# SUBSIDIES AND COSTS IN THE CALIFORNIA SOLAR MARKET: AN EMPIRICAL ANALYSIS

Cecilie Teisberg, M.Sc, Norwegian University of Science and Technology, Vardeveien 3 1444 Drøbak Norway, +4745688060, cecilite@stud.ntnu.no

Rakel Håkegård, M.Sc, Norwegian University of Science and Technology, Magnus Åldstedts vei 1 7024 Trondheim Norway, +4798650335, rakelhaa@stud.ntnu.no

*Environmental concerns have prompted governments around the world to subsidize renewable energy markets. One of the major risks associated with subsidizing is that it may inflate costs. Thus, understanding the drivers of costs, and specifically how subsidies affect costs is crucial for evaluating and designing good subsidy policies. In this paper, we identify and estimate the cost drivers of solar photovoltaic systems in the California market using a semi-parametric regression model, and further quantify the cost-inflationary effect by simulation using machine learning techniques. We find evidence for significant cost inflationary effects of subsidies. The regression results suggest that a 1% increase in incentives per kW installed is associated with nearly 0.1% increase in costs per kW installed. Furthermore, simulations indicate that cut-off of subsidies in 2012 would have saved the California government US\$1.15bn, while the extra costs imposed on end-customers would be only US\$0.30bn. Our results suggest that a cut-off in 2012 would not have lead to a substantial jump in costs to end-customers at the cut-off point, and that costs would only be slightly higher for end-customers than with subsidies. The results indicate that an accelerated subsidy down-scaling may be desirable, with minimal adverse implications for end-customers.*

**Keywords: solar power, renewable energy, subsidies, cost-inflation**

## Introduction

Solar power is the most rapidly expanding source of energy in California, and today the state is leading in the US in terms of electricity generation from solar photovoltaics (PV), accounting for nearly half of the US total. A statewide effort to promote growth of the market for solar PV was initiated in 2007, known as Go Solar California. The main component of the campaign was the California Solar Initiative (CSI) which subsidizes roof-top solar PV installations by providing rebates for end-customers. These subsidies have likely been one of the most important drivers of growth in the California solar PV market.

$CO_2$-emission is recognized as a major issue by most governments and efforts to reduce emissions have been initiated all over the world. To this end, subsidizing of renewable energy markets is commonly employed to minimize dependence on fossil fuels as a source of energy, and California is not alone in subsidizing solar PV systems. The US federal government provides tax credit for residential solar systems across the United States, and numerous other countries have introduced subsidies for solar power.

The costs of solar photovoltaics have decreased substantially in recent years, leading more countries to open up to solar power as a viable source of energy, decreasing the dependence on nonrenewable energy and thus reducing $CO_2$-emissions. Understanding the costs of solar PV systems and how subsidies may impact the costs will be crucial when evaluating and designing subsidy policies for emerging solar power markets. One of the major risks associated with subsidizing is that it may inflate costs. Contractors and

manufacturers may see potential to raise prices to end-customers when subsidies boost purchasing power. Thus, the "more is better"-principle does not necessarily apply in the case of subsidies. Thorough analysis of alternative subsidy policies is imperative for identifying policies with the desired properties, i.e. promoting market growth while minimizing cost inflationary effects.

As a first step in analyzing solar PV subsidy policies, the objective of this paper is to investigate any cost inflationary effects of subsidies in the context of the California solar PV market. We do this in two stages. First, we aim to identify the most important cost drivers of solar PV systems, and look specifically at the effect subsidies have on costs. This is achieved using a semi-parametric regression model, enabling us to model complex relationships in the data while also providing descriptive insight. Second, we aim to quantify any cost-inflationary effect by simulating costs under alternative subsidy policies. Cost simulations can be used to evaluate different policies in terms of minimizing cost inflationary effects, which would be valuable for governments or other institutions considering subsidizing solar PV systems in the future. We use machine learning techniques to build a prediction model that can generate simulations. As a benchmark to test the model against, we use our semi-parametric regression model on out-of-sample predictions. We use the prediction model to simulate costs under some simple, alternative subsidy policies, to quantify the cost inflationary effects. Although our analysis is limited to the California solar PV market, the findings are likely to be relevant also for other emerging renewable energy markets.

The rest of this paper is structured as follows. In the following section we present a brief review of existing literature in the field, highlight potential problem areas and place our research into the body of literature. Next, the Data section gives a description of the data and briefly presents the data pre-processing, as well as discuss any limitations of the data. The Methodology section presents a brief introduction to the methodologies used and provides the model specifications. The results are presented and discussed in the Results section. Finally, we summarize key findings and give some recommendations for further work.

## Literature review

Several empirical studies in economic literature have investigated the effects of subsidies on costs. Early work by Pucher and Markstedt (1983) and Feldstein og Friedman (1977) examine costs in the context of mass transit systems and health care systems, respectively, and find evidence for cost inflationary effects of subsidies. For the California solar PV market, Mauritzen (2017) and Wiser, et al. (2006) have conducted empirical studies investigating cost drivers, and a key finding from both of these studies is that higher subsidies are associated with higher cost, hence indicating that cost-inflationary effects of subsidies are present also in this market. However, we find several opportunities for improvement on these studies. Wiser, et al. (2006) use a linear regression model incapable of capturing any non-linear relationships between predictor and response variable. Mauritzen (2017) show that there are non-linear relationships and this should be taken into account. Also, the study by Wiser, et al. (2006) is concerned with the incentive program preceding the CSI program, and thus is out-dated. In the study by Mauritzen (2017), there is a lack of sufficient data pre-processing and identification of linearity in the relationships between response and predictor variables. Furthermore, both studies fail to address certain market effects that could explain the association between higher subsidies and higher costs, and we argue that their conclusions are somewhat pre-mature. In particular, end-users may be encouraged to buy higher quality systems as purchasing power is increased with subsidizing. This effect could be in line with the goals of the subsidy program, and should

not be confused with any cost inflationary effect. We take all these aspects into consideration and aim to improve on existing research.

Avato and Cooney (2008) study how to accelerate clean energy adoption, with focus on R&D. They state that while there are many promising clean energy technologies, most are very costly. Subsidy policies are introduced with the goal of market expansion, and the return on investment for end-customers must be increased in order to increase the number of solar PV installations by end users. Minimizing cost inflationary effects of subsidies is a part of this process. Wiser, et al. (2006) suggest that as the CEC program – predecessor of the CSI program – gradually reduced its incentive levels, system retailers absorbed some of the decrease by reducing prices. Wiser et al. (2006) state that as a result, the net cost to the end user was essentially unchanged as incentives scaled down. In this paper, we investigate the opportunity of accelerated downscaling of subsidies, using machine learning techniques and more recent data.

Whether or not subsidies in solar power markets have the desired effects, and what types of policies are most effective are other important questions. Chernyakhovskiy (2015) examines the effectiveness of policy incentives to increase residential solar PV capacity in the United States, and finds that financial incentives are an important driver of growth. Incentives that reduce up-front cost of adoption and that are subject to low uncertainty are found to have the largest impact. Hsu (2012) investigates the environmental impact of different combinations of promotion policies for solar PV installations in Taiwan. He finds that policies with higher capital subsidy and lower initial feed-in tariff price has the lowest average cost of $CO_2$-emission reduction, out of all the combinations studied. Astbury (2017) argues that the U.S. should focus government investment on R&D instead of on policy mechanisms, in order to most effectively achieve clean energy goals. Our analysis is based solely on the policies used by the CSI-program, which are capital subsidies targeting end-customers. The effect of other policy mixes is therefore not taken into account in this paper.

