

# Statistical Modelling of Individual Substations in a District Heating System

---

**LUND UNIVERSITY**

Simon Ingvarsson

Thesis for the degree of Master of Science in Engineering  
Division of Efficient Energy Systems  
Department of Energy Sciences  
Faculty of Engineering | Lund University





# Statistical Modelling of Individual Substations in a District Heating System

Simon Ingvarsson

13 June 2018



**LUND**  
UNIVERSITY

Thesis for the degree of Master of Science in Engineering  
Examiner: Assoc. Prof. Marcus Thern

ISRN LUTMDN/TMHP-18/5413-SE

ISSN 0282-1990

© 2018 Simon Ingvarsson

Division of Efficient Energy Systems  
Department of Energy Sciences  
Faculty of Engineering, LTH – Lund University  
Box 118, 221 00 Lund  
Sweden  
[www.energy.lth.se](http://www.energy.lth.se)

## Abstract

In the rise of digitalization, new possibilities are being discovered for district heating in areas like demand side management and fault detection. For these purposes it is necessary to have reliable models describing the substations. In this thesis, the aim is to develop a group of mathematical models to describe measured quantities and their progress through time, for a well performing substation. For the models to be relevant in applications, they must apply also to other substations. The study explores the possibilities to use the models for fault detection and to track slow drifts in the substations' performance. The method is chosen based on earlier studies on heat load modelling on the system level and all combinations of the endogenous variables heat power and delta-T, and the exogenous variables outdoor and supply temperature are tested. The results show that the best suited model is a SARIMAX  $(0, 1, 1) \times (0, 1, 1)_{24}$ , for any combination of variables. As heat load patterns of individual substations are random in nature it is impossible to create a model with high detail, but it fits the measurements reasonably well. The model for delta-T is applicable also to other substations than the the reference unit, but the heat power model does not perform as well. A sudden fault simulated on one of the substations could be detected as a deviation from the delta-T model and a slow change in the performance of another substation can be detected by re-estimating the model parameters over time. The results are discussed and some ideas for improvement are suggested.

**Keywords:** District Heating, Substations, Time Series Analysis, Statistical Model, ARMA, SARIMA, SARIMAX



## Sammanfattning

Den snabba utvecklingen av nya digitala verktyg leder till nya möjligheter för fjärrvärmebranschen inom områden som laststyrning och feldetektering. För dessa ändamål behövs tillförlitliga modeller som beskriver centralerna. Syftet med denna studie är att utveckla en grupp matematiska modeller som beskriver uppmätta värden och deras utveckling över tid, för en väl fungerande fjärrvärmecentral. För att modellerna ska vara relevanta för tillämpningar måste de fungera på ett flertal fjärrvärmecentraler. Studien undersöker möjligheterna att använda modellerna för feldetektering och att spåra långsamma förändringar i centralernas prestanda. Metoden väljs baserat på tidigare studier av värmelastmodellering på systemnivå och alla kombinationer av de endogena variablerna värmeeffekt och delta-T, och exogena variablerna utomhus- och tillflödestemperatur testas. Resultaten visar att den bäst lämpade modellen är en SARIMAX  $(0, 1, 1) \times (0, 1, 1)_{24}$ , för samtliga kombinationer av variabler. Eftersom värmelastens i enskilda centraler till hög grad är slumpmässiga är det omöjligt att skapa en modell med hög detaljnivå, men modellen överensstämmer trots detta väl med den uppmätta datan. Modellen för delta-T fungerar även för andra byggnader än referensenheten, men värmeeffektmodellen fungerar inte lika bra. Ett plötsligt fel som simulerats på en av centralerna kan detekteras som en avvikelse från delta-T-modellen och en långsam förändring av prestanda hos en annan central kan detekteras genom att skatta om modellparametrarna över tid. Resultaten diskuteras och några förslag till förbättringar föreslås.

**Nyckelord:** Fjärrvärme, Fjärrvärmecentraler, Undercentraler, Tidsserieanalys Statistisk Modell, ARMA, SARIMA, SARIMAX





# Preface

This master thesis was written at the Department of Energy Sciences at Lund University, in close collaboration with NODA Intelligent Systems. It is one of several pre-studies linked to the TEMPO project, within the EU Horizon 2020 framework. I would like to thank my supervisors Per-Olof Johansson Kallioniemi and Sara Månsson at LU as well as Johan Sjögren at NODA for the support during the semester. Valuable input, especially in the early stages of the work, was also given from Jens Brage and Christian Johansson at NODA.

