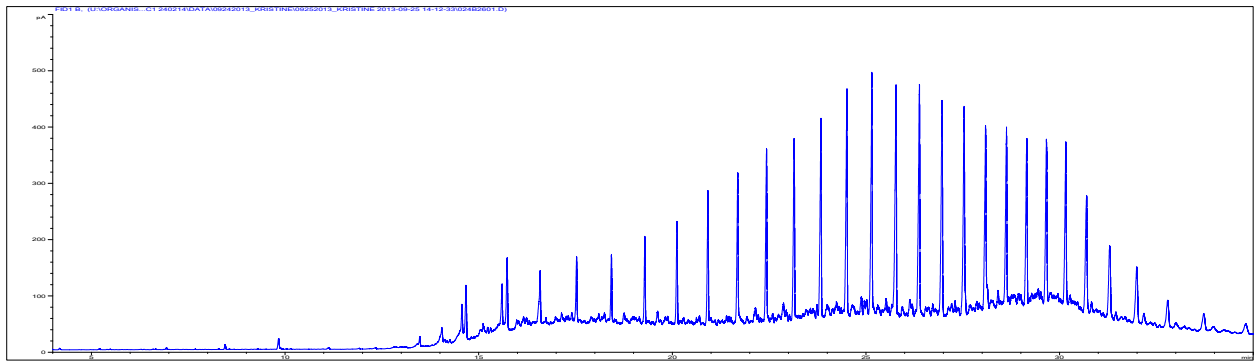
2011-AK01



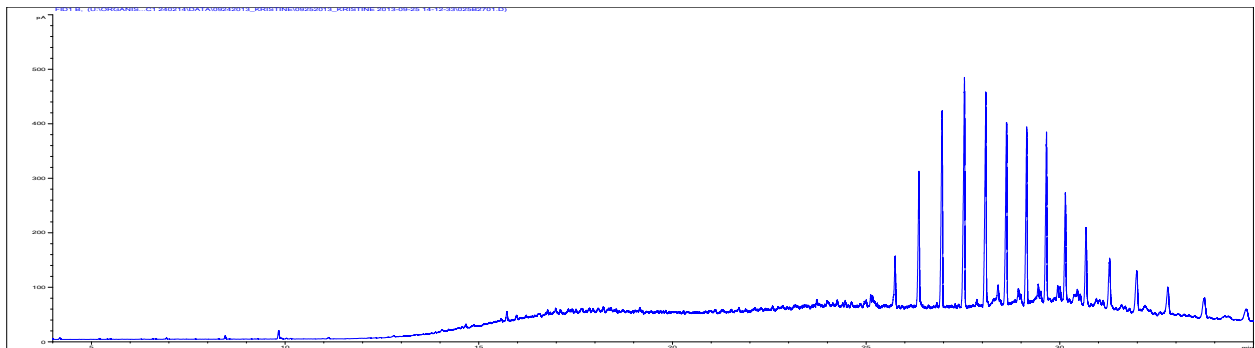
2011-AK06



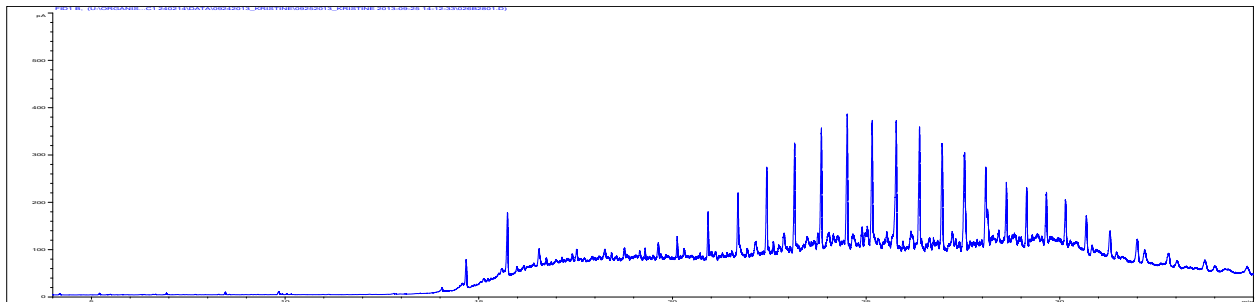
2011-AK08



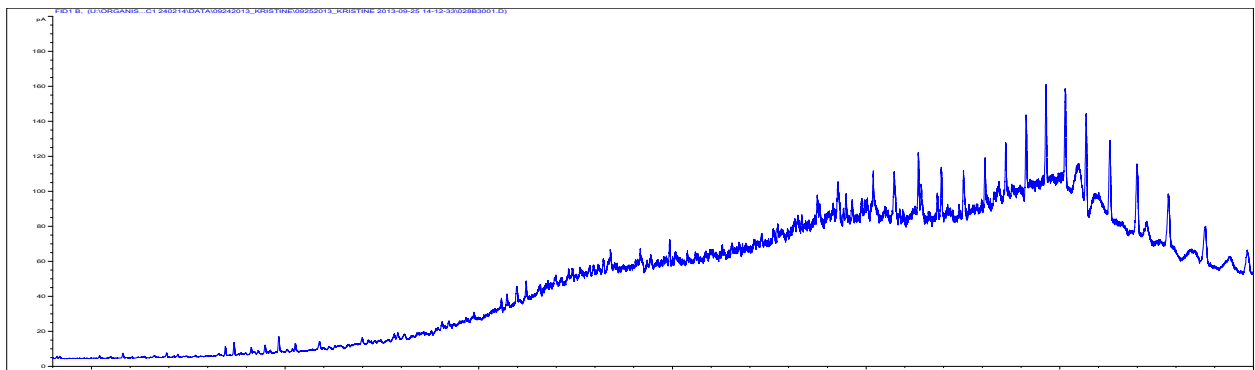
2011-AK09



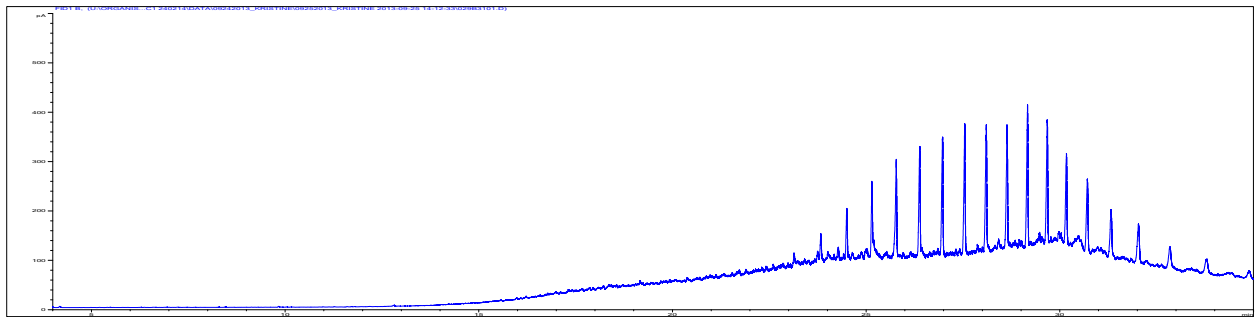
2011-AK17



2011-AVK02

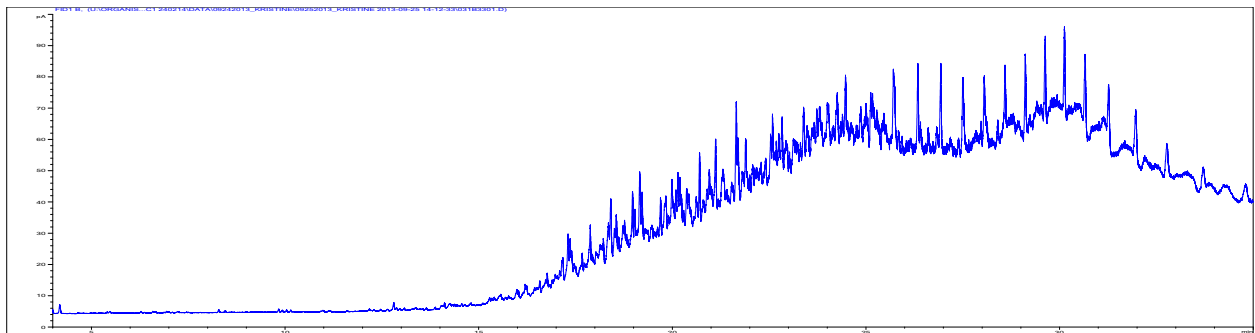


2011-AVK15

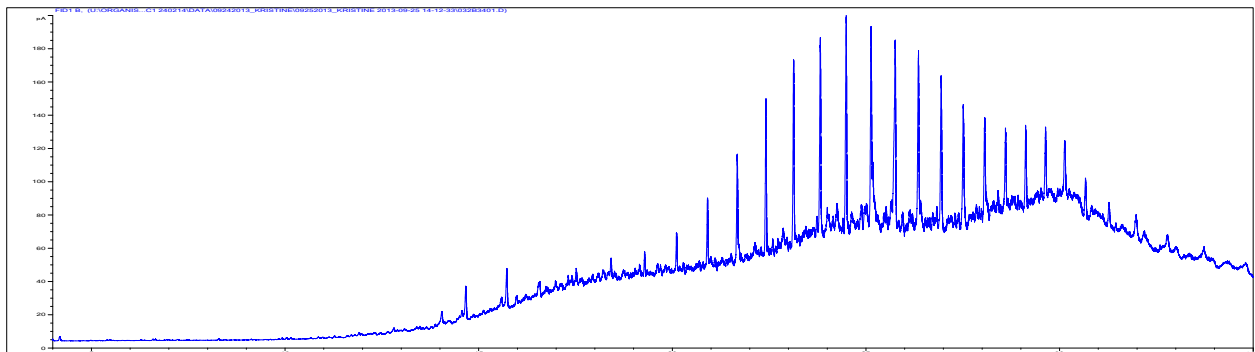


2011-AVK18

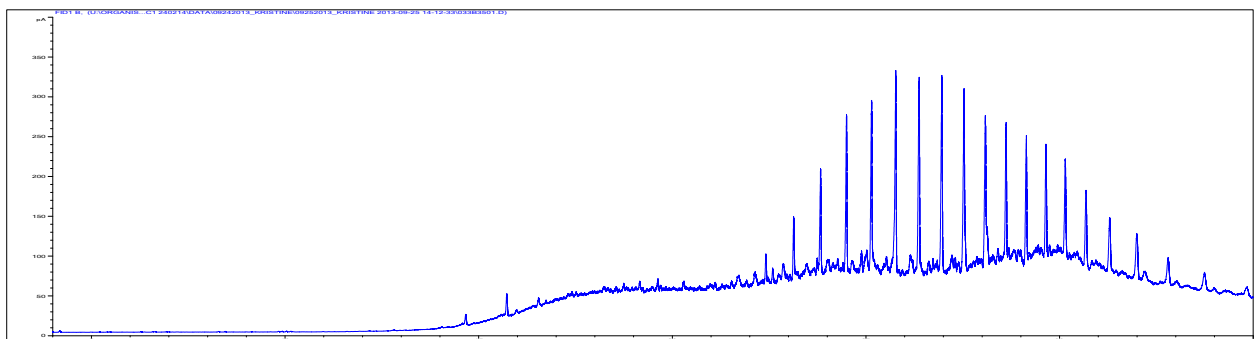
GC-FID 2012



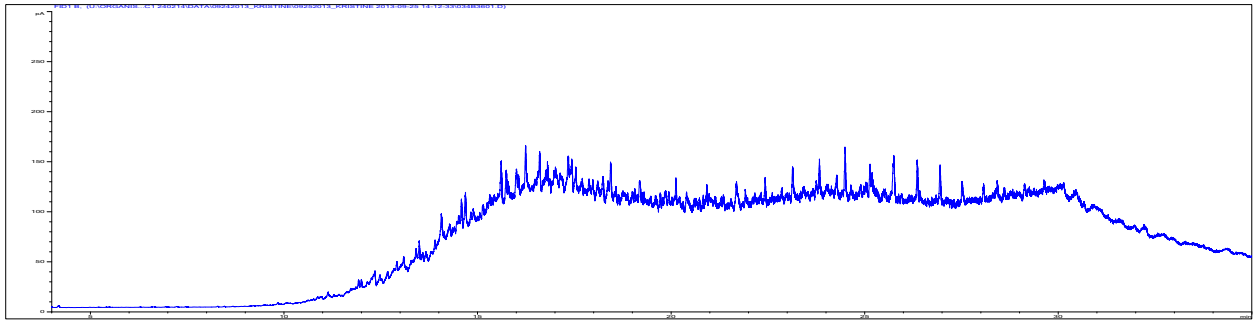
2012-0248



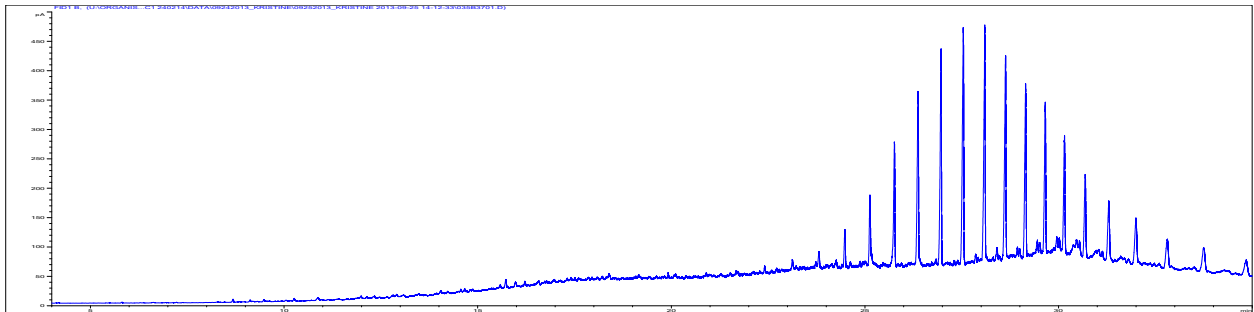
2012-249



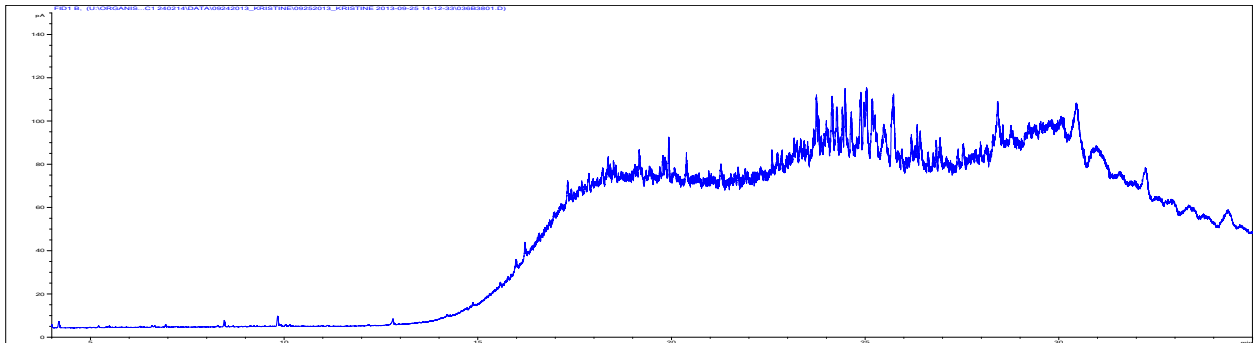
2012-0251



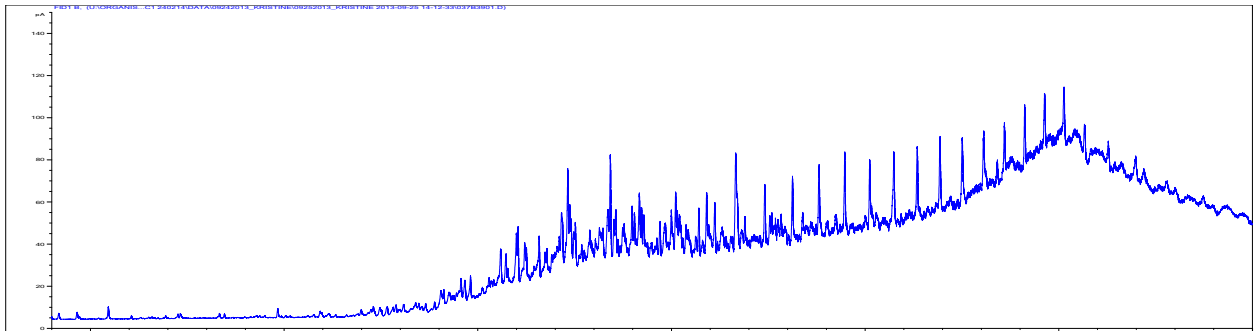
2012-0256