Several different methods have been employed for descriptive analyses in solar power markets. While Mauritzen (2017) and Wiser, et al. (2006) use semi-parametric and linear regression respectively, Hsu (2012) uses a more complex system dynamics model. For the purpose of this study we find that a system dynamics model is not necessary to capture the effects of interest, and we use a semi-parametric regression model for our descriptive analysis.

We make two main contributions to the existing literature. We improve on existing econometric models of the California solar PV market by enhancing data pre-processing and model identification. Next, we further examine and quantify the cost inflationary effects of the CSI policy using machine learning techniques for simulation.

## Data

For our analyses, we use publicly available data from the California Solar Initiative of more than 140,000 solar PV system installations across the state of California, from the start of the incentive program in 2007, until mid-2017. The data contains 124 attributes for each installation, including for example total cost, incentives received and nameplate capacity. All installations covered by the CSI program are included in the data set. Table 1 contains summary statistics for key variables. We identify 14 outliers with a cost per kW above US$40 000, which are excluded from the data. A brief analysis and justification for the exclusion of these data points is provided in *Appendix A*. Furthermore, 11 data points with a reported cost per kW of US$0 were also excluded.

| | Mean | Median | Min | Max | 1st Qu. | 3rd Qu. |
|---|---|---|---|---|---|---|
| **Installation date, years since 2007** | 5.177 | 5.436 | 0.137 | 10.704 | 3.841 | 6.559 |
| **Cost per kW, US$** | 6206.1 | 5811.1 | 537.3 | 106949 | 4940.0 | 7319.2 |
| **Incentive per kW, US$** | 639.99 | 289.02 | 31.69 | 5623.32 | 172.22 | 953.72 |
| **Nameplate capacity, kW** | 11.80 | 5.39 | 0.92 | 5945.94 | 3.85 | 7.50 |
| **Number of observations** | 142017 | | | | | |
| **% leased** | 48.62 | | | | | |
| **% with Chinese panels** | 22.94 | | | | | |

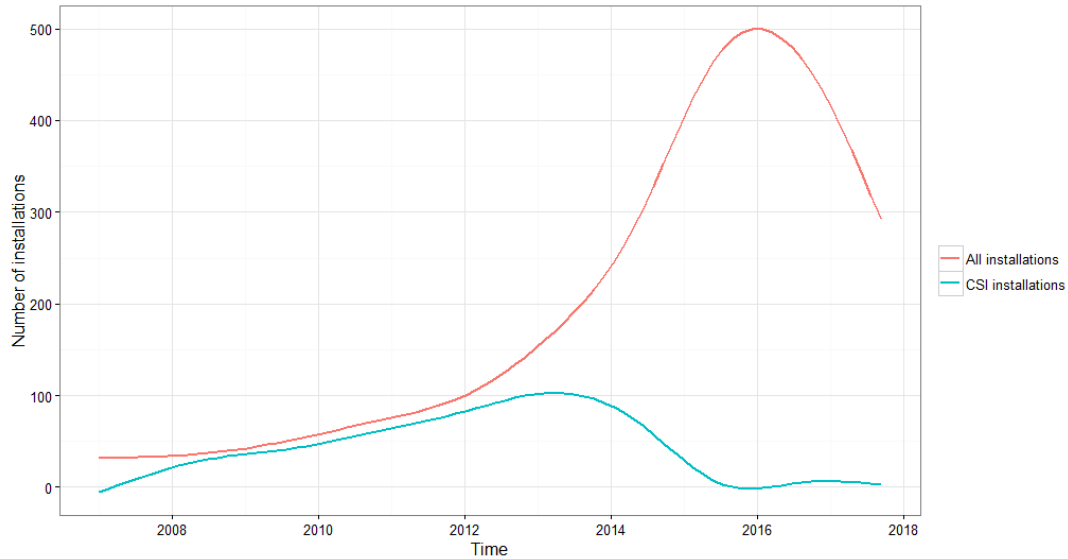*Table 1: Summary statistics for key variables*
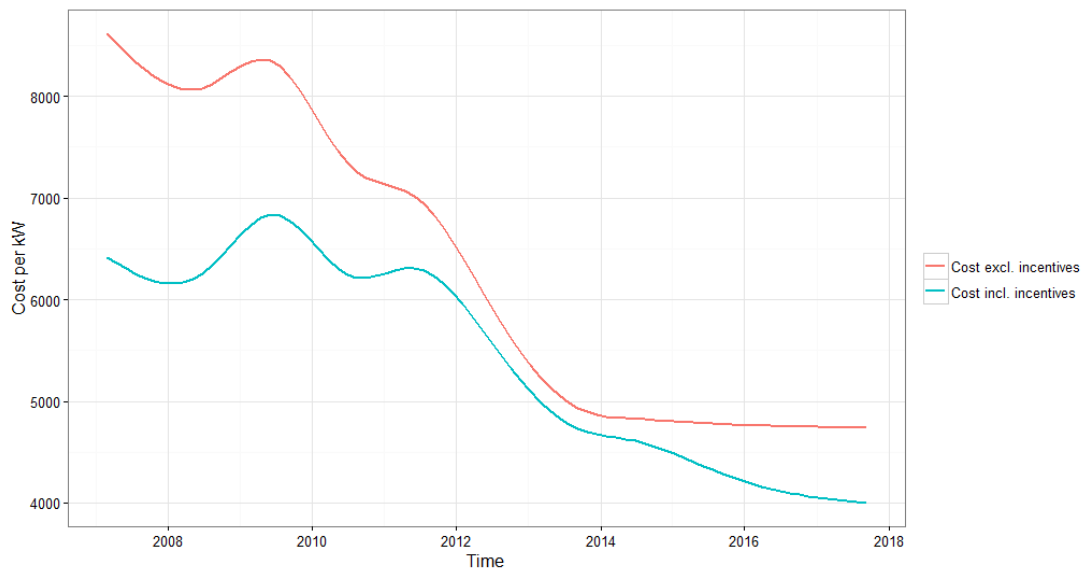


*Figure 1: Number of installations per day*



*Figure 2: Average cost per kW for PV system installations*

The California solar PV systems market has seen substantial growth over the course of the incentive program. Figure 1 shows a smoothed curve of the number of installations per day under the CSI program, along with a smoothed curve for all California installations. The graph shows an exponential increase for the total number of installations from 2007 to 2016, followed by a decrease. On the other hand, the number of installations being subsidized has decreased sharply after 2013. This is due to downscaling of the incentive program. Figure 2 shows smoothed curves of the average cost per kW of installed capacity under the CSI program, with and without incentives, from 2007 to 2017. The costs are reported by the system owners as a part of participating in the incentive program. The dates used for the data points is that of program application approval. It is evident that the overall trend in the market has been steadily declining costs, at least since the beginning of the subsidy program. Only installations covered by the CSI program report data on costs, thus only these are represented in the graph. As the downscaling of incentives has lead to very few data points from around 2015 onward, there is great variance in the average cost of installations in this interval.

Several of the variables of interest are highly non-normal in distribution, and are log transformed to exhibit normality. Although it is not strictly necessary to have normally distributed variables for the results to be meaningful, it minimizes the probability that the regression suffers from high-leverage points skewing the results. It is therefore preferential that the variables are at least approximately normally distributed. We log transform incentives per kW, nameplate capacity, yearly total of installed capacity at the zip code level, contractor size, as well as the response variable, cost per kW.