Simon Ingvarsson, june 2018



# Terms and Abbreviations

Term	Description
ACF	<i>Auto-Correlation Function.</i> Tool to estimate the order of an MA-process.
AR-process	<i>Auto-Regressive process.</i> In such a process every value depends linearly on past values in the sequence.
ARIMA-process	<i>Integrated Auto-Regressive Moving Average process.</i> Extended version of the ARMA-process. It includes a difference operator to fit non-stationary processes.
ARMA-process	<i>Auto-Regressive Moving Average process.</i> A combination of the AR and MA processes.
Customer Installation	The secondary system and the substation, combined.
Delta-T	The difference between supply and return temperature on the primary side of the substation.
Difference Operator, $\nabla$	Eliminates non-stationary behaviour of a process by transforming every value to the difference between itself and its neighbour.
Endogenous Variable	Variables whose values are determined by other variables in the modelled system.
Exogenous Variable	Variables that affect the modelled system without being affected by it.
Heat Power	Hourly values of the heat load.
MA-process	<i>Moving Average process.</i> In such a process every value depends linearly on past and present values of a random noise sequence.
PACF	<i>Partial Auto-Correlation Function.</i> Tool to estimate the order of an AR-process.

Primary System	The part of the district heating system that contains the heat generating facilities and the distribution pipes.
SARIMA-process	<i>Seasonal Integrated Auto-Regressive Moving Average process.</i> Extended version of the ARMA-process. It includes a seasonal difference operator and seasonal AR and MA polynomials to fit processes with a periodic cycle.
SARIMAX-process	A SARIMA process combined with linear regression against some exogenous variable.
Seasonal Difference Operator, $\nabla_S$	Eliminates periodic behaviour of a process by transforming every value to the difference between itself and its value one full period back.
Secondary System	The heating system inside the customer building.
Substation	The component that connects the primary and secondary systems.
$(p, d, q)$	Orders of an ARIMA process. $p$ and $q$ are the number of coefficients in the AR- and MA-polynomials respectively. $d$ is the number of difference operators applied.
$(P, D, Q)_S$	Additional orders of a SARIMA process. $P$ and $Q$ are the number of coefficients in the seasonal AR- and MA-polynomials respectively. $D$ is the number of seasonal difference operators applied and $S$ is the season length.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	1
1.2	Aim . . . . .	2
1.3	Limitations . . . . .	2
1.4	Thesis Outline . . . . .	3
<b>2</b>	<b>District Heating</b>	<b>5</b>
2.1	District Heating Systems . . . . .	5
2.2	Heat Load Variations . . . . .	7
<b>3</b>	<b>Mathematical Tools</b>	<b>11</b>
3.1	Linear Regression . . . . .	11
3.2	Stochastic Variables . . . . .	12
3.3	Stationary Stochastic Processes . . . . .	13
3.4	ARMA-processes . . . . .	13
<b>4</b>	<b>Earlier Work</b>	<b>17</b>
<b>5</b>	<b>Method</b>	<b>19</b>
5.1	Model Creation . . . . .	19
5.2	Model Evaluation . . . . .	21
<b>6</b>	<b>Modelling Results</b>	<b>23</b>
6.1	Modelling Heat Power . . . . .	23
6.2	Modelling Delta-T . . . . .	27
6.3	Summary of Models . . . . .	28
<b>7</b>	<b>Model Performance</b>	<b>29</b>
7.1	Heat Power Models . . . . .	29
7.2	Delta-T Models . . . . .	33
7.3	Other Buildings . . . . .	37
<b>8</b>	<b>Detecting Drifts and Faults</b>	<b>39</b>
8.1	Detection of Slow Changes in Performance . . . . .	39
8.2	Detection of Sudden Faults . . . . .	41
<b>9</b>	<b>Discussion</b>	<b>43</b>
9.1	Model Structure and Performance . . . . .	43
9.2	Applicability to Other Buildings . . . . .	44
9.3	Applicability in Fault Detection . . . . .	44
9.4	Suggestions for Future Work . . . . .	45
<b>10</b>	<b>Conclusions</b>	<b>47</b>
<b>11</b>	<b>References</b>	<b>49</b>



# 1 Introduction

## 1.1 Background

Central heating solutions have been around for centuries and since the 1960's modern district heating is widely used in some countries, mostly in northern and eastern Europe. Although the technology has been used for a long time it has not been given much attention from academia until the last decades, and little exchange has been done between countries and with other scientific disciplines. A new wave of interest is seen now, in Sweden and in the rest of Europe, and not least in China. In search of efficient and sustainable heating solutions many look to district heating as a way to make use of excess energy from industries and data centers, or to raise efficiency in their local energy systems by using combined heat and power plants.

Energy companies that provide district heating to the customers face a number of challenges to keep costs and emissions as low as possible, while delivering the heat their customers expect. To optimize the system performance the heat production must be controlled to match the demand from customers, and the temperature of the circulating water must be kept low to avoid losses and increase efficiency. One of the most common problems is that the substations, that make up the interfaces between the network and the customer buildings, often perform sub-optimally. Commonly, district heating systems are optimized only by controlling the production unit, and a lot of the fault detection at customer site is done manually.

Now, however, in the emerging era of Digitalization and Internet of Things, a range of possibilities are being discovered. In a smart district heating system, substations are able to provide high-resolution data on flows, temperatures and energy use to the network owner. This gives possibilities to cut peak loads and optimize network performance by demand side management. It also opens up for new ways to detect faults using data analysis. For both these purposes it is necessary to have reliable mathematical models describing the substations.