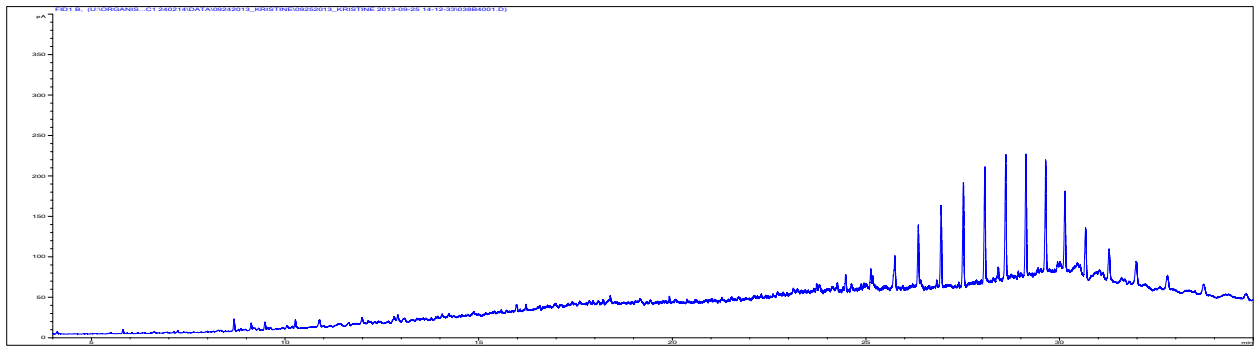
2012-0258



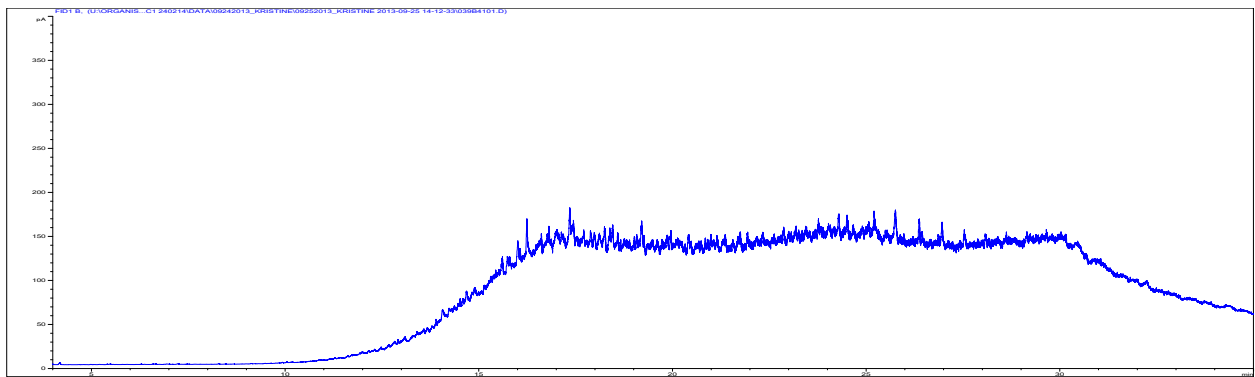
2012-0260



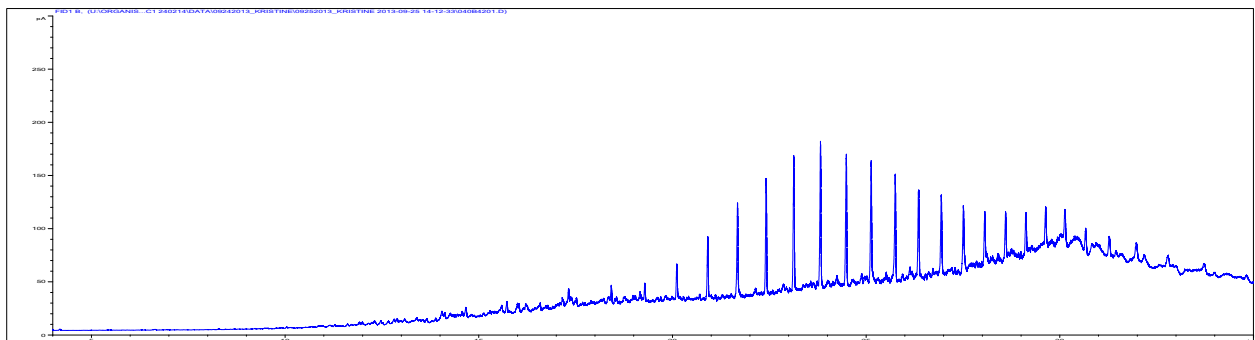
2012-0262



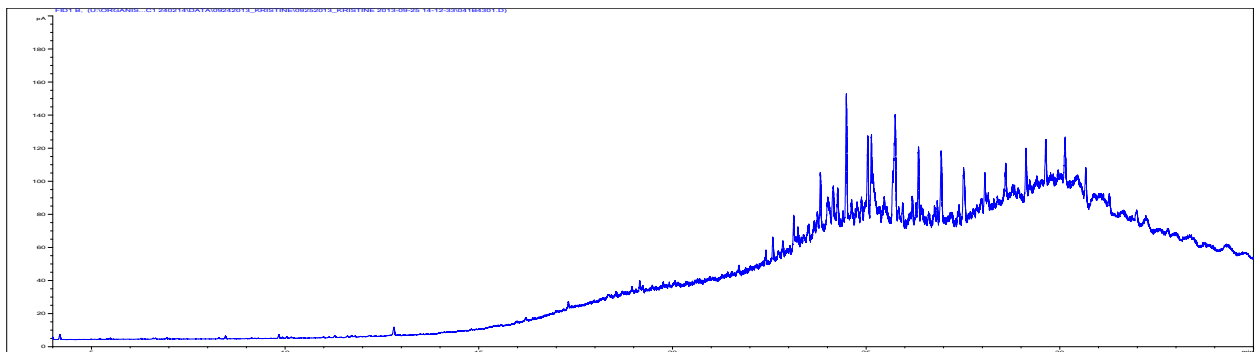
2012-0267



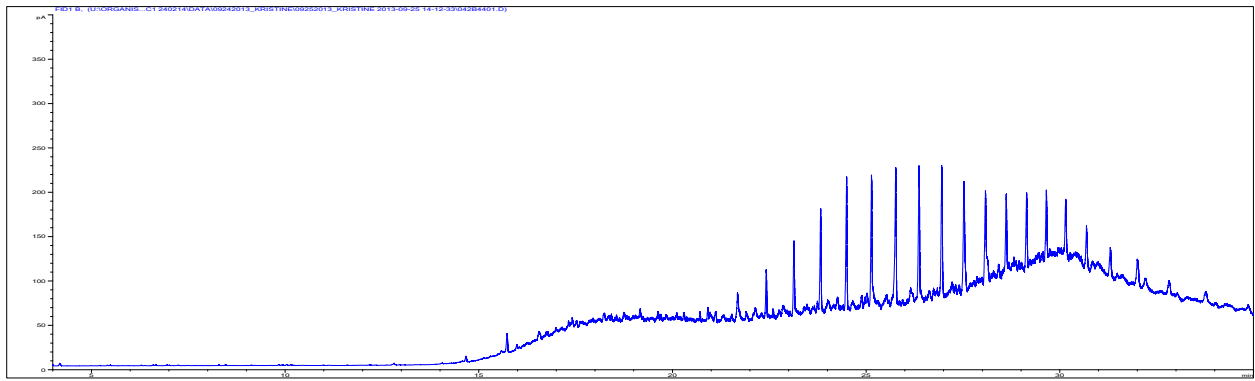
2012-0269



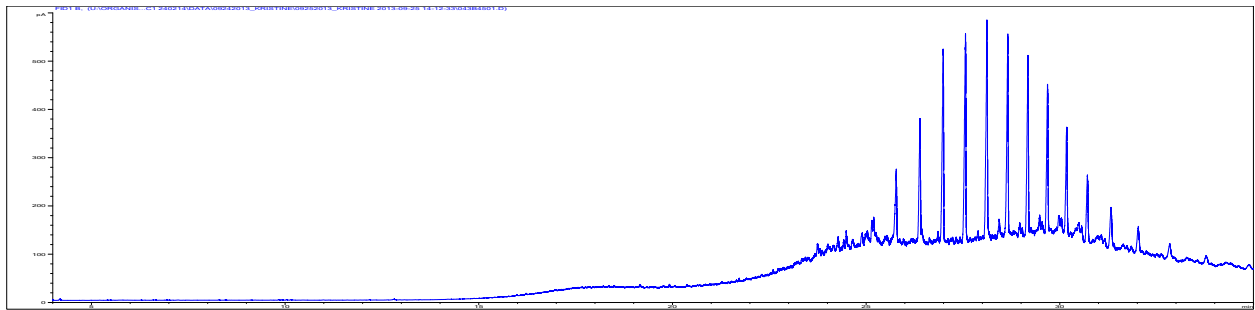
2012-0272



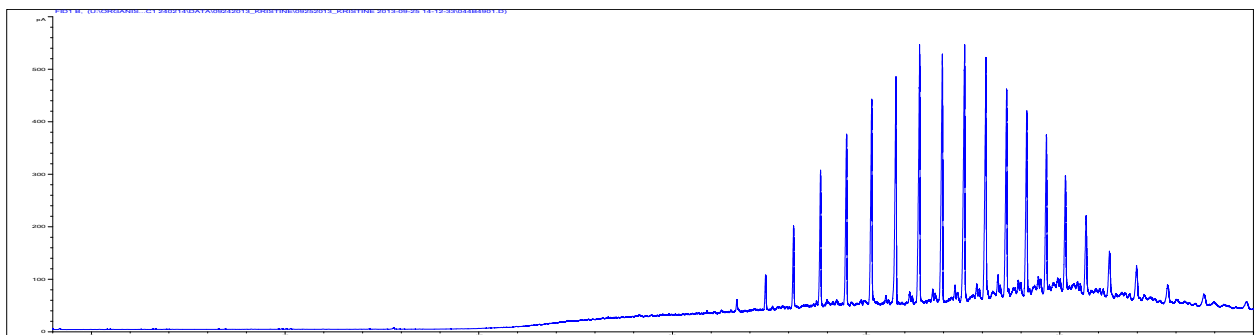
2012-0273



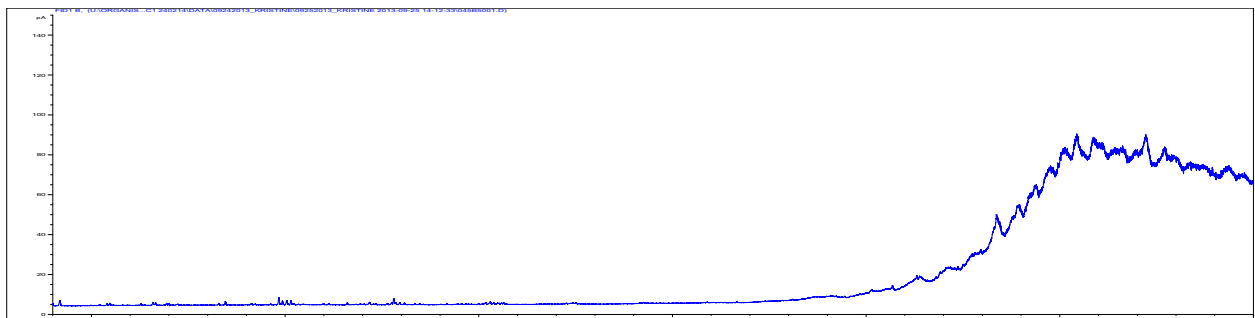
2012-0276



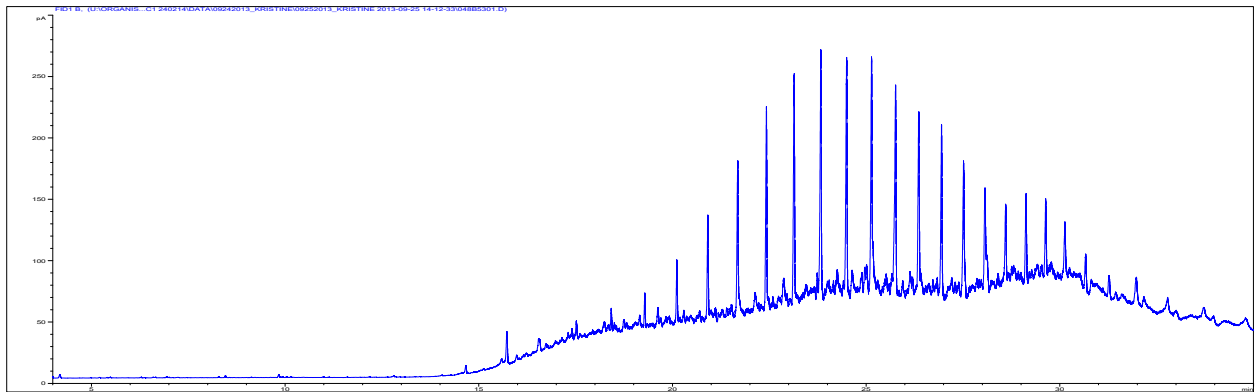
2012-0278



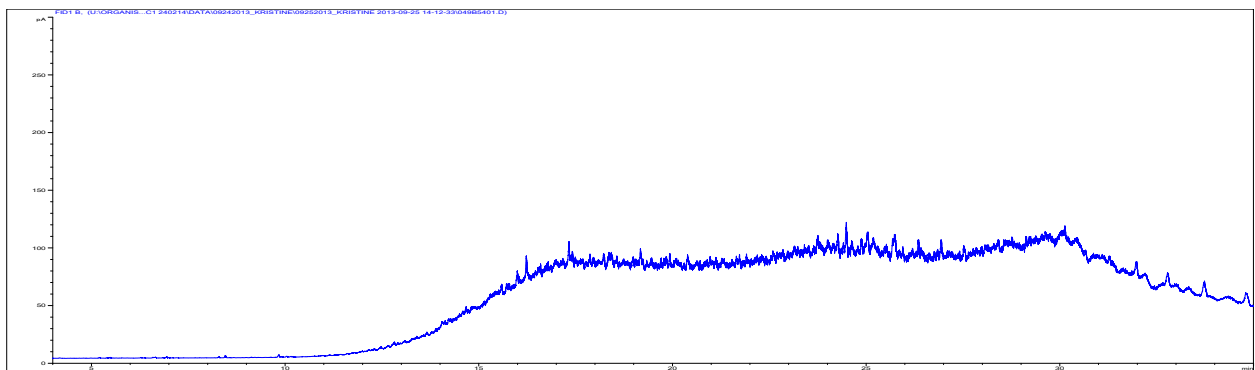
2012-0290



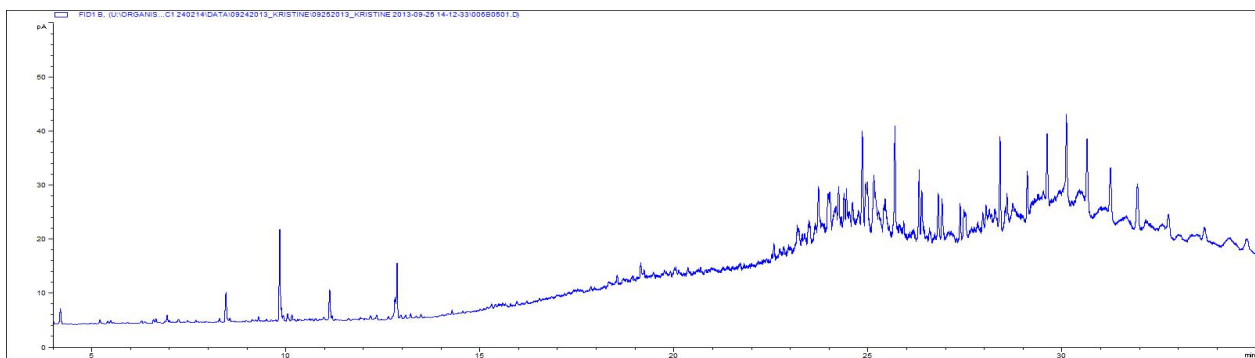
2012-0295



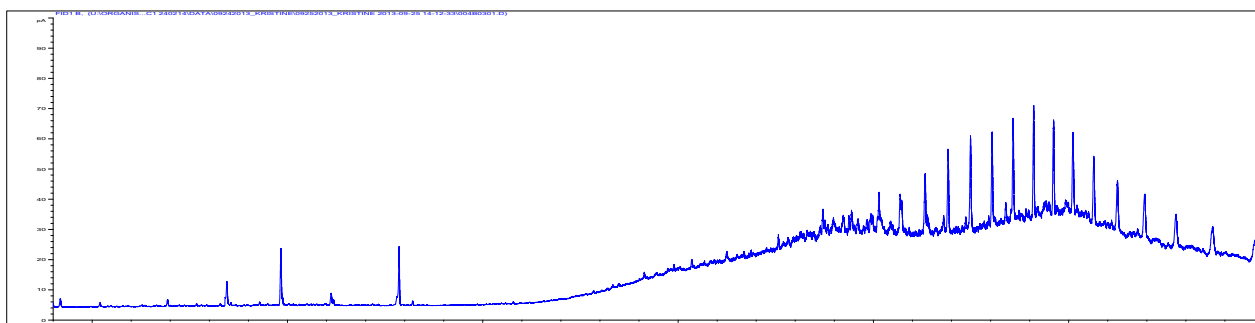
2012-314



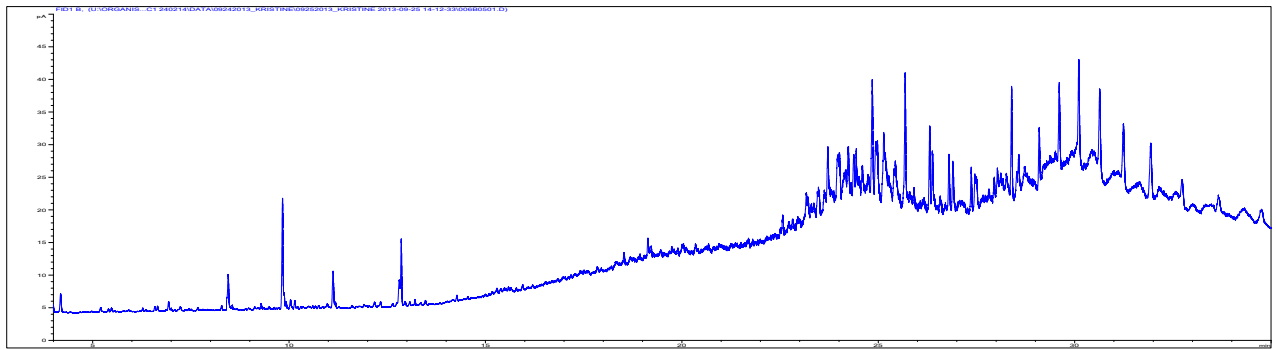