During data pre-processing for artificial neural networks it is customary to rescale, or standardize, the input variables. This is not strictly necessary when using a Multi-Layer Perceptron (MLP), as we are, however, it can make training faster and reduce chances of getting stuck in local optima. We standardize the input variables by removing the mean and scaling to unit variance. The scaling parameters are computed using the training set only, and scaling is applied to both the training and test set. When making predictions on new input data we scale these using the same scaling parameters computed on the training set.

We identify three main limitations of the data. Firstly, the data only includes installations covered by the CSI incentive program, and as shown in Figure 1 there is a substantial excess amount of installations. Adding baseline data of unincentivized installations would greatly enhance our analyses of the impact subsidies have on costs. Unfortunately, there is currently no data available on the costs of solar PV system installations not covered by the program. Secondly, no data for prices of PV modules is included. Wiser, et al. (2006) argue that these solar module prices are exogenously specified, and is significant in explaining the total cost of solar PV systems. For the GAM regression the effect of the PV module prices are likely absorbed by other variables and will thus not cause any trouble in the descriptive analysis. However, for prediction purposes, the price of PV modules at some time lag is likely to play a significant role, and it would be valuable to add this information in the prediction model. Lastly, the data set contains only promotion policies used by the CSI program. These are capital subsidies based on either expected system performance or realized performance over five years. It is therefore infeasible to test for the effect of other promotion policies such as feed-in-tariffs, net metering, tradable green certificates and tax credits. The US government tax credit level is held constant at maximum 30 % of installation costs during the time period of the data. Thus, we are not be able to measure any effect the federal subsidy might have on costs.

## Methodology

The aim of this paper is to investigate any cost inflationary effects of subsidies in the context of the California solar PV market. Subsidies are predetermined by the California government and are therefore exogenous to the system. Wiser, et al. (2006) and Mauritzen (2017) use linear regression and semi-parametric regression, respectively, to model the costs in the California solar PV market. Potential endogeneity of the explanatory variables, i.e. bidirectional causality with the response, are not discussed in either of the studies. We believe it is plausible that some of the control variables, like nameplate capacity, could have a bidirectional causal relationship with cost per kW. A lower cost per kW could motivate a larger installation (increased capacity), thus forming a feedback loop between cost per kW and nameplate capacity. This would necessitate a model that allows for specification of several endogenous variables, such as simultaneous equations systems, to identify the correct coefficient parameters. However, we believe that this feedback effect in the data is minimal, for example for nameplate capacity other factors like roof area available will likely limit the total capacity of an installation. Thus, we follow the existing literature and use single equation regression to model costs.

Data characteristics must also be accounted for when modelling costs. Mauritzen (2017) find evidence for non-linear relationships in the data, which would make linear regression and Generalized Linear Models (GLM) inappropriate. Hence, we adopt a Generalized Additive Model (GAM), which is a semi-parametric regression that provides descriptive power while allowing for complex, non-linear relationships to be modelled. Another advantage of GAMs over linear regression and GLMs is that model specification and identification is simplified. The non-parametric components of the GAM are able to automatically identify appropriate polynomial terms and transformations of the predictors.

As a second step in investigating any cost inflatory effects, we aim to quantify the impact of subsidies by simulating costs under alternative subsidy policies. In order to obtain plausible cost simulations, a method of high accuracy in terms of out-of-sample predictions should be applied. We adopt a deep neural network, a method that has proven useful in out-of-sample predictions in various applications in recent years, like stock market predictions. The neural network regression approach is completely non-parametric, and hence does not provide any direct descriptive power. However, our aim in this part of the study is not to directly interpret the relationships in the data, but rather simulate alternative market scenarios and quantify the effects of cost inflation.

### *Generalized additive model*

A generalized additive model is a semi-parametric regression model on the form:

$$g(\mu_i) = X_i\beta + f_1(x_{1i}) + f_2(x_{2i}) + \cdots + f_m(x_{mi}) \qquad (1)$$

$X_i\beta$ constitutes the linear component of the model, and the predictor variables, $X_i$, are strictly parametrically specified. The smooth functions $f_{1i}, f_{2i}, \dots f_{mi}$ constitute the nonparametric component of the model, and applies to the covariates, $x_{1i}, x_{2i}, \dots x_{mi}$. In Equation 1, $g$ is a smooth monotonic link function, $\mu_i = \mathcal{E}[Y_i]$, and $Y_i$ is the response variable which follows some exponential family distribution. The distribution of $Y_i$ must be predetermined together with the function $g$.

The smooth functions allow for rather flexible specifications of the dependence of the response variable on the covariates, and is what separates GAMs from GLMs. We estimate the smooth functions using a cubic regression spline. Cubic polynomials are

fitted to the shape in segments and connected at points called knots, such that the function is continuous up to the second derivative.

We use the mgcv package in R to construct the GAM. The model we implement can be written as the following equation:

$$
\begin{aligned}
Log(cost\_per\_kW_i) \\
= \delta_{sector} + \beta_0 + \beta_1 Log(incentive\_per\_kW_i) \\
+ \beta_2 Log(zip\_year\_total_i) + \beta_3 Log(contractor\_size_i) + \zeta_1 lease_i \\
+ \zeta_2 china_i + f_1(time\_years_i) + f_2\big(Log(nameplate_i)\big) + \epsilon_i
\end{aligned}
$$

$$(2)$$

The left-hand side of Equation 2 is the response variable, log cost per kW of installed nameplate capacity. The right-hand side is composed of several terms of predictor variables. $\delta_m$ represents fixed effects, $\beta_m$ represent coefficients of the linear predictors, $\zeta_m$ are the coefficients of dummy variables and $f_k(\cdot)$ are the non-parametric smooth functions.

We are mainly interested in the variable *incentive_per_kW*, in order to analyze the impact of subsidies on costs and investigate any cost inflationary effects. This variable is log transformed and enter the model as a linear component. The other variables are merely control variables in this study. There are four sectors of end-customers covered by the CSI program: residential, commercial, governmental and non-profit. Fixed effects for each sector are captured by $\delta_{sector}$. The variable *zip_year_total* represents the total installed capacity within the zip code of an installation, in the given year, and is log transformed to exhibit normality and included linearly in the model. The variable *contractor_size* captures the market share of the contractor responsible for the installation, in the given year. It is also log transformed and included as a linear component. Two dummy variables are included, *lease* representing whether a system is leased (as opposed to owned by the host) and *china* representing whether the PV modules are from a Chinese manufacturer. Finally, two smooth terms are included, *time_years* which is the number of years since 2007, and *nameplate* which is the nameplate capacity of the installation. Both of these were found to have non-linear relationships with the cost and were therefore included as smooth functions. The nameplate capacity is log transformed as we found this to give a better fit.

### *Artificial Neural Network*

| Model parameter | Value |
|---|---|
| **Number of hidden layers** | 8 |
| **Number of neurons in hidden layers** | 50 |
| **Activation function** | Rectified Linear Unit (ReLU) |
| **Solver** | ADAM |
| **Learning rate** | Set by ADAM solver |
| **Max. iterations** | 500 |
| **L2 regularization term** | 0.0001 |

*Table 2: Specification of model parameters for Artificial Neural Network*

Table 2 shows the specifications for our deep ANN model. We have used the MLPRegressor model from the scikit-learn package in Python to construct the neural network. After testing the model with different numbers of hidden layers and numbers of neurons, using 10-fold cross validation, we find that the best performance is achieved with 8 hidden layers, each with 50 neurons. For training the weights in the ANN we use

the ADAM solver. This is a variation of stochastic gradient descent outlined by Kingma and Ba (2014). The method computes individual adaptive learning rates for different parameters from estimates of first and second moments of the gradients; the name Adam is derived from adaptive moment estimation. For the activation function of the nodes we use the rectified linear unit (ReLU), given by $f(x) = max(0, x)$. This is a very common activation function in ANNs as it makes training easy. The curious reader is encouraged to refer to Zeiler, et al. (2013) for some of the advantages of the ReLU. For regularization, i.e. penalization of complexity in an effort to reduce overfitting, we use the default value of 0.0001.

## Results

Through the descriptive analysis we identify and estimate the cost drivers of solar PV systems in the California market, using the GAM. In the simulation analysis, we quantify cost inflationary effects of subsidies using the deep ANN model.

| | GAM | GAM by Mauritzen (2017) |
|---|---|---|
| (Intercept) | 8.0776*** | 8.6171*** |
| | (0.0104) | (0.0067) |
| Government sector | 0.0885*** | 0.1169*** |
| | (0.0070) | (0.0125) |
| Non-Profit sector | −0.1035*** | −0.0960*** |
| | (0.0083) | (0.0145) |
| Residential sector | 0.0082 | 0.0159* |
| | (0.0046) | (0.0065) |
| incentive_per_kW | 0.0997*** | 0.0012*** |
| | (0.0016) | (0.0000) |
| zip_year_total (MW/year) | 0.0008 | −0.0347*** |
| | (0.0013) | (0.0031) |
| county_year_total (MW/year) | - | 0.0010*** |
| | - | (0.0001) |
| contractor_size (%) | 0.0066*** | 0.0016*** |
| | (0.0003) | (0.0002) |
| nameplate | - | −0.0005*** |
| | - | (0.0000) |
| lease | 0.0387*** | 0.0205*** |
| | (0.0015) | (0.0021) |
| china | −0.0699*** | −0.0575*** |
| | (0.0015) | (0.0024) |
| EDF: s(time_years) | 8.9961*** | 8.9142*** |
| | (9.0000) | (8.9979) |
| EDF: s(nameplate) | 8.9698*** | - |
| | (8.9996) | - |
| AIC | -48886.7751 | 1876602.0511 |
| BIC | -48522.1902 | 1876802.3335 |
| Log Likelihood | 24480.3505 | −938280.1113 |
| Deviance | 5888.9110 | 277851294007.7000 |
| Deviance explained | 0.5202 | 0.3758 |
| Dispersion | 0.0415 | 2608171.0511 |
| $R^2$ | 0.5201 | 0.3757 |
| GCV score | 0.0415 | 2608658.6051 |
| Num. obs. | 141991 | 106551 |
| Num. smooth terms | 3 | 1 |

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$

*Table 3: Statistical data for our GAM model and the model of Mauritzen (2017) for comparison*

### Descriptive Analysis

The results of the GAM regression are shown in Table 3, along with the results of Mauritzen (2017) for comparison. The comparison model was fitted using data from 2007 to 2014. Coefficient estimates for the linear variables are given with the corresponding standard error in parenthesis. For the smooth terms, estimated degrees of freedom are given, where the p-values are from F-tests of whether the smooth terms significantly improve the fit of the model. At the bottom are summarizing statistics. Note that our GAM model has log transformed the response variable, along with incentives per kW, the contractor size and the nameplate capacity, and used an identity link function. Mauritzen (2017) has not transformed any variables and used a log link function. The GAM we present has a substantially higher $R^2$ value and deviance explained than that of Mauritzen (2017), at 0.52 and 0.38 respectively. This indicates that we have succeeded in improving the model, resulting in a better goodness-of-fit.

The coefficient for incentives per kW is significant and positive. According to our results, a 1% increase of incentives per kW installed is associated with nearly 0.1% increase in costs per kW installed. For the model of Mauritzen (2017), a US$1 increase of incentives per kW installed is associated with 0.1% increase in costs per kW installed. The difference in the form of the results is due to the log transformation of incentives per kW in our model, and accounting for this the results are quite similar. Wiser, et al. (2006) find that for the CEC program, predecessor of the CSI program, a US$1 increase in incentive levels yield a US$0.55-US$0.80 change in pre-incentive installed costs.

Wiser, et al. (2006) and Mauritzen (2017) conclude that their results are evidence of cost inflationary effects of the respective subsidies. However, we believe that their conclusion needs more justification, as an important aspect of the market dynamics are ignored. As a result of increased purchasing power due to subsidizing, end-users may invest in higher quality systems, which naturally have higher costs. This effect could be in line with the aims of the subsidy program, and would then not be an adverse effect that should be minimized. The quality of a solar PV system is mostly characterized by its efficiency and its life span. The system efficiency is accounted for in the model, as the costs are on a per kW nameplate capacity basis. Therefore, the life span of the system is the main factor that could contribute to the observed higher costs. Though there is some variabililty in the life expectancy of different solar PV systems, most manufacturers – at least eight of the largest in the California market, which together account for more than half of all systems installed under the CSI program – all have the same industry standard 10-year product warranty and 25-year power output warranty. Thus, there is reason to believe that the average system life-span is not remarkably impacted by the subsidy. Other factors, such as design and customer service, could also cause end-users to choose a higher cost installation, however, we then argue that any increase in average system cost based on such factors is counterproductive to the aims of the CSI program, and thus is comparable to cost inflation by contractors and manufacturers. With these aspects in mind we advise some caution in interpreting the coefficient for incentives per kW as a pure cost inflationary effect, however we find it reasonable to assume that at least part of the effect captured is due to cost inflation.

All of the control variables are significant and as expected, except for *zip_year_total* which is not significant at the 5% level. To verify the validity of the model we have conducted residual tests. Approximate normality of the resiudals was found, and a Harrisson-McCabe test verifies no heteroscedasticity. Complete model robustness analysis of the GAM is presented in *Appendix B*.

### Simulation analysis

The ANN prediction model achieves an average $R^2$ of 0.60 when using 10-fold cross validation to test the precision of out-of-sample predictions. For comparison, the GAM achieves an $R^2$ of 0.51 for out-of-sample predictions. Thus, the ANN model well outperforms the benchmark. To verify the validity of the ANN prediction model, we have conducted residual tests which show that the model is robust. Detailed model validation is presented in *Appendix B*. In the following we use this prediction model to simulate a few simple, alternative subsidy policies, in order to analyze the resulting market scenarios and quantify the effect of cost inflation.