2012-0316



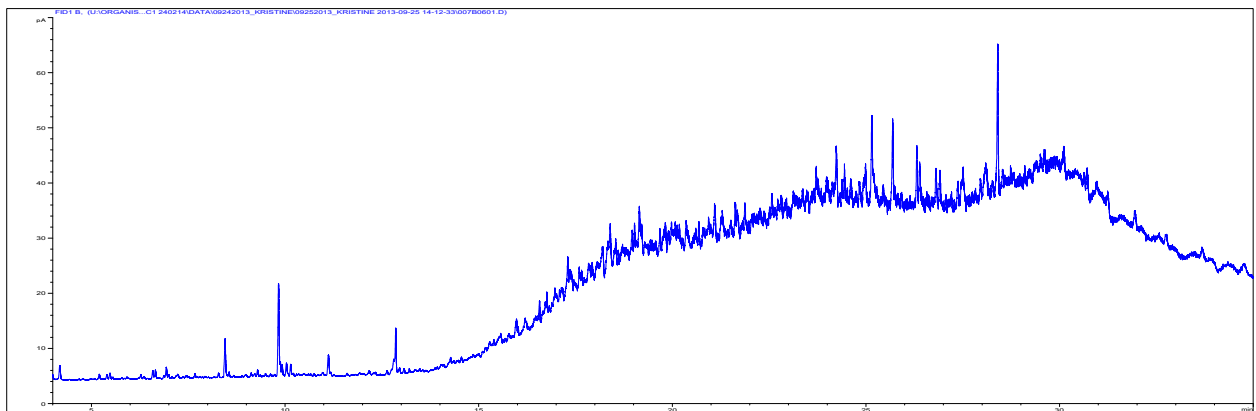
2012-AS02



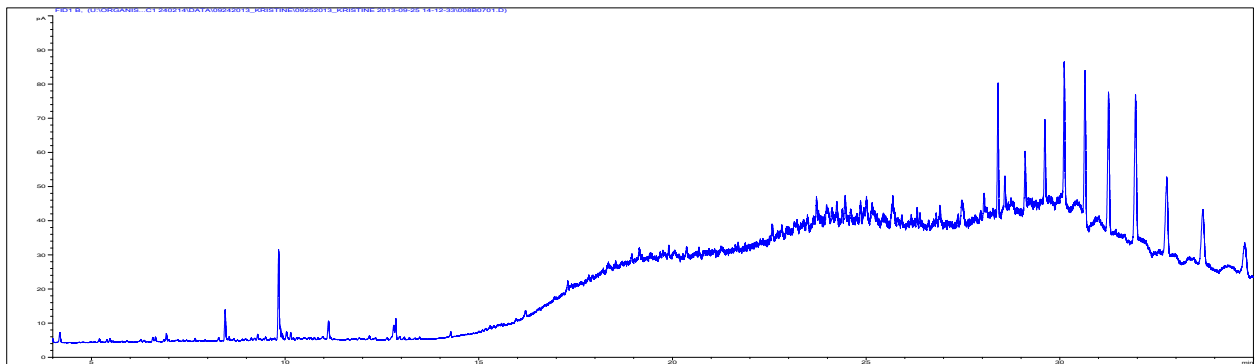
2012-AS06



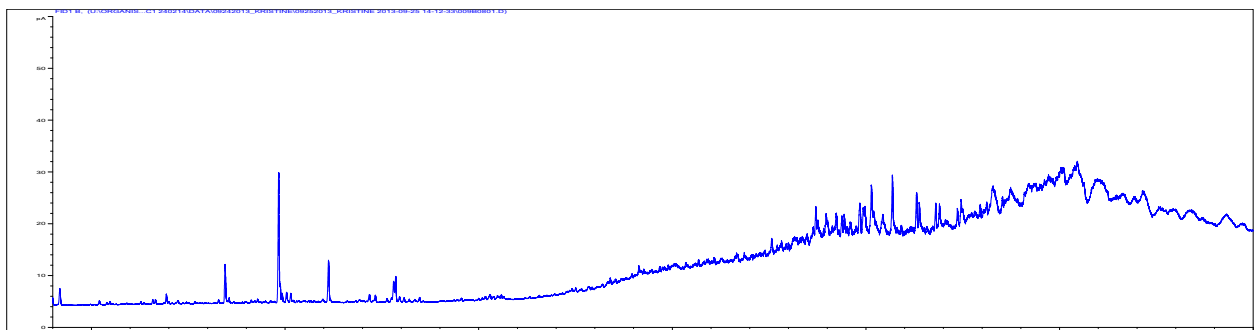
2012-AS08



2012-BS01

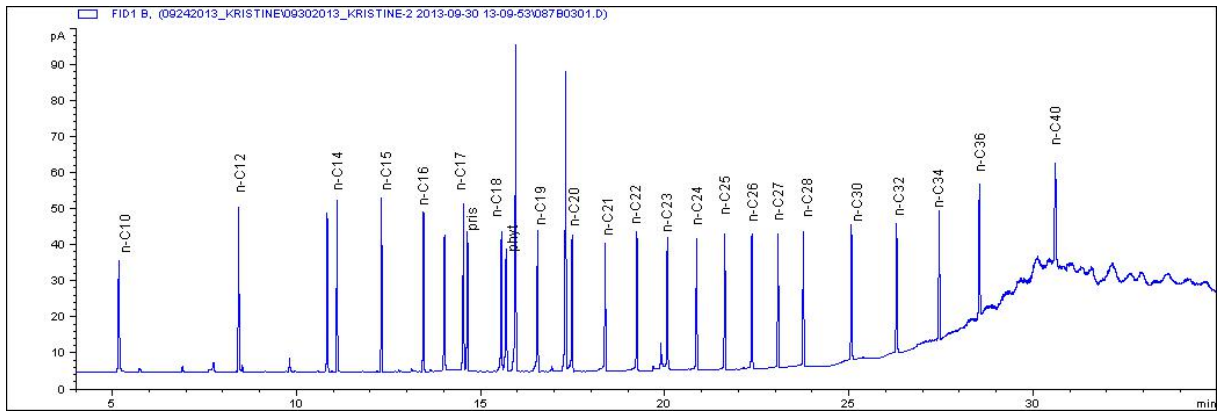


2012-BS03

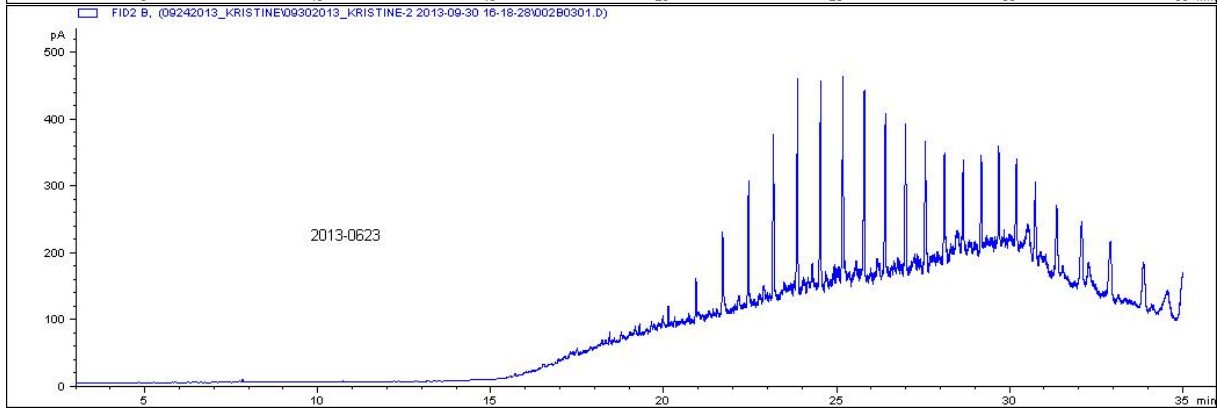
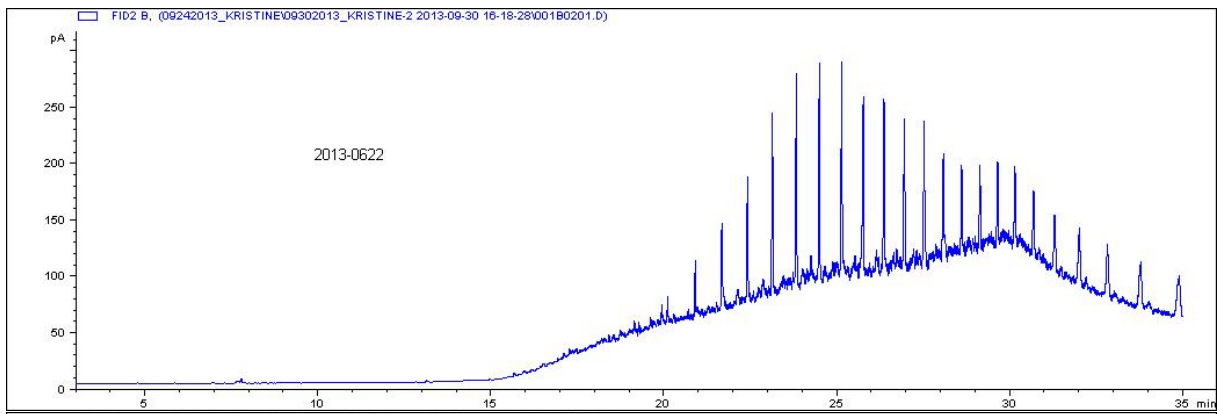


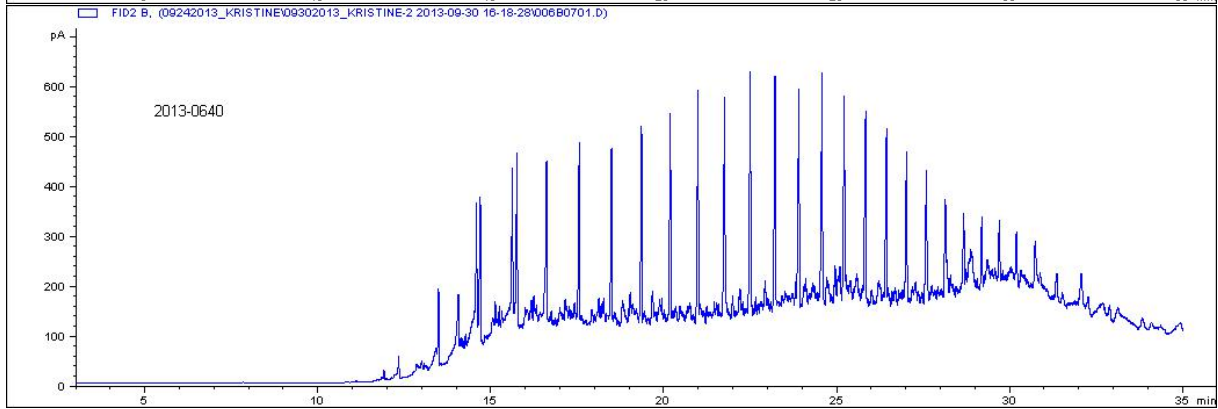
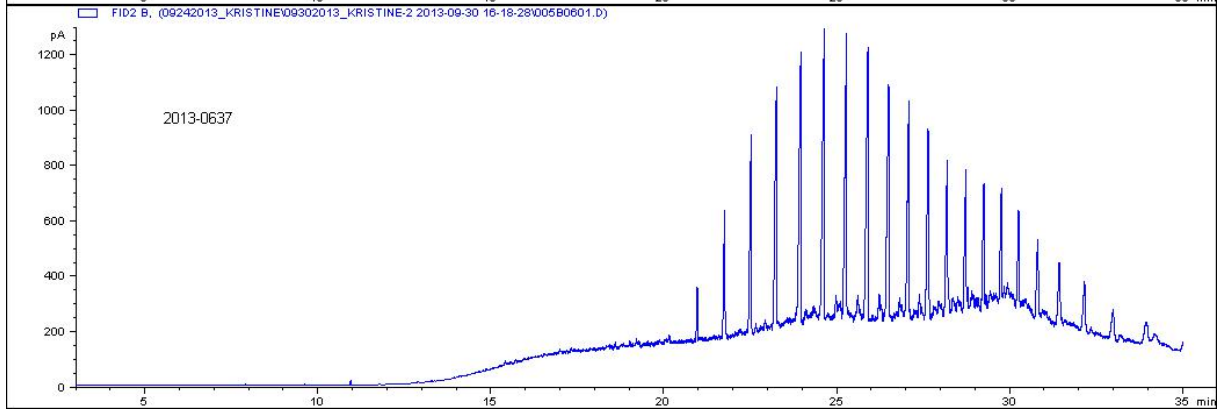
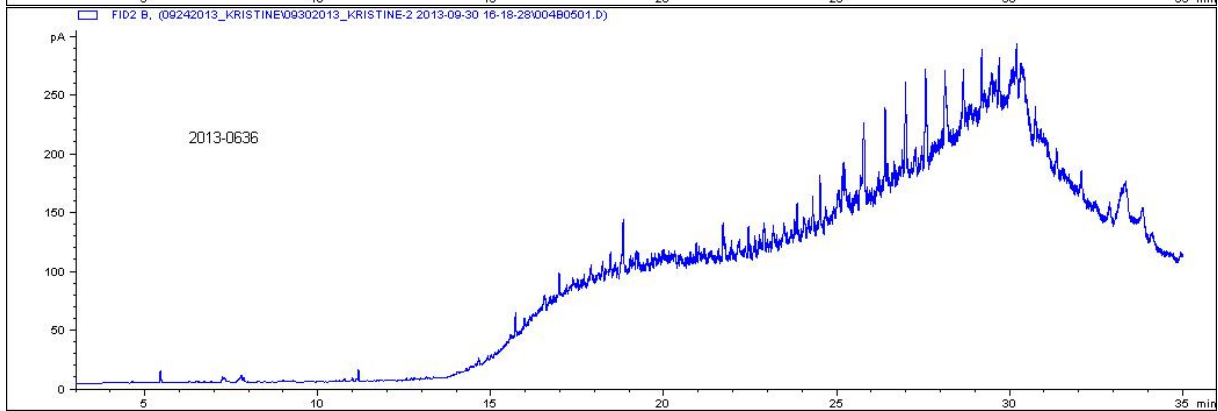
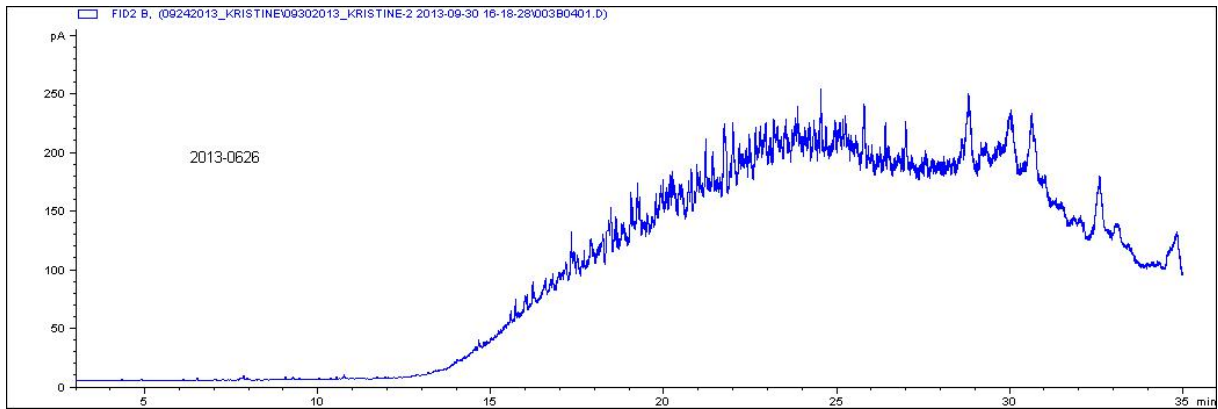
2012-BS05

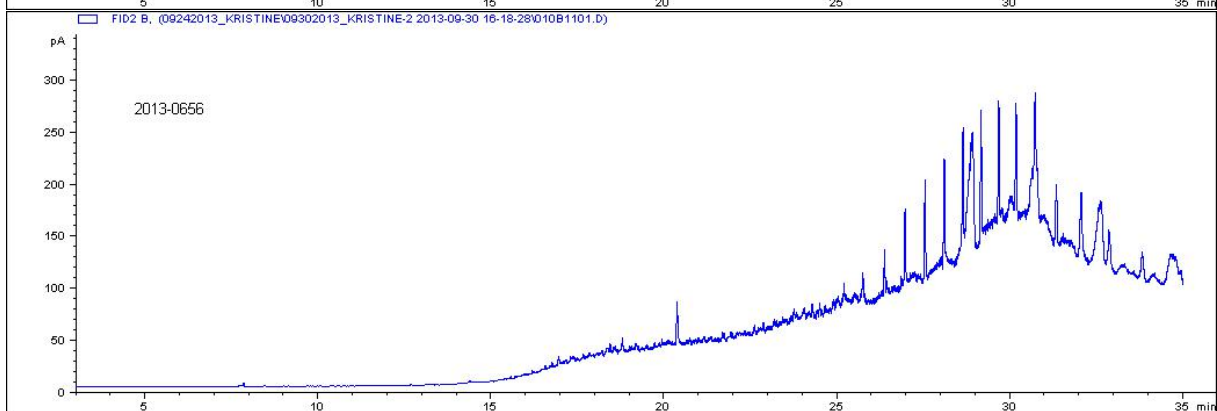
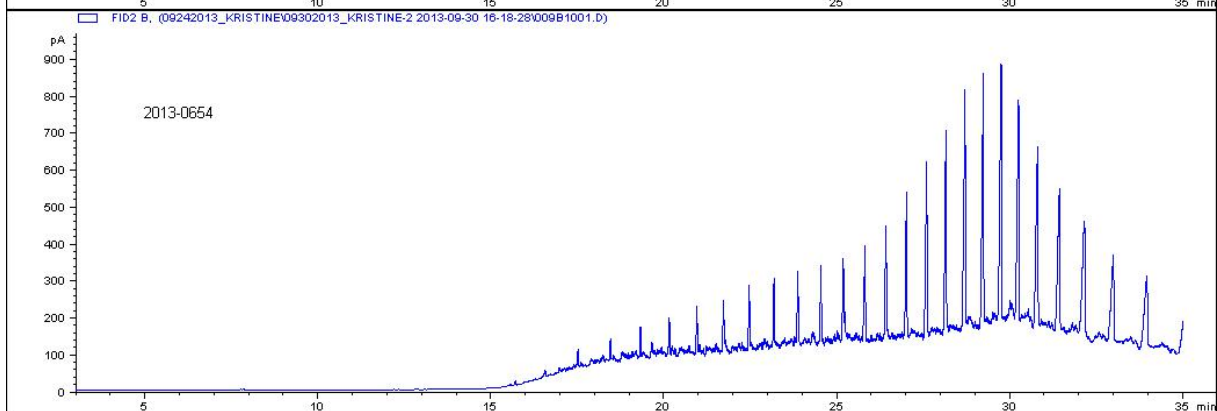
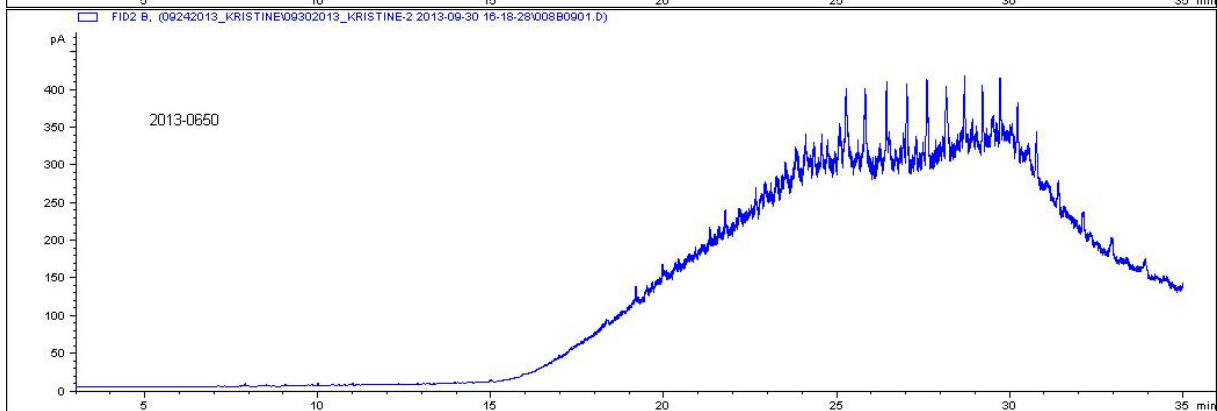
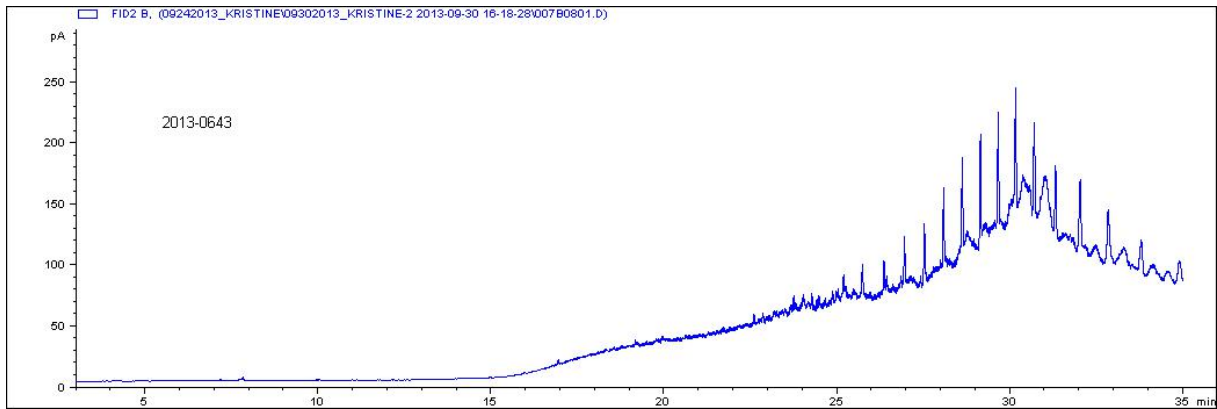
GC-FID 2013