Firstly, we simulate costs under a zero subsidy policy and a 2010 cut-off subsidy policy, which follows the CSI-policy up until 2010 and then drops to zero. Figure 3 and Figure 4 show smooth functions of the simulated average costs per kW for the two scenarios, along with the true and fitted costs per kW for the CSI-program. Both simulations indicate that lower subsidies are associated with lower costs, and are thus in line with the results from the descriptive analysis above. The resulting costs for the zero incentive simulation, given in Figure 3, should be interpreted with some care. The simulation predicts a sudden drop in costs in 2007, when the CSI program was initiated, which is not likely. A more likely scenario would be for the cost to start from the actual cost in 2007 and then exhibit a downward deviation from actual costs over time. We try to incorporate this behavior with the cut-off policy. Figure 4 shows the simulation of this alternative subsidy policy, with cut-off in 2010, exhibiting the expected behaviour of cost per kW, with no sudden jump. The results are in line with the findings of Wiser, et al. (2006), suggesting that intermediaries absorb some of the price increase to end-customers when subsidies are reduced. As discussed in the descriptive analysis, there are other possible explanations for higher costs associated with the subsidy. Nonetheless, the simulations may indicate that subsidies should be scaled down quickly to minimize inflation of costs.
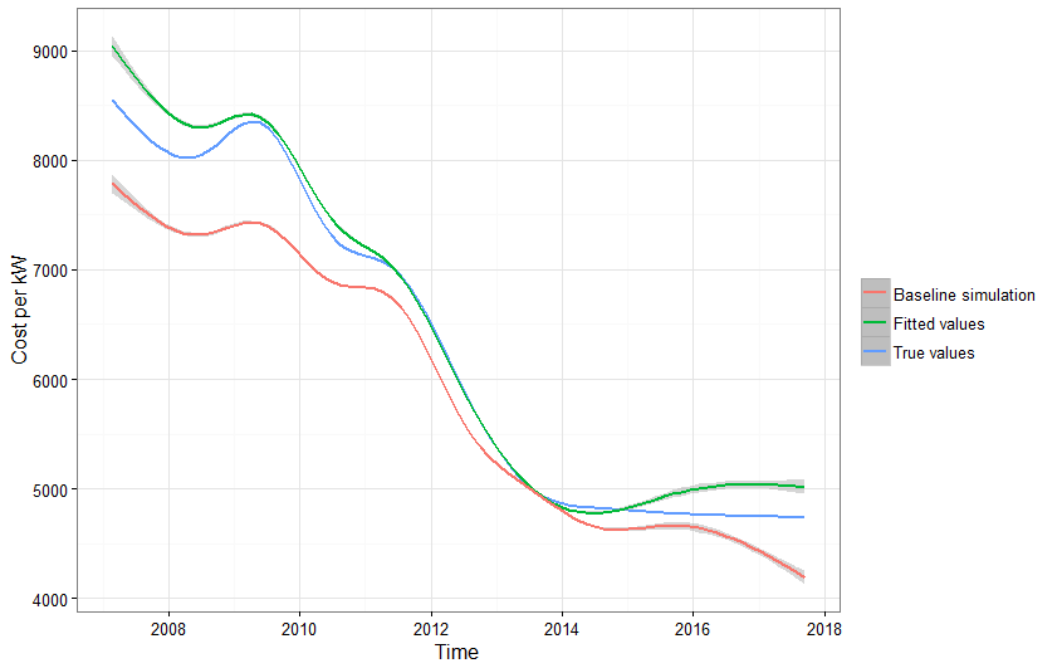


*Figure 3: Cost per kW before incentives over time: Simulated values for zero incentives policy, along with true and fitted values for the CSI policy*
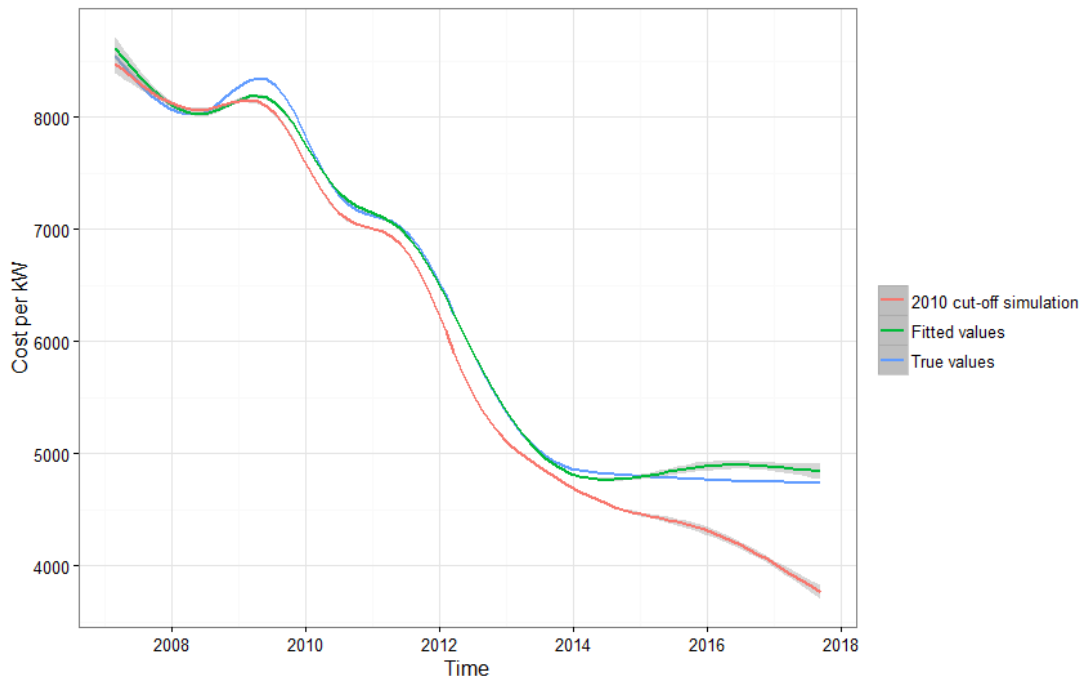
*Figure 4: Cost per kW before incentives over time: Simulated values for 2010 cut-off policy, along with true and fitted values for the CSI policy*