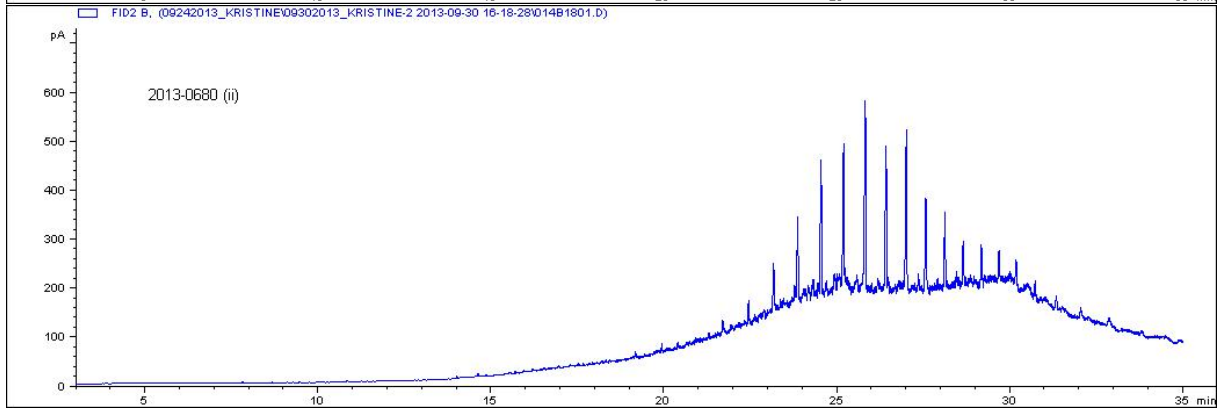
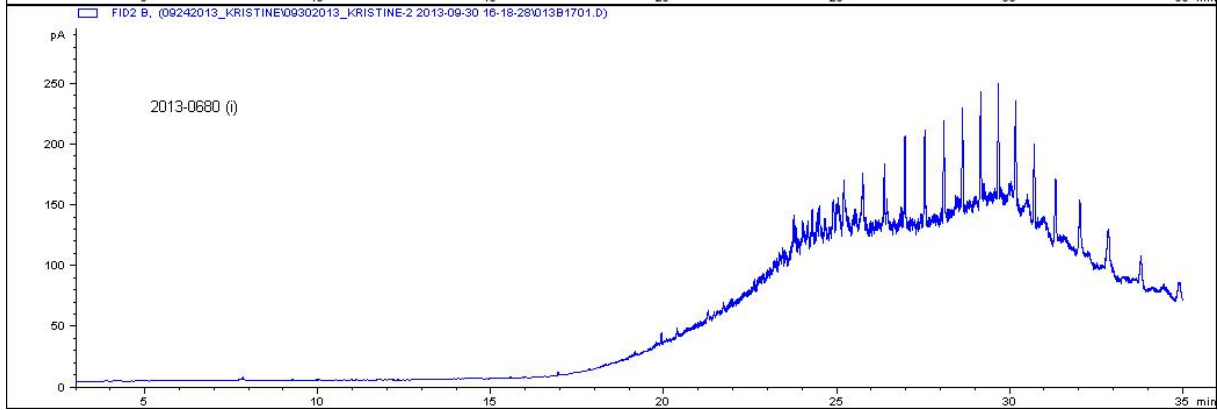
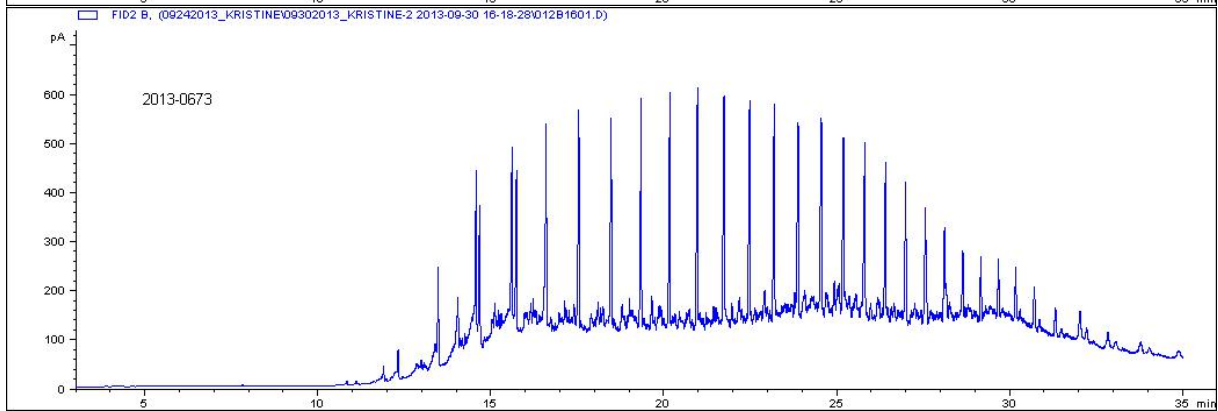
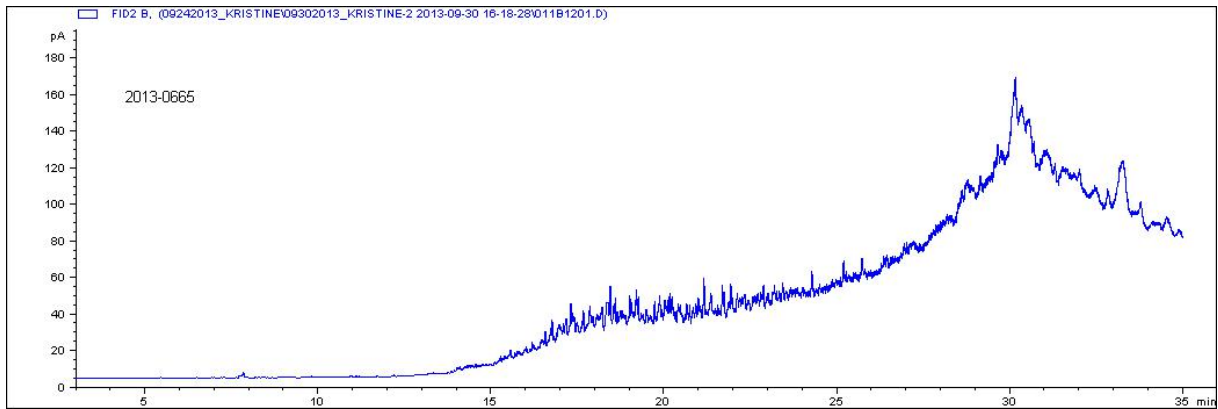


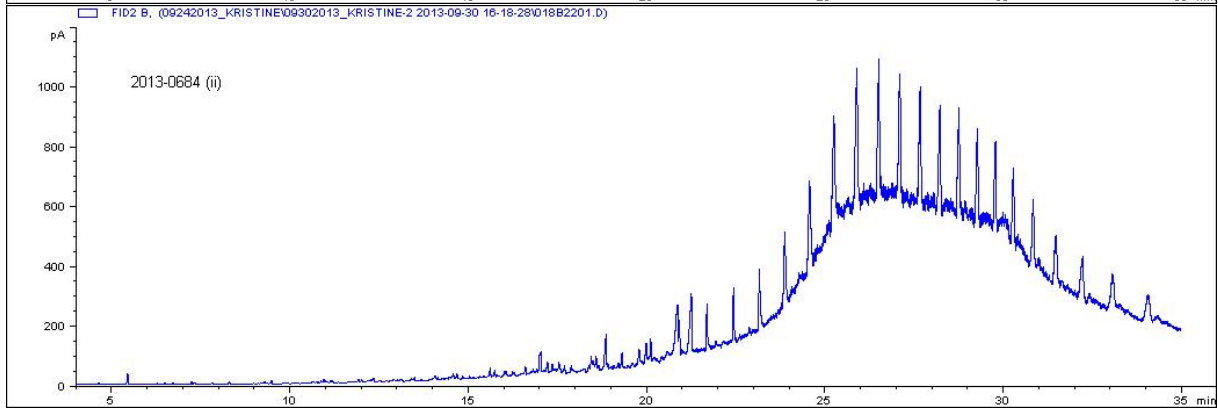
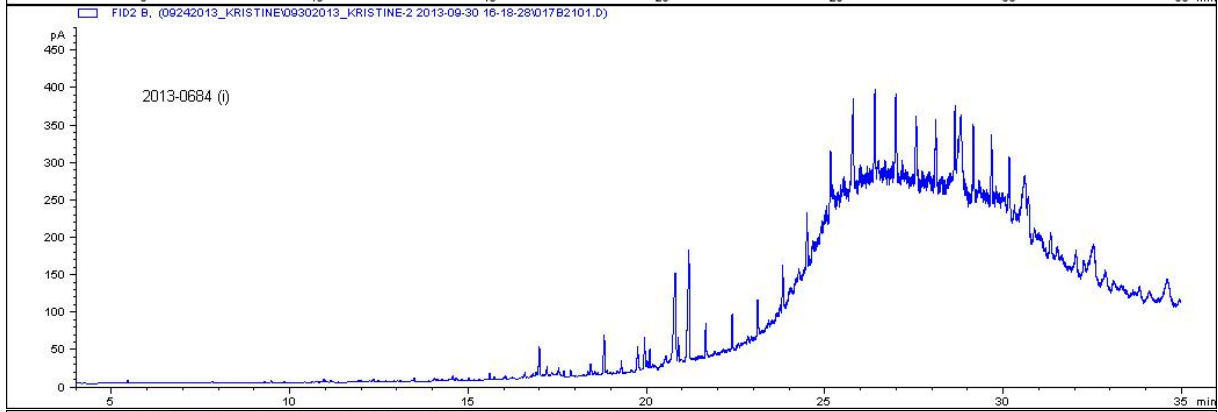
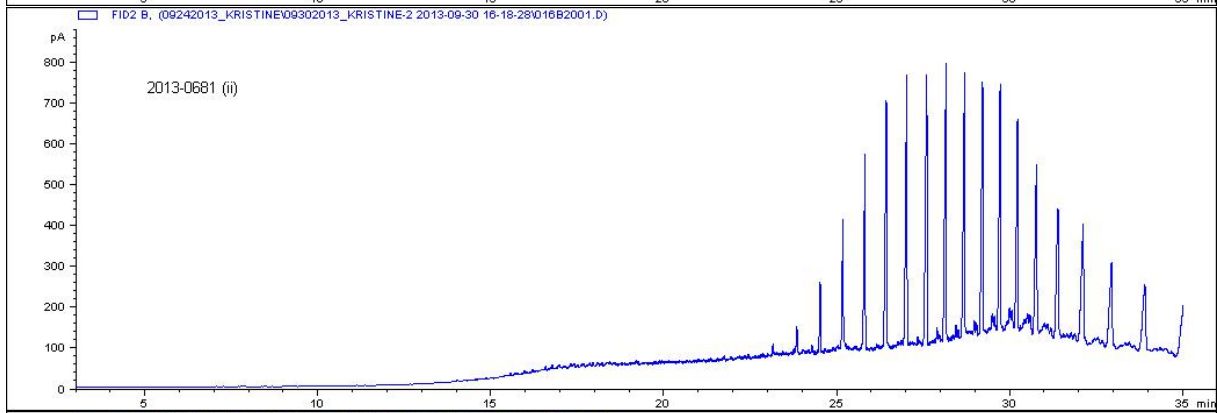
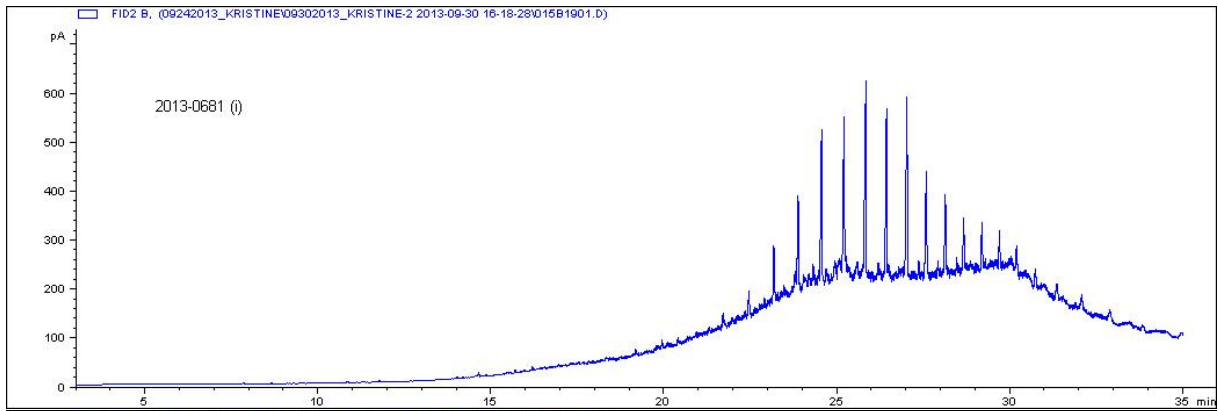
Alkane standard

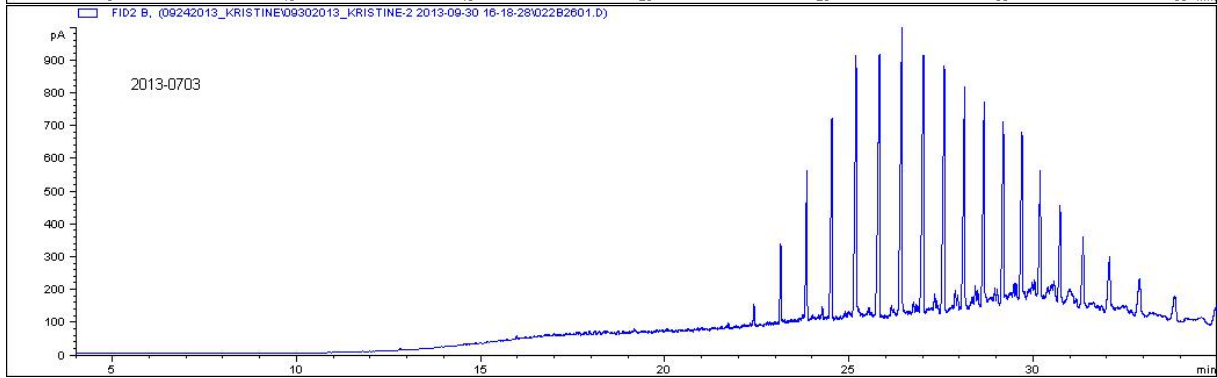
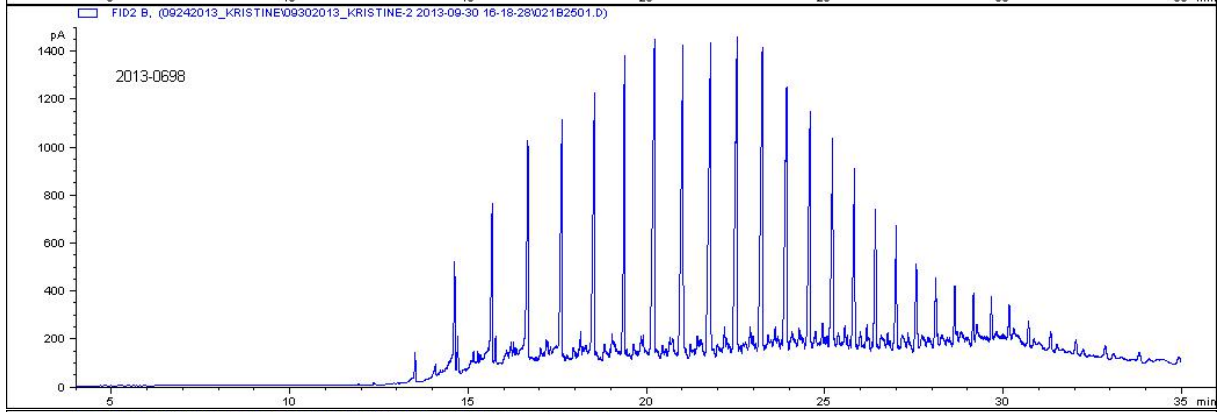
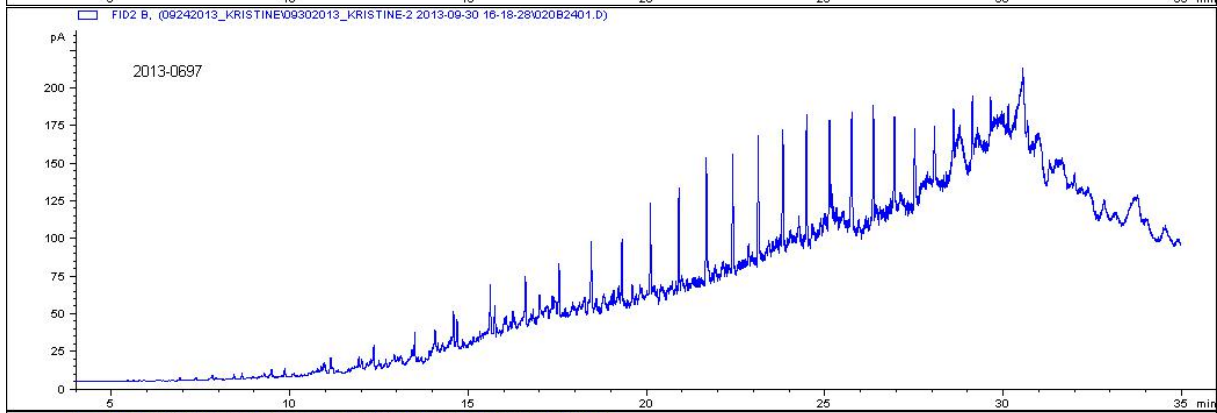
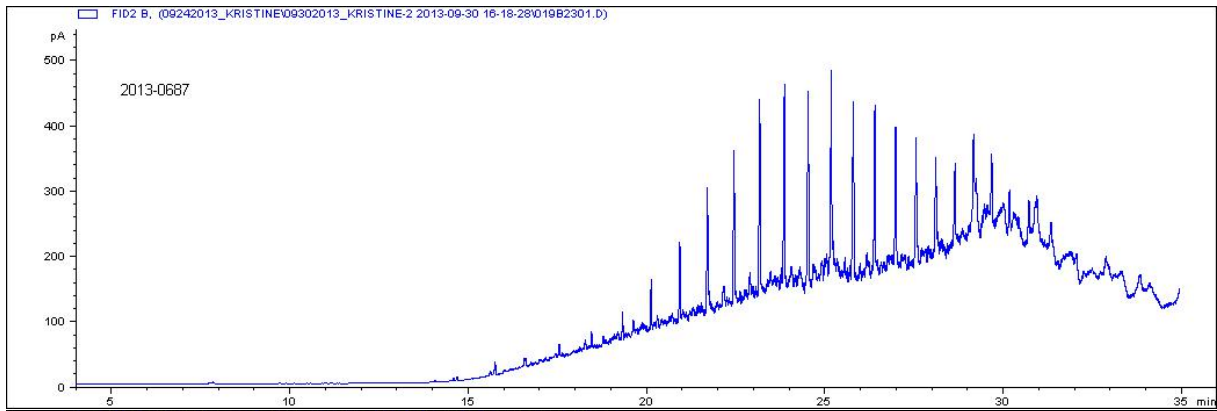




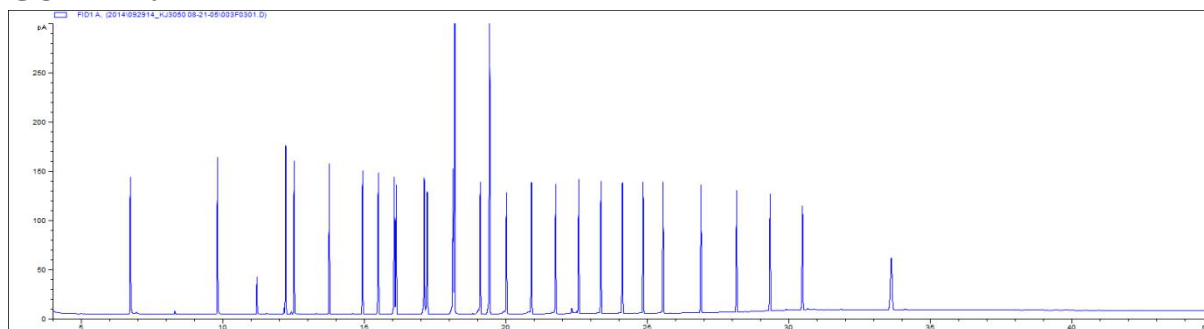




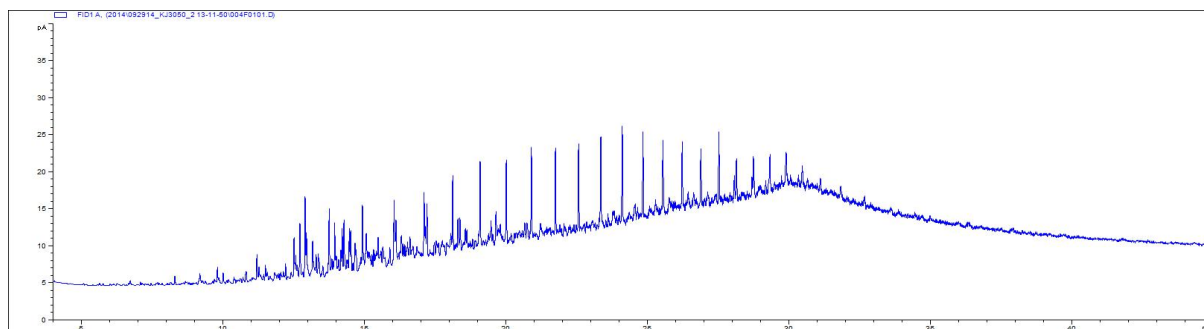




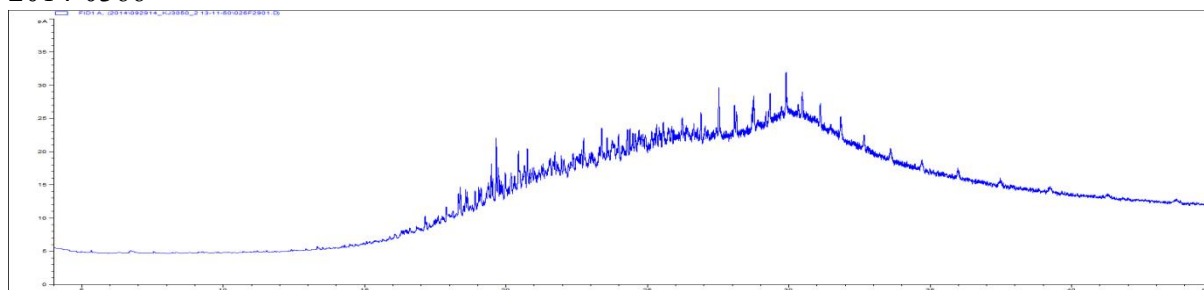
GC-FID 2014



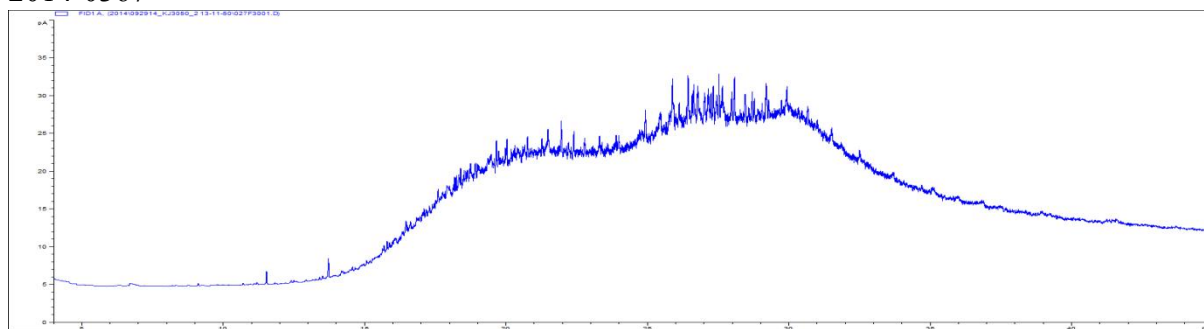
Alkane standard



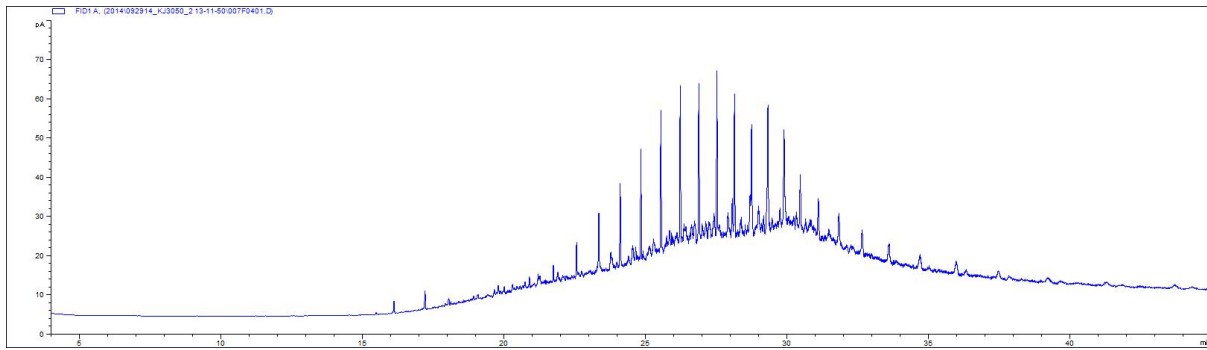
2014-0366



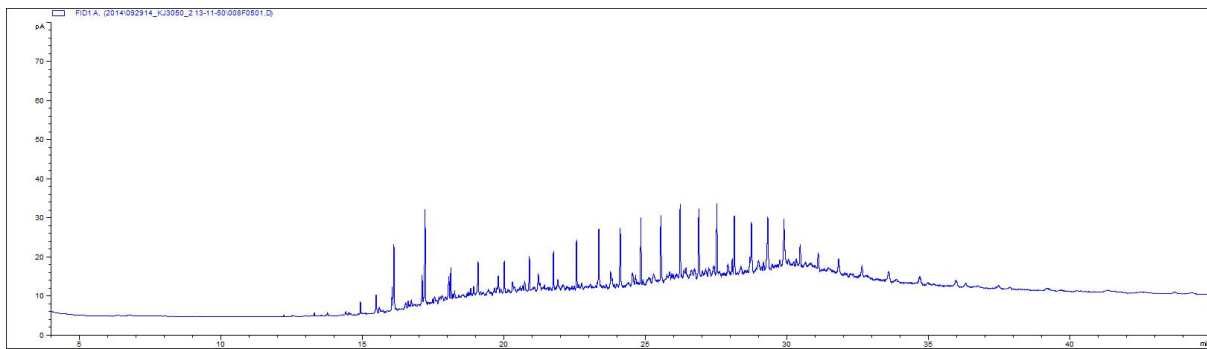
2014-0367



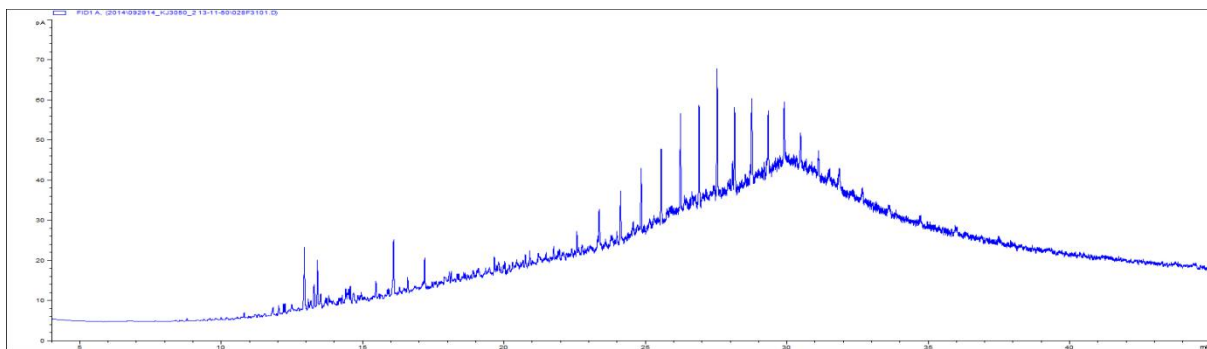
2014-0368



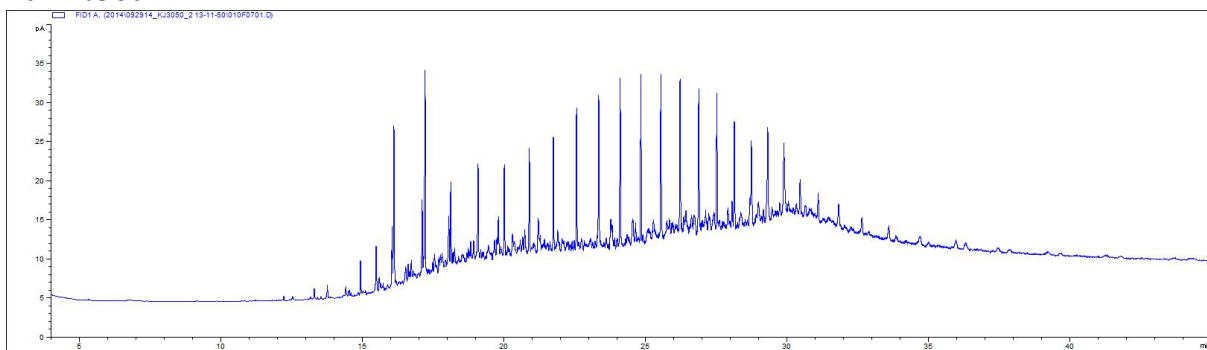
2014-0371



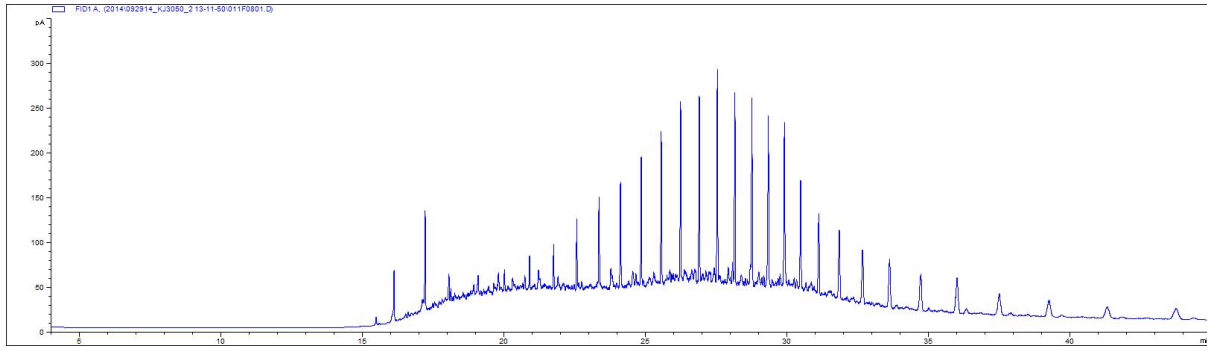
2014-0384



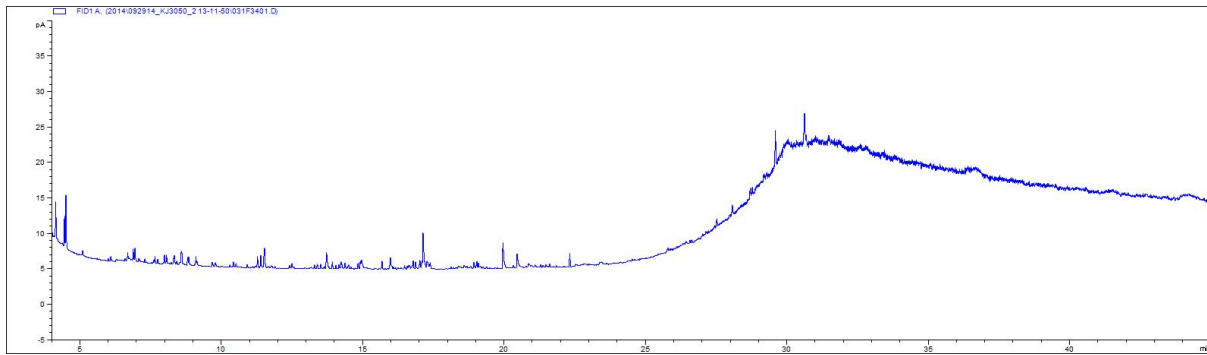
2014-0386



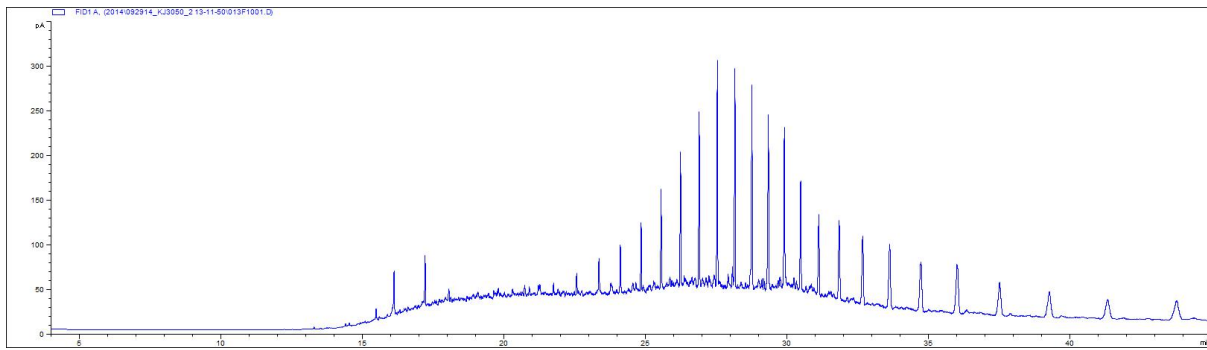
2014-0387



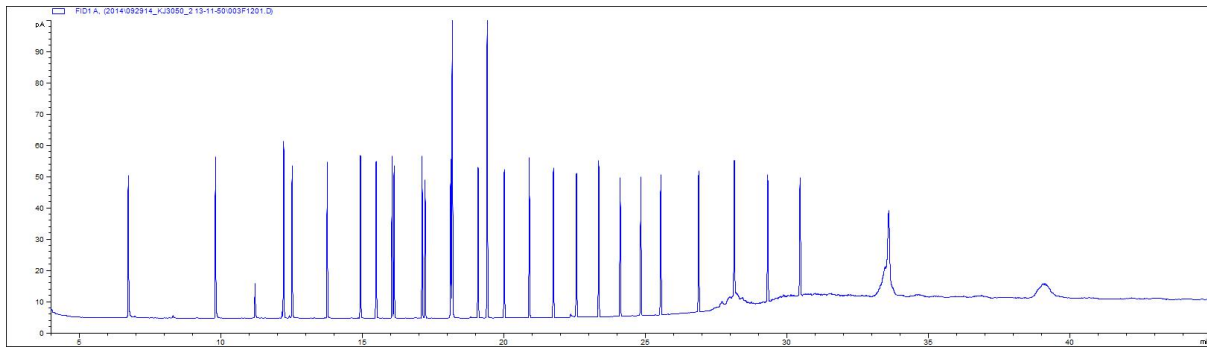
2014-0390



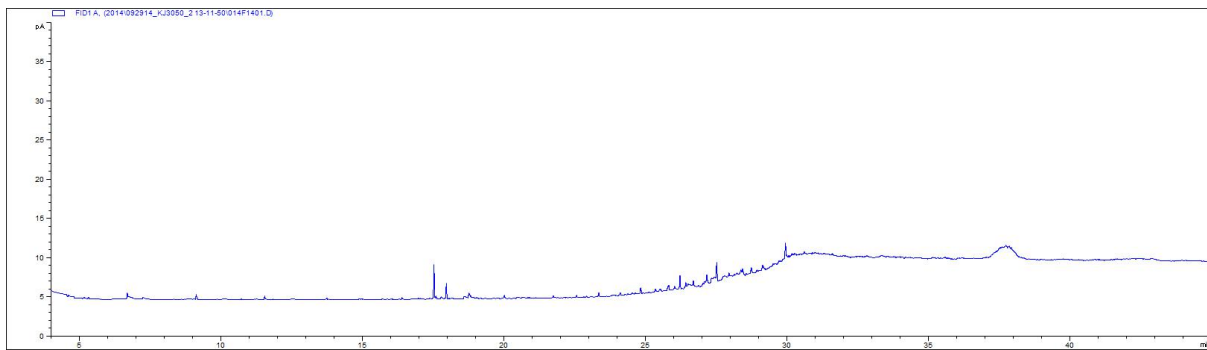
2014-0394



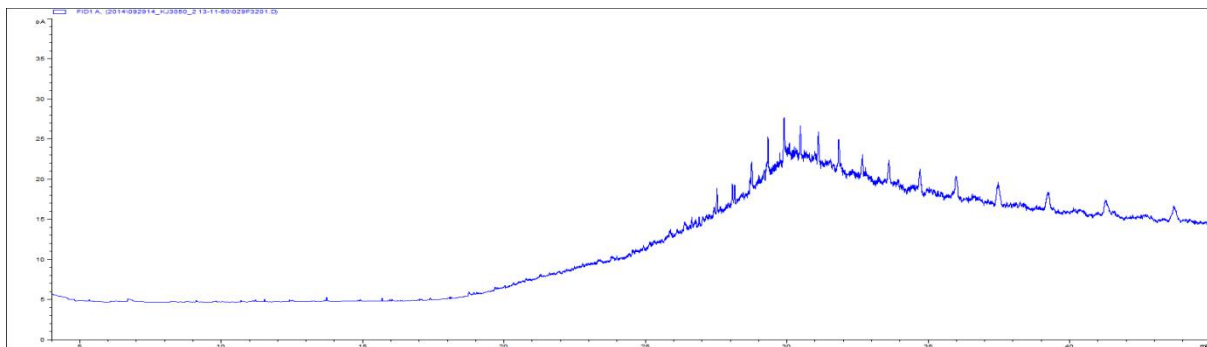
2014-0395



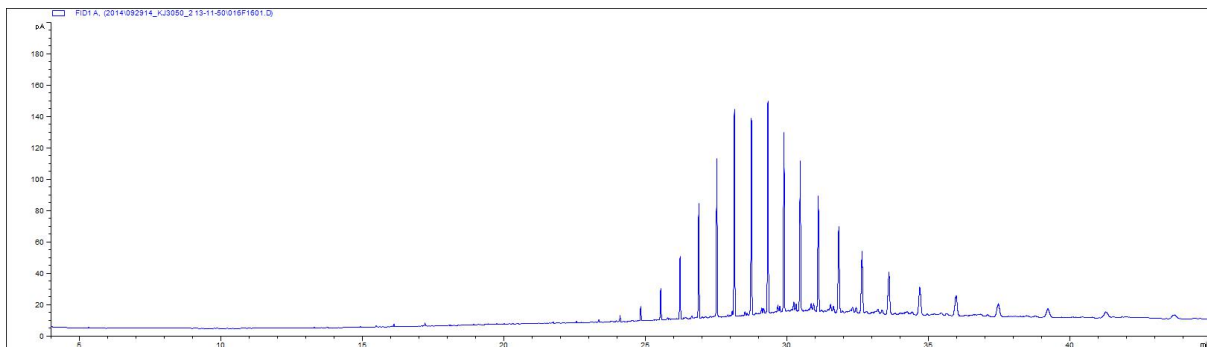
n-alkan std



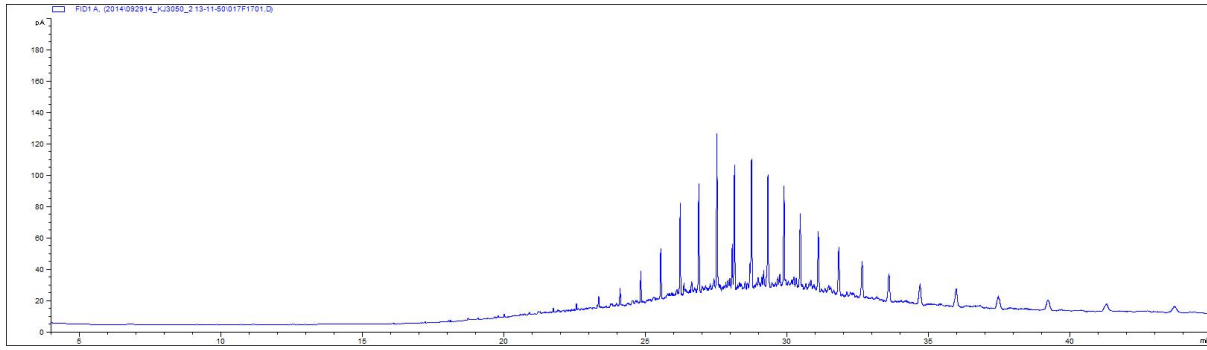
2014-0397



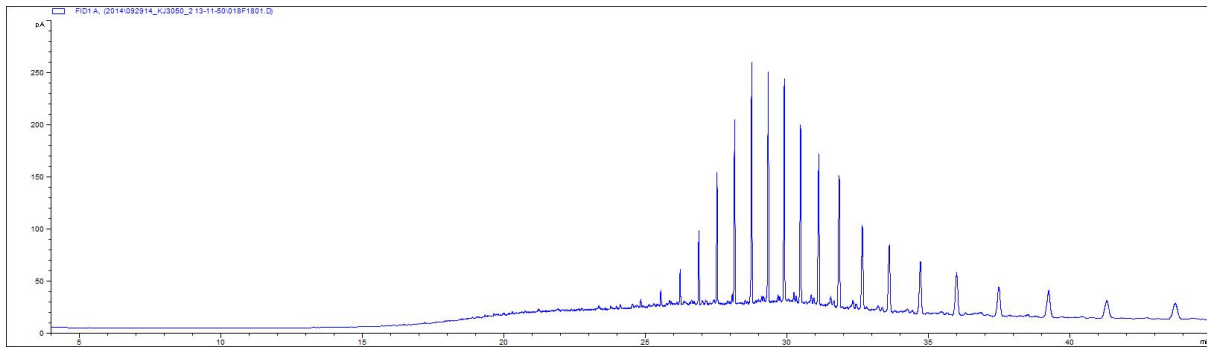
2014-0401



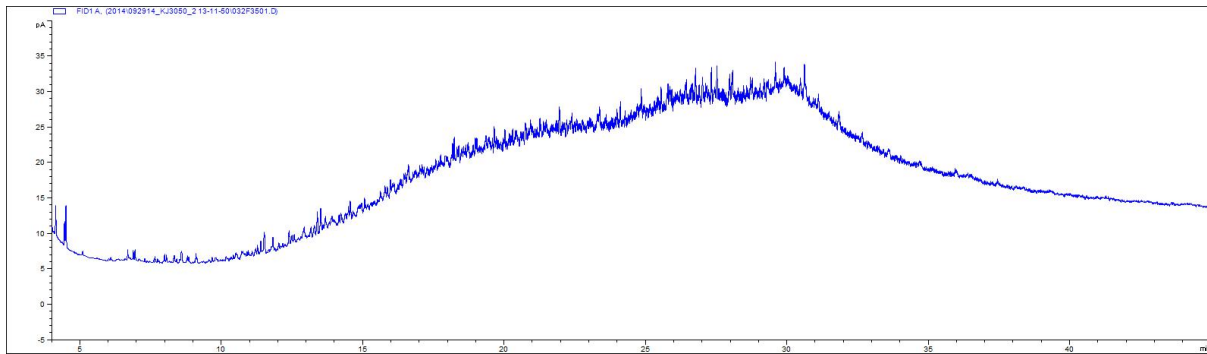