Secondly, it is important to examine how end-users are affected by a drop in subsidies. Although total costs decrease, the cost to end-users may increase, which will likely affect market growth. We assume that the return on investment for end-users drives the demand for solar PV installations. Thus, in order to promote market growth, costs to end-users should be steadily decreasing with time. We examine the cost effects to end-users under a 2010 cut-off policy and a 2012 cut-off policy. Figure 5 and Figure 6 show smooth functions of the simulated average costs per kW including incentives under the two subsidy policies, along with the realized costs per kW including incentives under the CSI policy. The timeline starts at the cut-off point, since costs prior to this would be equal. The results indicate that the 2012 cut-off policy performs best in terms of assuring market growth through steadily decreasing costs. Figure 5 shows that under the 2010 cut-off policy there would be a considerable jump upward in costs to end-customers. Additionally, the costs under this subsidy policy exceed the costs exhibited under the CSI policy by a significant amount between 2010 and 2012. This could potentially damage the market growth, as the return on investment for end-users is substantially lower. Figure 6 shows that these problems are not as prominent for the 2012 cut-off policy, as the difference in costs between the cut-off policy and the CSI-program is smaller and no significant jump occurs in costs at the cut-off point. Our simulations show that the total extra costs imposed on end-users under a subsidy cut-off in 2010 or 2012 would be only US\$0.50bn or US\$0.30bn, respectively. On the other hand, the cost of the incentives for the California government from 2010 until mid-2017 has been nearly US\$1.29bn, or US\$1.15bn from 2012. This indicates a substantial cost inflationary effect of the subsidy, and suggests that accelerated subsidy down-scaling may have been desirable.
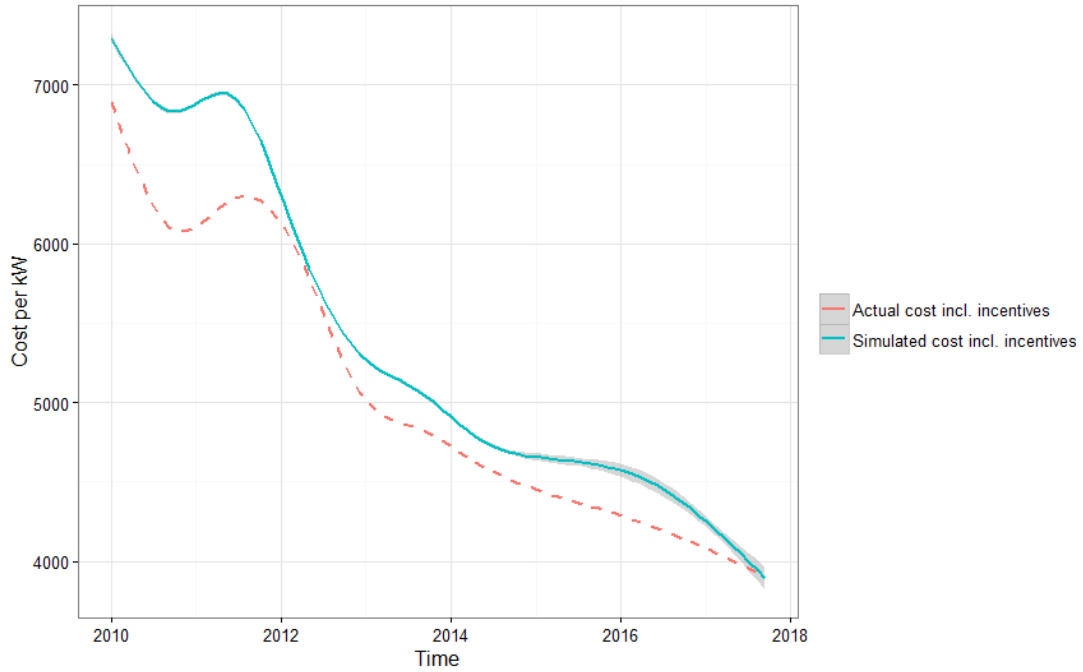
*Figure 5: Cost per kW after incentives over time: Simulations of 2010 cut-off policy, along with fitted values for the CSI policy*
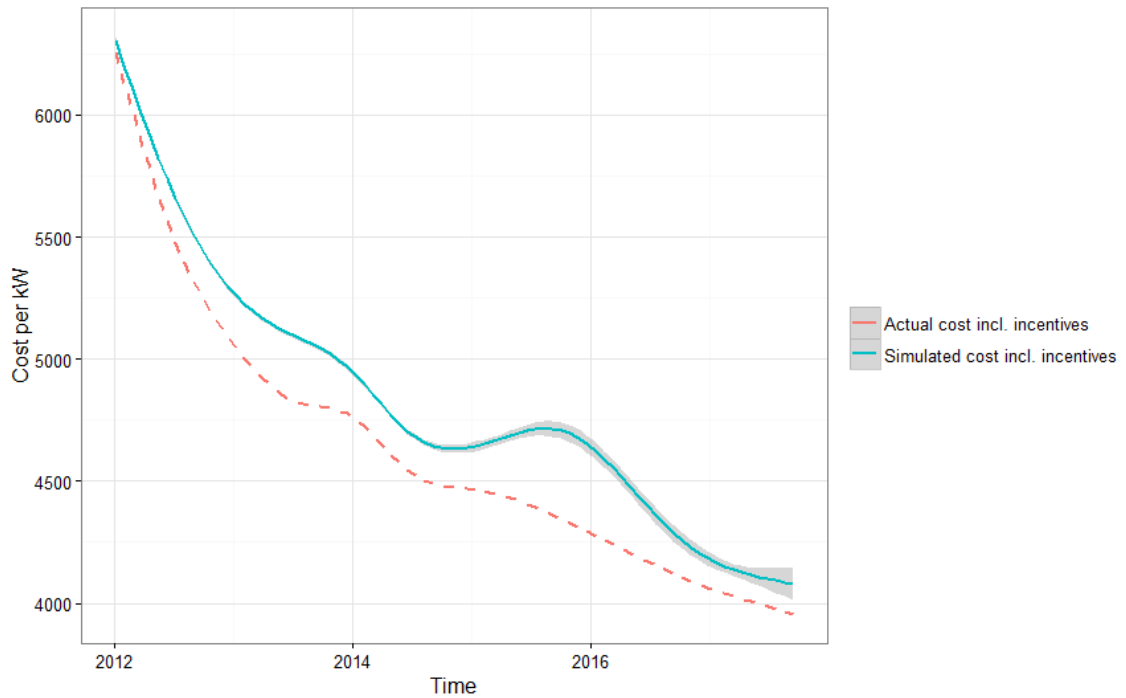


*Figure 6: Cost per kW after incentives over time: Simulations of 2012 cut-off policy, along with fitted values for the CSI policy*

## Conclusions

In this paper, we estimate and analyze cost inflationary effects of subsidies in the context of the California solar PV market. Using a semi-parametric regression technique we model the costs of solar PV installations and analyze the central cost drivers, with focus on subsidies. Furthermore, using machine learning techniques we build a prediction model that well outperforms the benchmark on out-of-sample-predictions. We use the prediction model to simulate costs under a zero subsidy policy and two simple cut-off policies, which follow the CSI policy up until some specified point in time, and then drops to zero.

The results of the GAM provide evidence for some cost inflationary effects of subsidies, confirming the results of existing literature on cost drivers of California solar PV systems. We find that a 1% increase in incentives per kW installed, is associated with nearly 0.1% increase in costs per kW installed. The results for the ANN model also suggest cost inflationary effects of the subsidy. Market simulations for zero incentives and for two different cut-off subsidy policies show lower total cost per kW compared to actual costs exhibited under the CSI policy. The cost effect to end-users, i.e. cost after incentives, is estimated under the 2010 and 2012 cut-off policies. Examination of the cost curves suggest that the 2012 cut-off policy is superior, as costs per kW after incentives do not exhibit any significant jump at the cut-off point, and the costs to end-customers are not notably higher than under the realized scenario. The costs saved by the California government under the 2012 cut-off policy would be US\$1.15bn, while the total extra costs imposed on end-users would be only US\$0.30bn. This suggests that accelerated subsidy downscaling may be desirable.

Our study is concerned only with the effect subsidies have on end-customers, and not how it may affect intermediaries and suppliers. We have found evidence for suppliers "inflating" costs under the California subsidy, meaning that although the subsidy is aimed at end-customers, suppliers are indirectly being subsidized. However, whether this is bad or good is not straightforward. Cutting the indirect subsidy to suppliers may result in adverse market effects. Fewer suppliers entering the market can lead to weaker competition and less investment in R&D. This, in turn, could impair technological improvements needed to enable steadily decreasing costs over time. In order to design good subsidy policies, it is crucial to first evaluate what type of subsidy is preferred. If it is found that the goal should be to subsidize end-customers only, minimizing cost inflationary effects is central. If it is found desirable to subsidize suppliers and intermediaries as well, the cost inflationary effect may give the intended result. Nonetheless, it is likely that a better option for subsidizing suppliers and intermediaries would be a direct subsidy. In that case, any end-customer subsidy, like the CSI, should aim to minimize cost inflationary effects, as we have aimed to in this study.

We note that the increased purchasing power of end-users under the subsidy may encourage higher quality investments, thus explaining some of the increase in average costs found in our analyses. However, we argue that the main quality factor the subsidy program should be concerned with is the system life-span, and as most manufacturers have the same industry standard warranties there is reason to believe the average system life-span in California is not greatly impacted by the subsidy. We conclude that although some caution is warranted in interpreting the results as a pure cost inflation by contractors and manufacturers, it is reasonable to assume that at least part of the effect captured is due to cost inflation, and that our simulation results therefore are significant.