2014-0401B



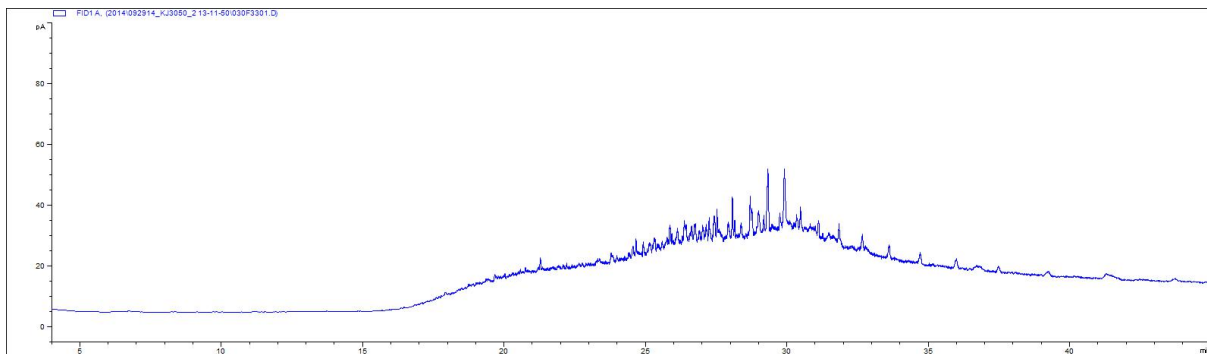
2014-0408



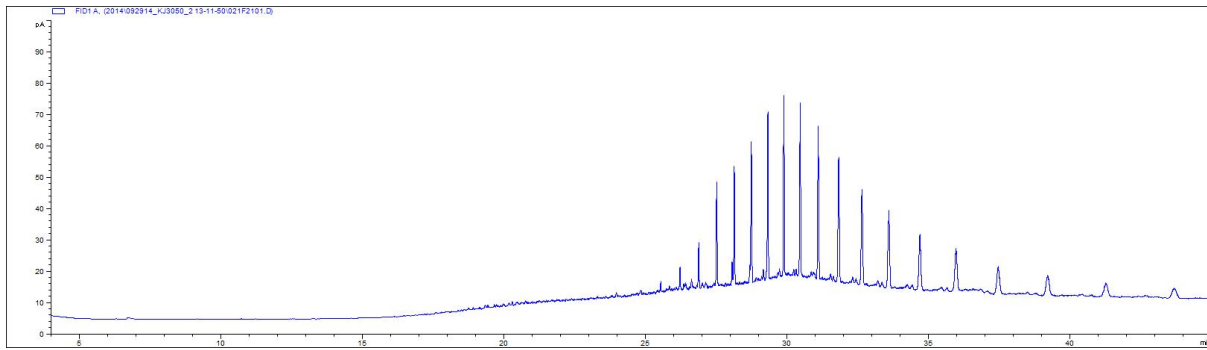
2014-0413



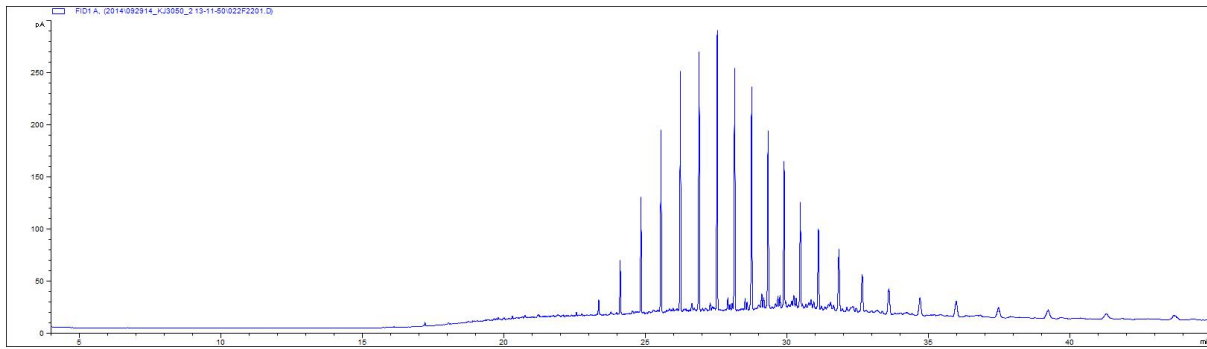
2014-0414



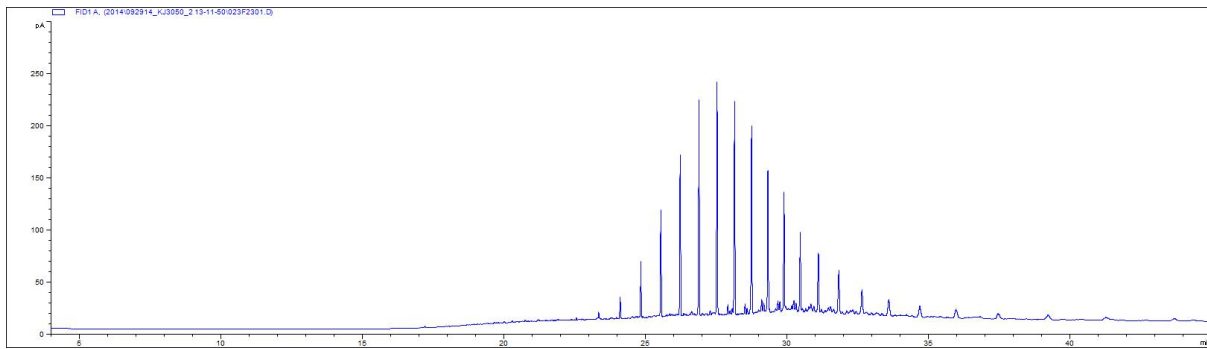
2014-0420



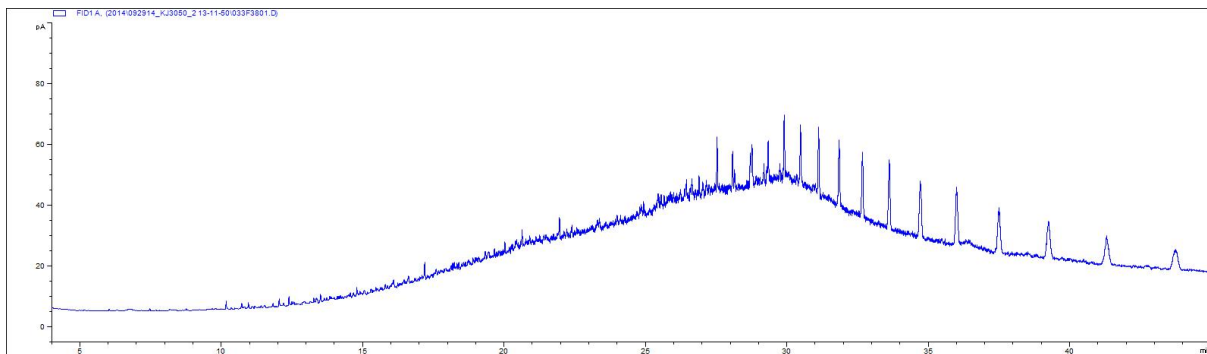
2014-0421



2014-0424

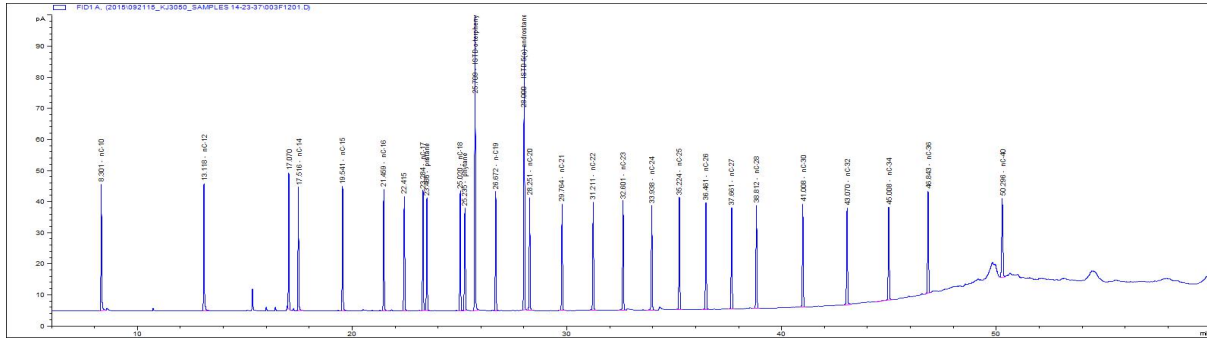


2014-0424B

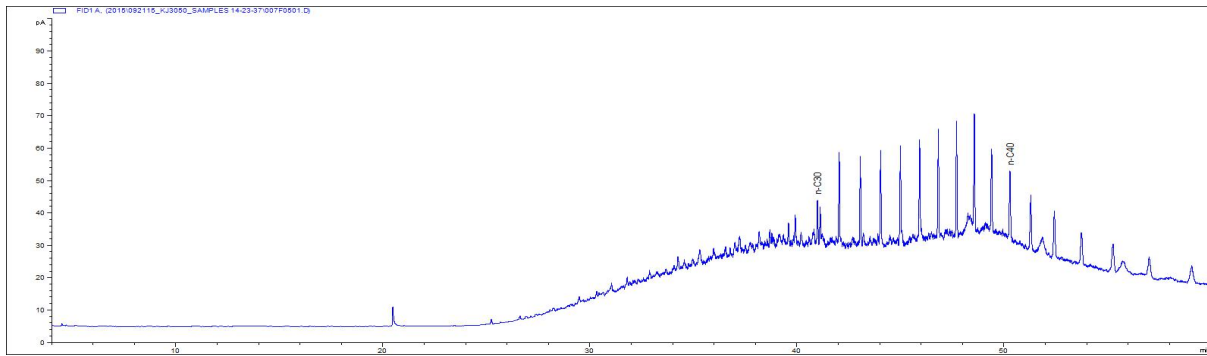


2014-0425

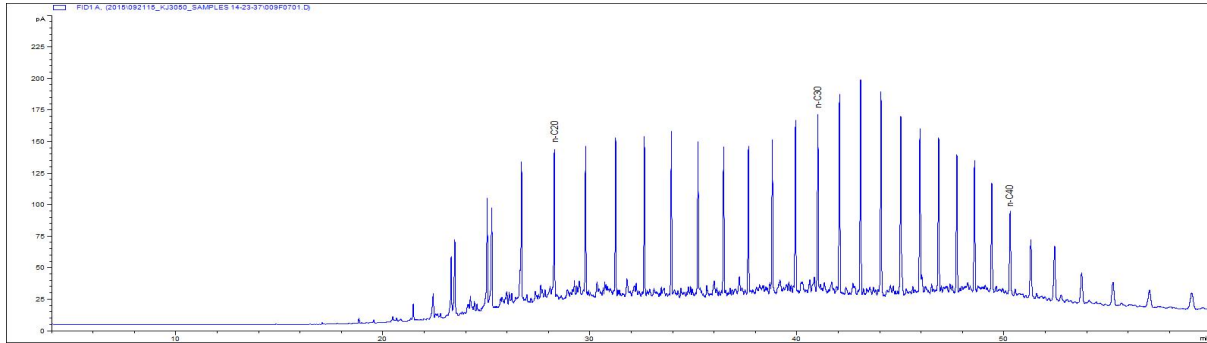
GC-FID 2015



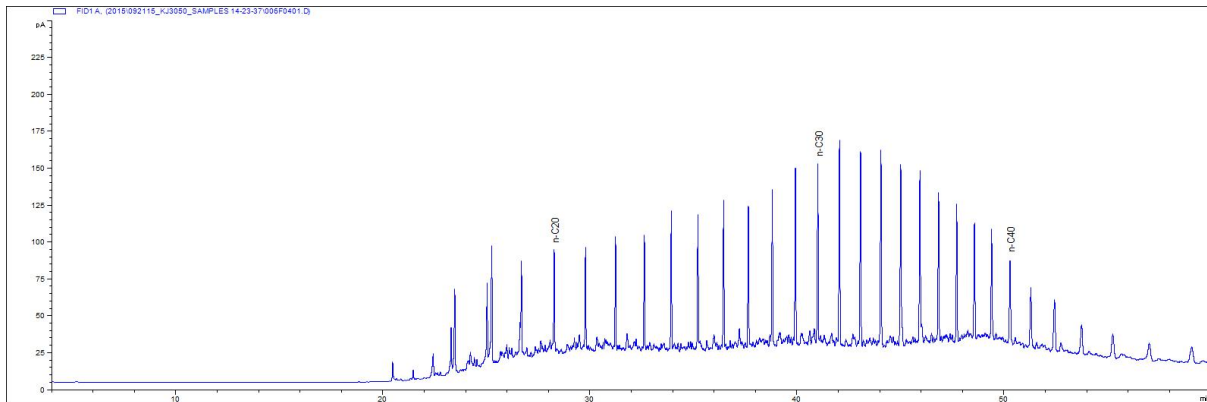
n-alkane standard



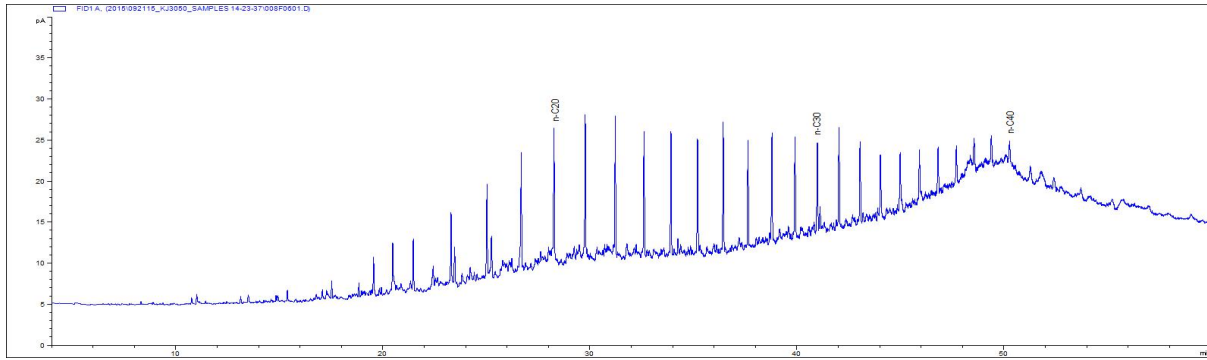
2015-837



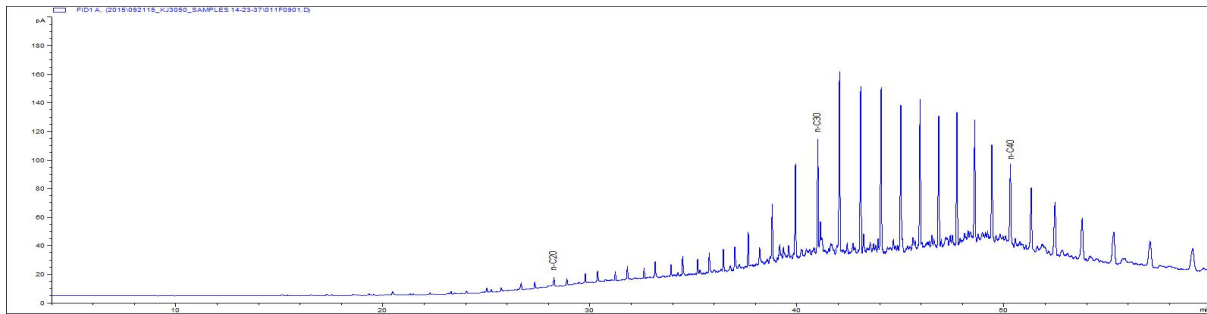
2015-839



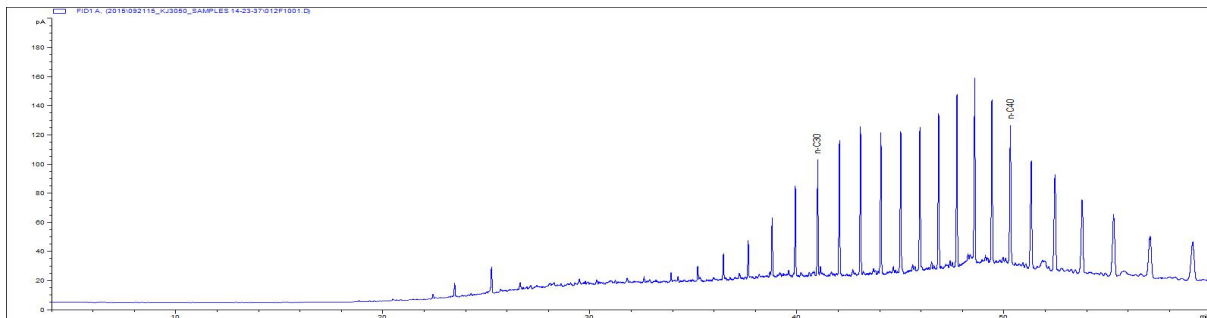
2015-844



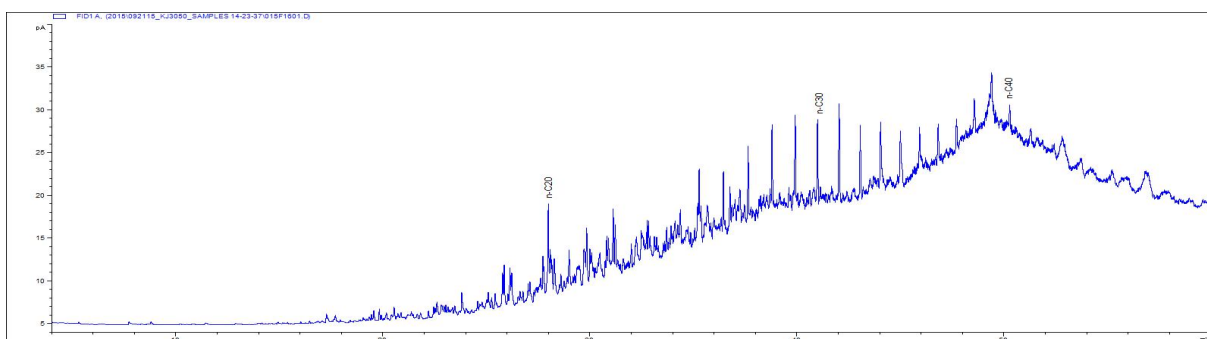
2015-846



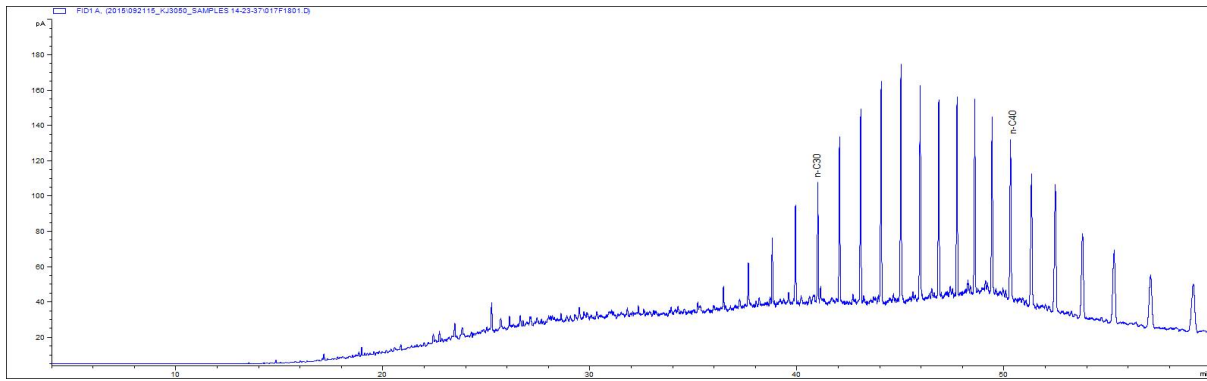
2015-847



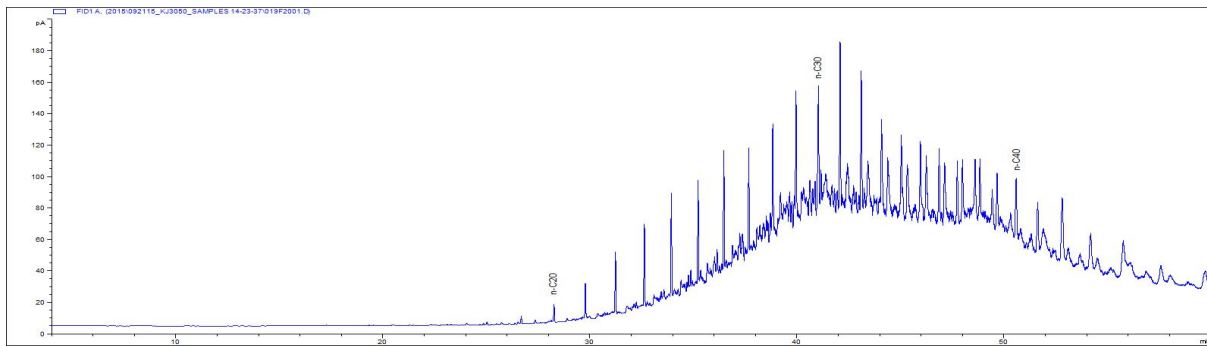
2015-848



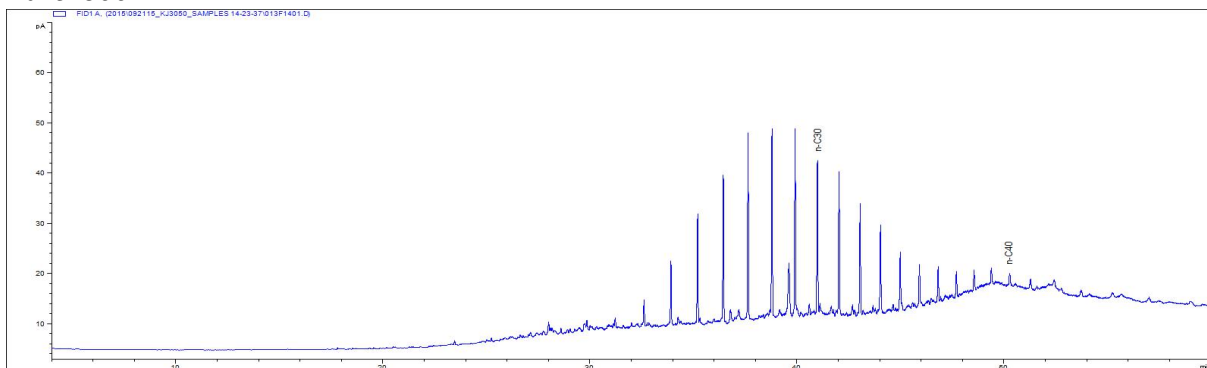
2015-857



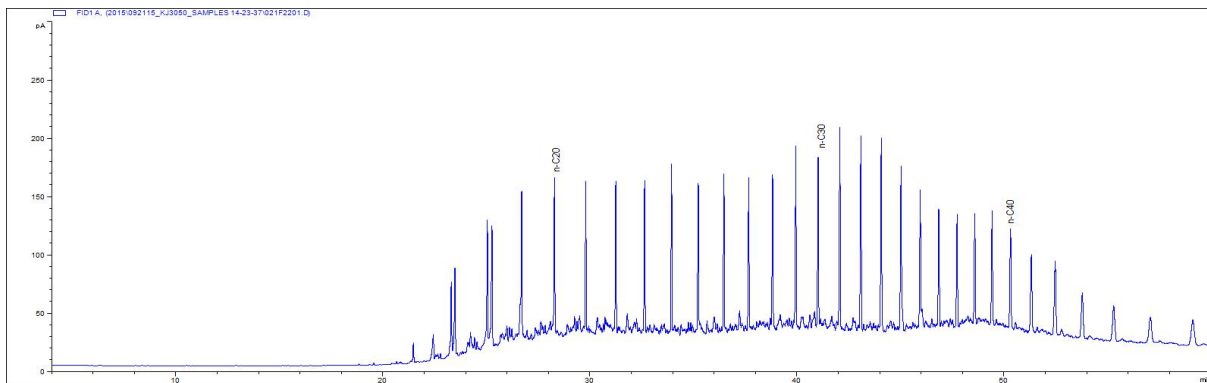
2015-864



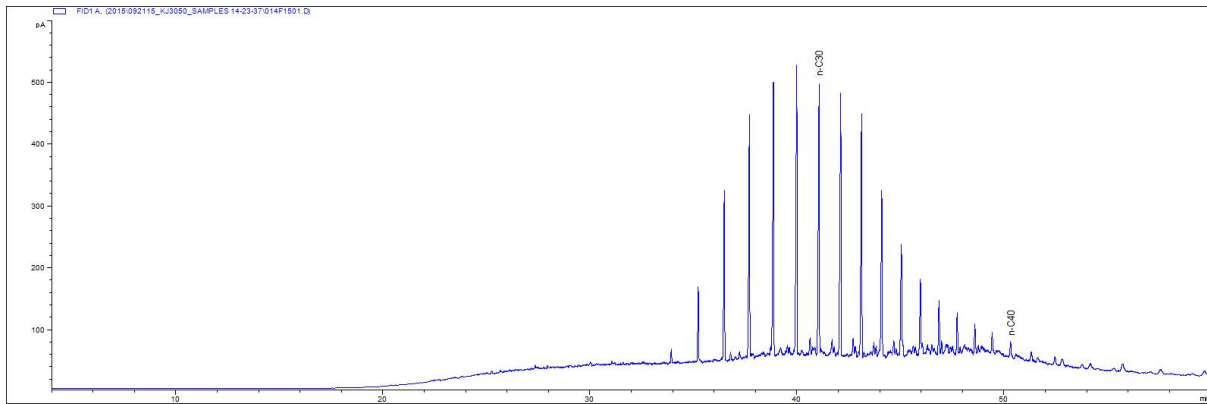
2015-866



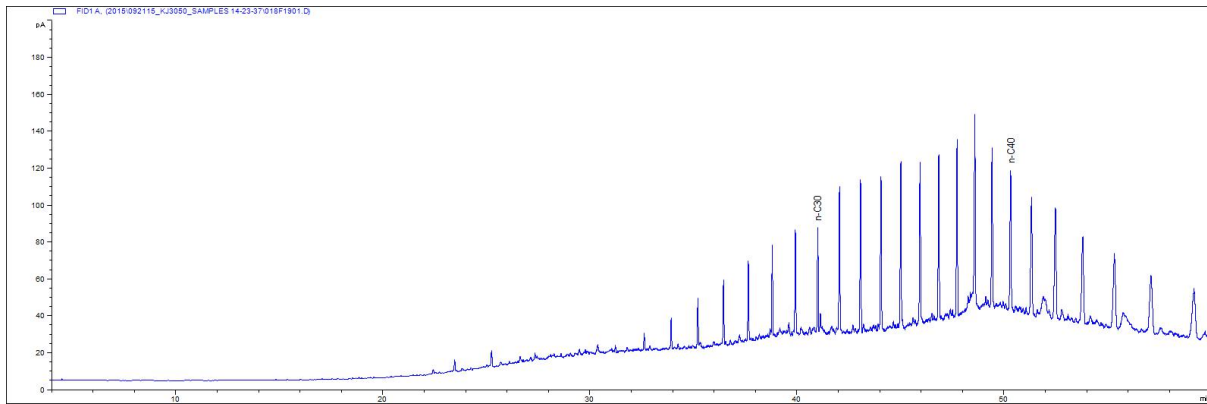
2015-874



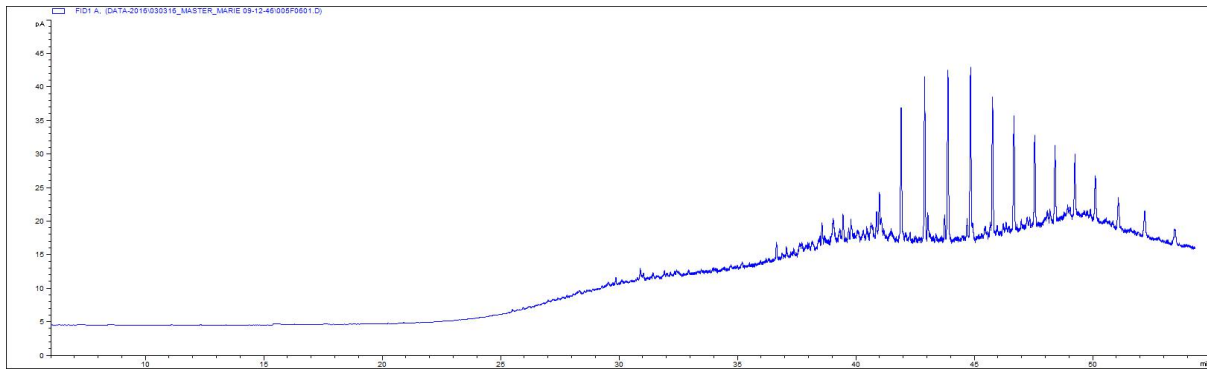
2015-848



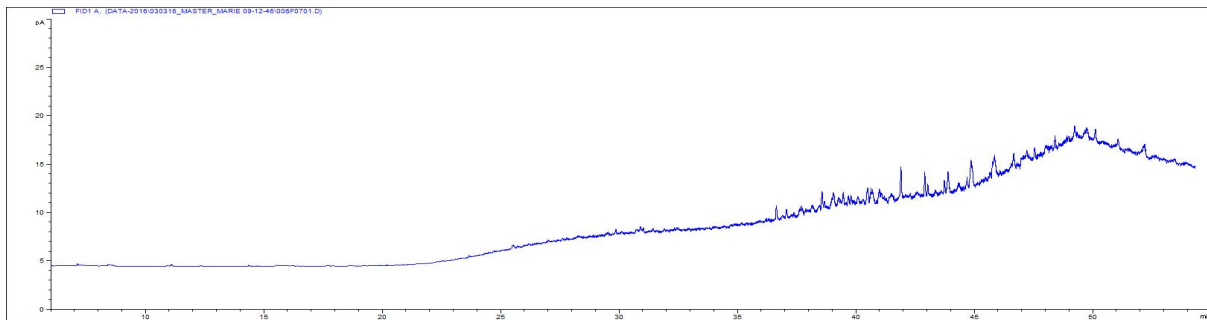
2015-881



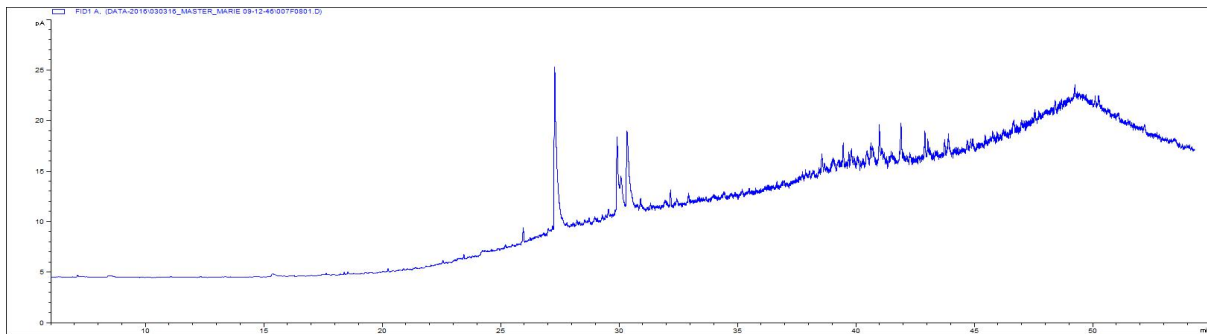
2015-882



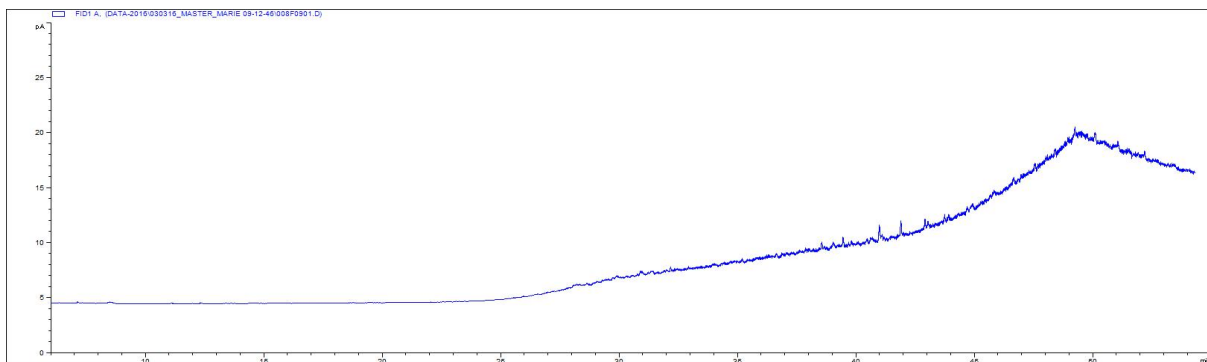
2016-0152



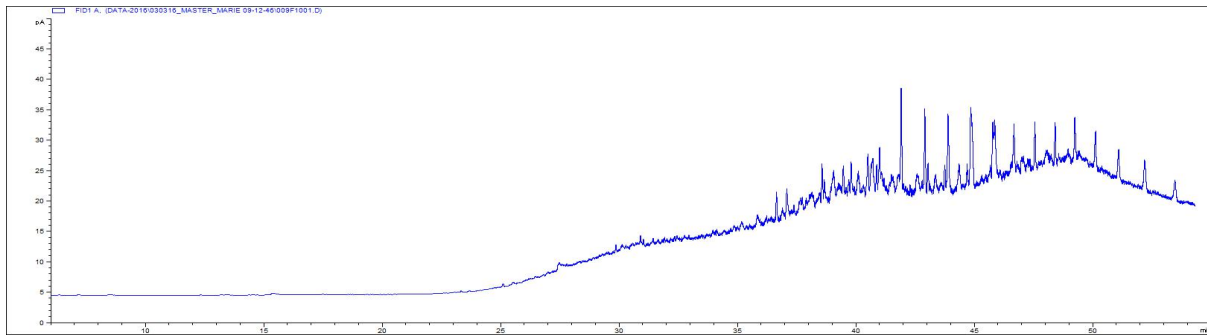
2016-0153



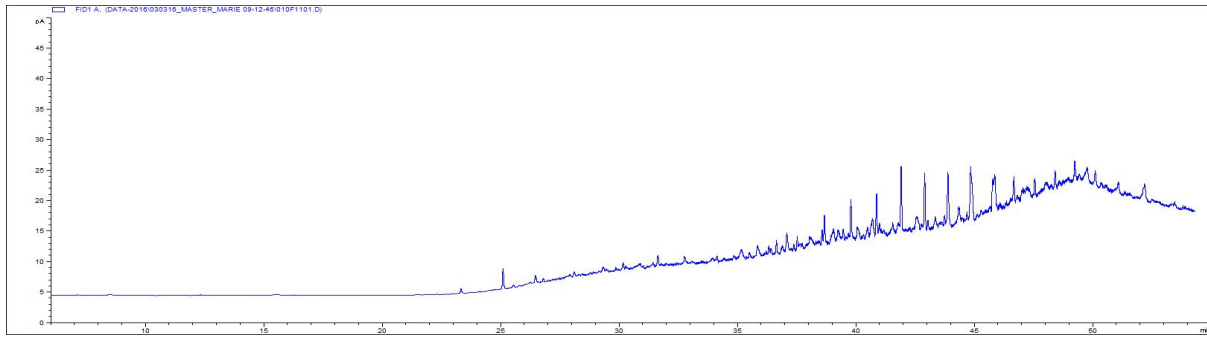
2016-0154



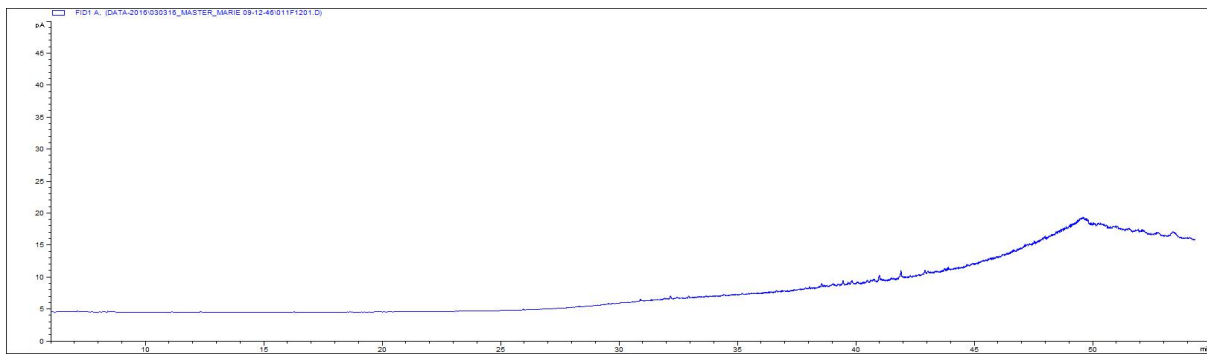
2016-0155



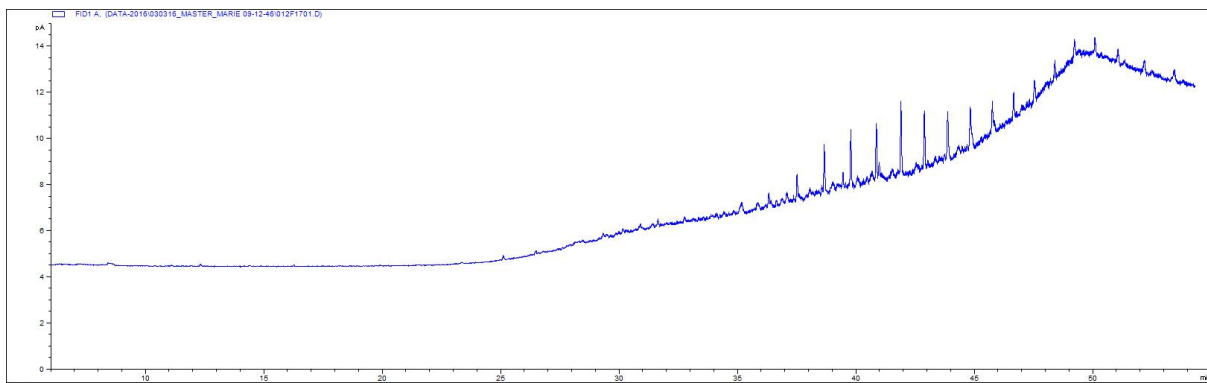
2016-0156



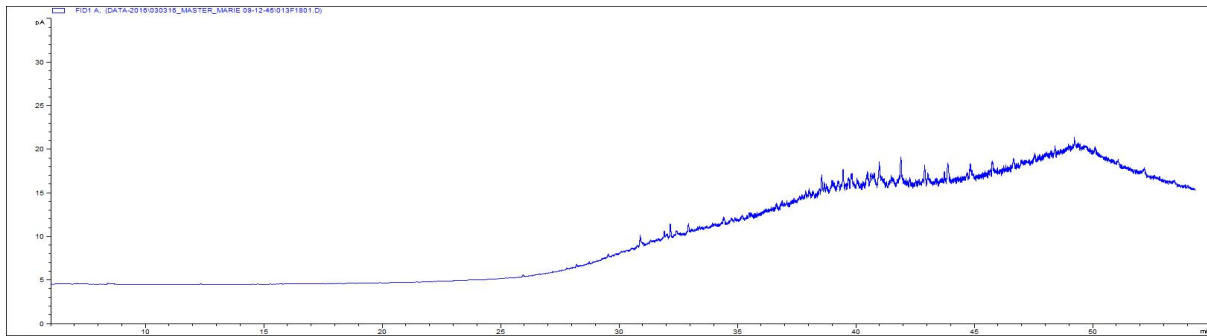
2016-0157



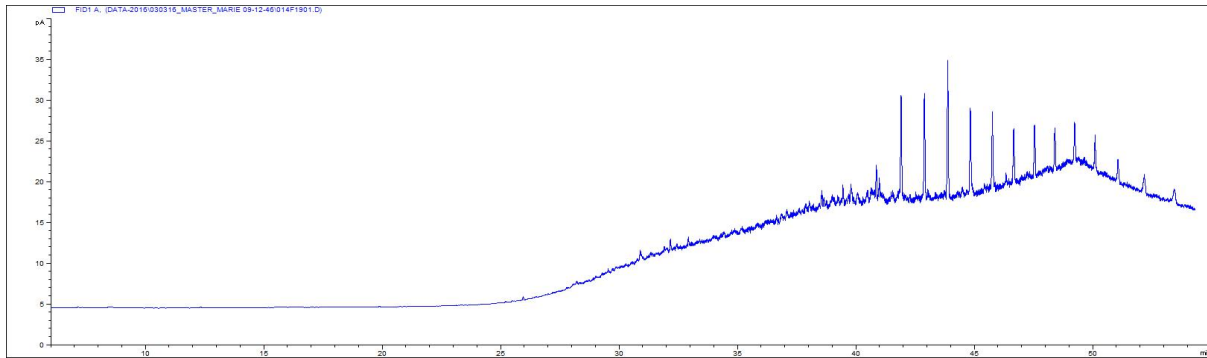
2016-0158



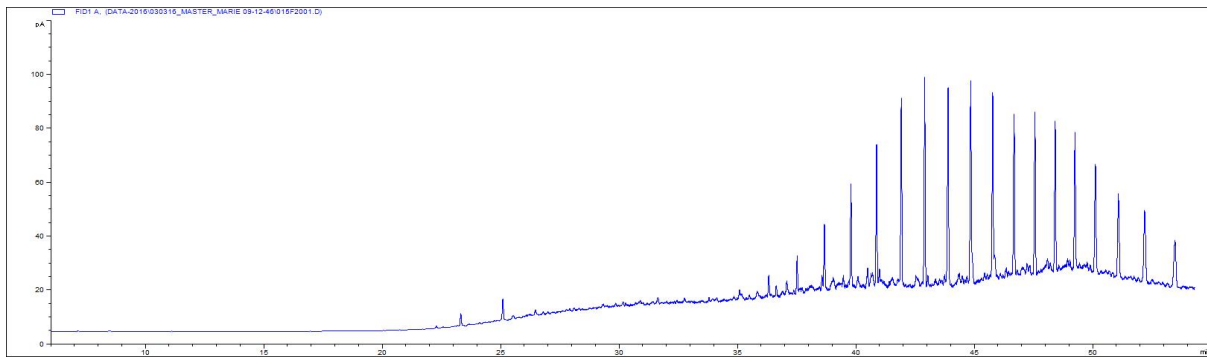
2016-0159



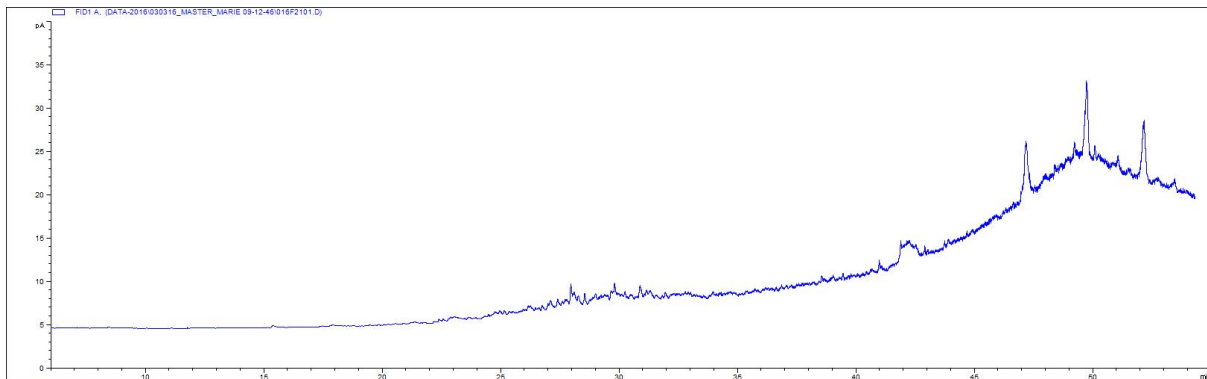
2016-0160



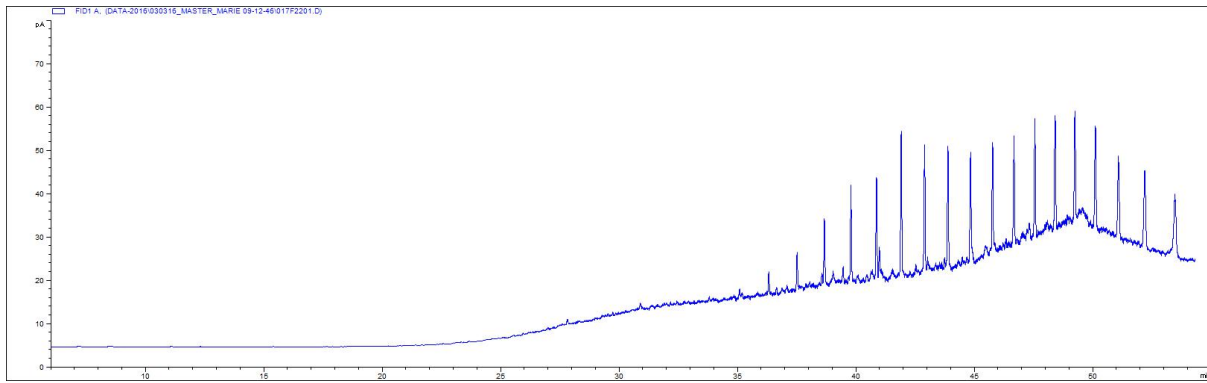
2016-0161



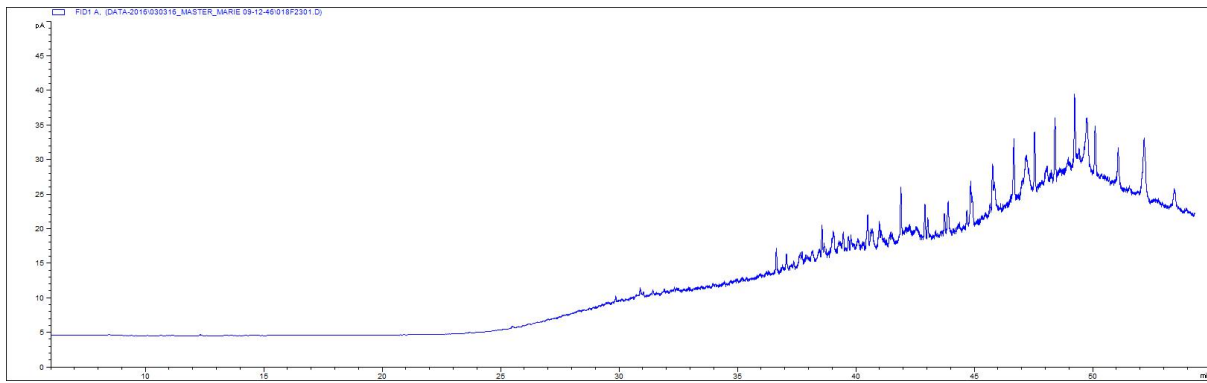
2016-0162



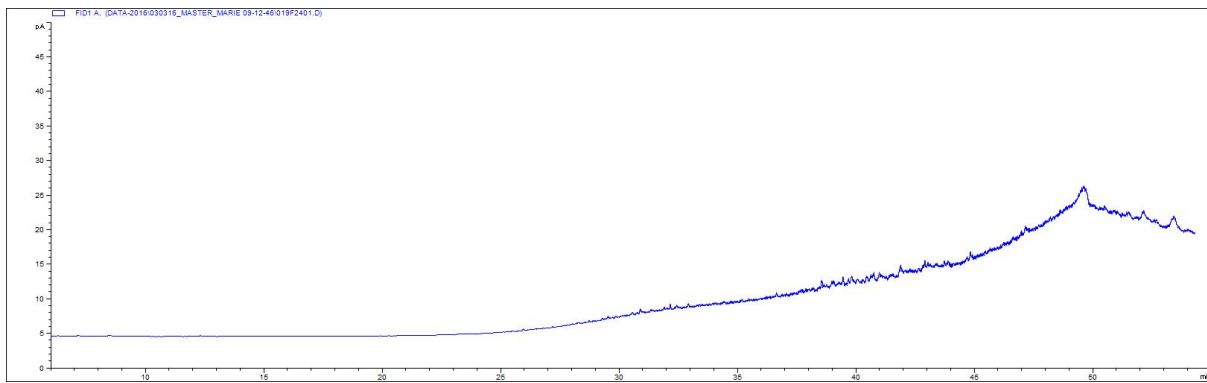
2016-0163



2016-0164



2016-0165



2016-0166

Appendix E