There are several interesting areas of future research on the topic of designing optimal subsidy policies in solar power markets. To find optimal policies in the context of the California solar PV market, more complex subsidy policies than the cut-off

policies we have tested should be evaluated. Moreover, the prediction model could be adapted to new markets and simulations extended to future scenarios, to enable guidance to the choice of subsidy policies for emerging solar PV markets. Lastly, further research on the quality of the systems purchased is needed to be able to verify the true cost inflationary effect, separating out any effect of end-customers investing in better, higher quality systems as their purchasing power is increased, which could be in line with the aims of the subsidy policy.

## References

Astbury, Chrissy. 2017. "How America's Solar Energy Policies Should Follow (and Stray) from Germany's Lead: Working Towards Market Parity Without Subsidies." *Indiana International & Comparative Law Review* 209-245.

Avato, Patrick, and Jonathan Cooney. 2008. *Accelerating Clean Energy Technology Research, Development, and Deployment: Lessons from Non-Energy Sectors.* World Bank Working Paper, Washington DC: World Bank. https://openknowledge.worldbank.org/handle/10986/6528.

Bollinger, Bryan K., and Kenneth Gillingham. 2012. "Peer Effects in the Diffusion of Solar Photovoltaic Panels." *Marketing Science* 31 (6): 900-912. Accessed 4 4, 2018. http://pubsonline.informs.org/doi/abs/10.1287/mksc.1120.0727.

Chernyakhovskiy, Ilya. 2015. "Solar PV Adoption in the United States: An Empirical Investigation of State Policy Effectiveness." University of Massachusetts Amherst.

Feldstein, Martin, and Bernard Friedman. 1977. "Tax subsidies, the rational demand for insurance and the health care crisis." *Journal of Public Economics* 7 (2): 155-178. Accessed 4 4, 2018. http://sciencedirect.com/science/article/pii/0047272777900445.

Folkman, Jim, Katharine Larson, Joseph Omoletski, and Patrick Saxton. 2016. *Guidelines for California's Solar Electric Incentive Programs (Senate Bill 1), Sixth Edition.* California Energy Commision.

Guisan, Antoine, Thomas Jr. C. Edwards, and Trevor Hastie. 2002. "Generalized Linear and Generalized Additive Models in Studies of Species Distributions: Setting The Scene." *Ecological Modelling* 157: 89-100.

Hastie, Trevor, and Robert Tibshirani. 1986. "Generalized Additive Models." *Statistical Science* 1 (3): 197-318.

Hsu, Chiung-Wen. 2012. "Using a system dynamics model to assess the effects of capital subsidies and feed-in tariffs on solar PV installations." *Applied Energy* 100: 205-217. Accessed 4 4, 2018. https://sciencedirect.com/science/article/pii/s0306261912001389.

Kingma, Diederik P., and Jimmy Ba. 2014. "Adam: A Method for Stochastic Optimization." *Computing Research Repository* abs/1412.6980. http://arxiv.org/abs/1412.6980.

Lesser, Jonathan A., and Xuejuan Su. 2008. "Design of an economically efficient feed-in tariff structure for renewable energy development." *Energy Policy* 36: 981-990.

Mauritzen, Johannes. 2017. "Cost, contractors and scale: An empirical analysisof the california solar market." *The Energy Journal* 38: 177-198.

Nyberg, Michael. 2017. *Total System Electric Generation.* 23 06. Accessed 04 06, 2018. http://www.energy.ca.gov/almanac/electricity_data/total_system_power.html.

Pucher, John, and Anders Markstedt. 1983. "Consequences of public ownership and subsidies for mass transit: Evidence from case studies and regression analysis."

*Transportation* 11 (4): 323-345. Accessed 4 4, 2018.
https://link.springer.com/article/10.1007/bf00150722.

Shahan, Zachary. 2016. *California Solar Incentives, Solar Installers, \& Solar Costs.* 14 06. Accessed 04 06, 2018. https://cleantechnica.com/2016/06/14/california-solar-subsidies-solar-installers-solar-costs/.

Wiser, Ryan, Mark Bolinger, Peter Cappers, and Robert Margolis. 2006. "Letting the Sun Shine on Solar Costs: An Empirical Investigation of Photovoltaic Cost Trends in California." *Lawrence Berkeley National Laboratory.* Accessed 4 4, 2018. http://escholarship.org/uc/item/22p2m6vr.

Wood, Simon N. 2006. *Generalized Additive Models: an introduction with R.* Chapman & Hall/CRC.

Zeiler, Matthew D., Marc'Aurelio Ranzato, Rajat Monga, Mark Z. Mao, K. Yang, Quoc Viet Le, Patrick Nguyen, et al. 2013. "On rectified linear units for speech processing." *ICASSP.* Vancouver: IEEE. 3517-3521. http://dblp.uni-trier.de/db/conf/icassp/icassp2013.html#ZeilerRMMYLNSVDH13.

## Appendix A – Outlier Analysis

There are 14 observations in the data with a reported cost per kW above $40 000. As the summary statistics in Table 1 shows, the average cost per kW is $6206, and 75% of the data has cost per kW in the range [$4940, $7319]. Thus, the identified observations have extremely high values for this variable. We note that all these outliers are from the period 2008-2010. As average cost has declined significantly since 2010 this could possibly help explain the high cost values, however, as top panel of Figure 7 shows the outliers have abnormally high costs even for the time period. Comparing values for a the other variables we find that the outliers are within normal ranges: slightly low for *nameplate*, *zip_year_total* and *contractor_size*, and slightly high for *incentive_per_kW*, but given the time period of the observations this is not unexpected. Figure 7 plots the data with three key variables against the *cost_per_kW*. The 14 outliers (circled in the plots) are clearly separated from the rest of the data for all three variables, thus further underpinning the hypothesis of noisy or erroneous data. We therefore remove them from the data before performing further analyses.
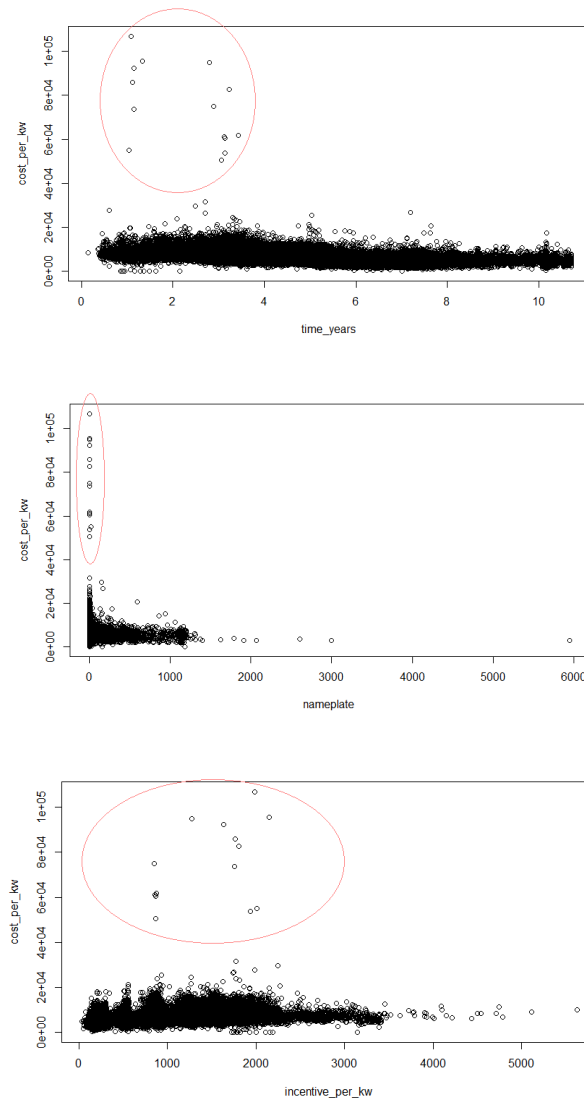