COSIWeb

Table E.1: Diagnostic ratios applied in COSIWeb

Ratio name	Defintion	(m/z)
TS	27Ts/30ab	191
TM	27Tm/30ab	191
28ab	DR-28ab/30ab	191
29ab	29ab/30ab	191
30O	30O/30ab	191
31abS	31abS/30ab	191
30G	30G/30ab	191
27dbR	27dbR/27dbS	217
27bb	$27\beta\beta(S+R)/29\beta\beta(S+R)$	218
TASC26	SC26TA/RC26TA +SC27TA	231
TASC28	SC28TA/RC26TA+SC27TA	231
TARC27	RC27TA/RC26TA +SC27TA	231
TARC28	RC28TA/RC26TA+SC27TA	231
C17/pr	n-heptadecane/pristane	85
C18/ph	n-octadecane/phytane	85
pr/ph	pristane/phytane	85
2MP	2-methylphenanthrene/1-methylphenanthrene	192
MA	Methylanthracene/11-methylphenanthrene	192
4MD	4-methyldibenzothiophene/1-methyldibenzothiophene	198
2MF	2-Methylfluoranthene/4-methylpyrene	216
B(a)F	Benzo(a)fluorene/4-methylpyrene	216
B(b+c)F	Benzo(b+c)fluorene/4-methylpyrene	216

2MPy	2-methylpyrene/4-methylpyrene	216
1Mpy	1-methylpyrene/4-methylpyrene	216
Retene	Retene/Tetra-methyl-phenantrene	234
BNT	Retene/ Tetra-methyl-phenantrene	234