*Figure 7: Plots of three key variables against cost_per_kW, outliers are circled*

## Appendix B – Model Validation

### *Generalized additive model*

The results of the residual tests indicate that the model is robust. Figure 8 shows residual tests for the GAM regression. The QQ-plot compares the distribution of the residuals against the theoretical quantile values of the normal distribution. The distributions of the residuals and the theoretical quantiles should be linearly related, and hence the QQ-plot should be a straight line. Although heavy-tails are present, normality is approximately achieved for the residual interval $[-5000, 2500]$. The histogram of the residuals confirm that almost all residuals lie in this interval. From the plot of the residuals vs. the linear predictors we see that the residuals are distributed quite evenly around zero. The bottom right plot shows the response variable vs fitted values. If all fitted values are correct the points will lie on the line $y = x$, and we can see from the figure that the points form a cluster which follows this line. The residual plot in the upper right shows no sign of heteroscedasticity. This is further underpinned by a Harrison-McCabe test which cannot reject the null hypothesis of homoscedasticity, even at the 30% significance level.
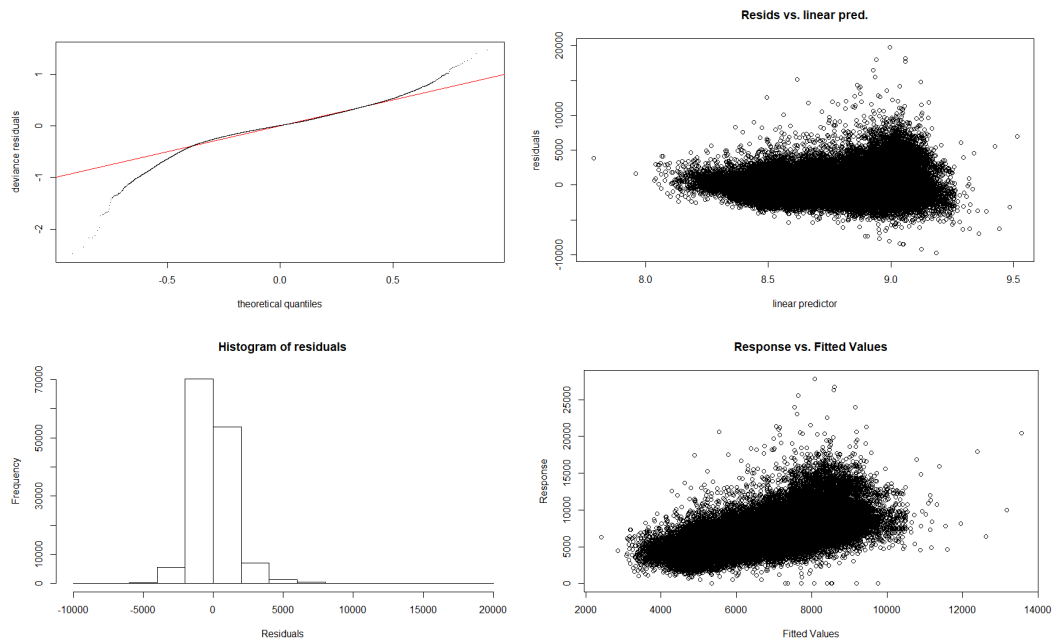


*Figure 8: GAM model residual tests, from upper left: QQ-plot, residuals vs. linear predictor, histogram of residuals, response vs. fitted values*

### *Artificial Neural Network*

The residual plots in Figure 9 show that the ANN regression model is quite robust, with residuals approximately normally distributed. From the QQ-plot it is evident that like the residuals of the GAM model, there are distinct heavy-tails. This means that there are many extreme valued residuals compared to a normal distribution. Within the critical range of 2-3 standard deviations from the mean, the error distribution follows the normal distribution quite closely. The bottom left panel of Figure 9 shows that the distribution is slightly skewed, but has the desired bell-shape. From the residual plot in the upper right and the prediction plot in the bottom right we can see that there seems to be some heteroscedasticity present, as the residuals "fan out" as the predicted cost per kW

increases. The heavy-tails and heteroscedasticity is not likely to present any significant problems in the model as it is solely used for predictions, and because the ANN is more flexible than a GAM and does not make the same assumptions about the residual distribution. On a final note, there is no clear bias in the residuals, as they are approximately symmetric about zero.
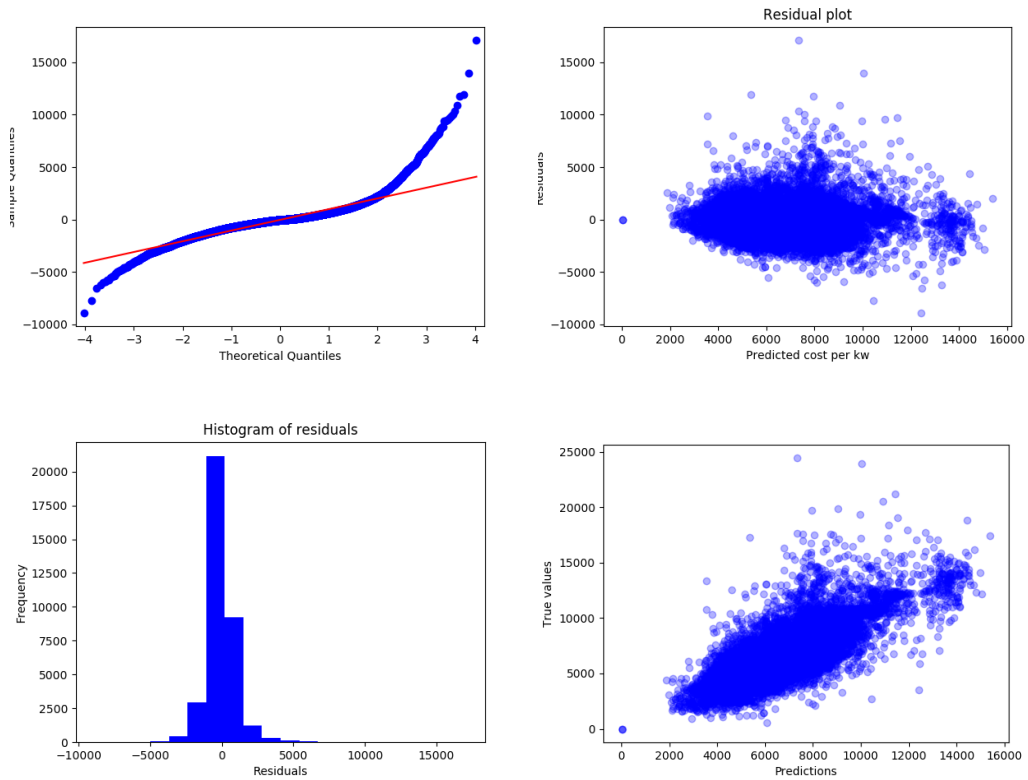


*Figure 9: ANN model residual tests, from upper left: QQ-plot, residuals vs. linear predictor, histogram of residuals, response vs. fitted values*