A lot of work has earlier been done on the optimization and modelling of district heating systems. Since detailed data from individual substations have not been accessible typically such research has focused on optimizing the heat plant. As we now have access to more detailed data, it is possible to zoom in and study the characteristics of each individual substation.

## 1.2 Aim

The aim of this thesis is to explore new ways to examine the characteristics of individual customer installations in a district heating system. This will lead to insights supporting further research on system optimization and fault detection.

Given data from a well performing substation, relationships between measured physical quantities and their progress through time will be studied. From these insights, the aim is to create one or several mathematical models that describe the substation and its connected system. For the models to be relevant for future applications, it is necessary that they are applicable also to other substations than the one used in the model development. Therefore it will be investigated how the model performs when compared to measurements from several buildings.

This thesis will also explore the possibilities to use the models for fault detection and to track slow changes in the substations' performance.

The developed models will be evaluated by answering the following research questions:

- Do the model predictions match the measured data?
- Are the models applicable to several buildings?
- Can the models be used to detect faults or slow changes in the substation's performance?

## 1.3 Limitations

The study is limited to use only the physical quantities available from the substation metering, and hence does not include solar gains and wind chill or any measurements from the customer side. This is partly to limit the time and work needed, but also because it is interesting to see what can be done with the data that is usually available to the network owners.

To further limit the size of the study, only multi-family houses are considered. These constitute a large part of most district heating networks and customer behaviour is quite regular with little difference between weekdays and weekends. Thus the results from this study cannot be extrapolated to other kinds of buildings. When examining if the model is applicable to many buildings, the study only investigates a handful of buildings. All the buildings are from the same system, so it is not covered how well the models apply to substations in other cities.



## 1.4 Thesis Outline

In the following two sections some background theory will be presented, first on district heating technology and then on statistical methods for time series modelling. Readers well familiar with either of the subjects may skip one of the two sections without losing context. Section four gives an overview of earlier studies on modelling of district heating systems that have influenced the choice of methods for this study and provides a context for the results achieved here. Section five describes the methods used for this work. The results are presented in sections six to eight. In section nine the results are analyzed and discussed, relating to the aim of the thesis as well as earlier and future studies. The choice of methods is also discussed. In section ten conclusions are drawn and the questions posed in the introduction are answered.



## 2 District Heating

This section gives a background theory on District Heating technology. Special emphasis is placed on seasonal and daily variations in the heat load. The content is based on (Frederiksen & Werner, 2013).

### 2.1 District Heating Systems

Demands of heat with lower or higher temperature appears in the industrial, residential, service and agricultural sectors. Historically the heat has mostly been produced locally where it is needed. During the 20th century however, centralized heating solutions have become an option for low temperature heat demands such as space heating and hot water preparation in residential, office and service buildings.

The idea behind district heating is to use central heat production to support entire areas with heat, instead of having separate solutions in every single building. The heat can be supplied from a variety of sources, including plants burning fossil fuels, biofuels or waste, large scale electric heat pumps or as waste heat from industries. By using either combined heat and power plants or waste heat from industries, energy is used that would otherwise be wasted.

The heat is carried by some medium that circulates between the production site and the customers through underground pipes, as illustrated in figure 1. Different pipe types and transporting media have been tried in different countries over time. The so called third generation of heat distribution technology, which is the present generation in Europe, transfers water through pipes with foam insulation and plastic casing. In an average Swedish system the supply temperature from the production unit is 86 °C and the return temperature after delivery is 47 °C. The flow around the system is driven by pressure differences in supply and return pipes. Even with modern pipes some heat loss is inevitable. This, in combination with installation costs, limits the usefulness of district heating to densely populated areas.

The part of the network so far described is called the primary system. This is connected to the customer installations where it supports space heating and hot water preparation in one of several possible ways. There are examples of installations where the water from the primary system flows directly through the buildings. Much more common in Sweden though is that every building has its own secondary system that circulates heat around the building for space heating and hot water preparation, connected to the primary side system by a heat exchanger. This allows for a high pressure on the primary side without needing expensive equipment on the secondary side and reducing the risk of accidents.

The space heating circuits distribute heat in the building through circulating water to radiators and underfloor heating, if any. Hot water for taps, showers etcetera



Figure 1: A conceptual image of a simple district heating system. The red and blue lines represent the hot and cold water that is circulated from production site to the customer buildings and back.

can either be heated instantaneously by letting incoming cold water pass by a separate heat exchanger, or by using a coil around a storage tank. In Sweden instantaneous heating is the more common of the two, especially in apartment buildings.

### 2.1.1 Substations

The heat exchanger that links the two systems is part of a larger unit called a substation, symbolized with a grey box in Figure 1. Apart from the heat exchanger the substation also features sensors, meters, valves and controls to regulate the secondary system. The substation meters measure supply and return temperatures and water flow on the primary side, as well as the outdoor temperature. These are connected to a control box that controls the flow through the heat exchangers to keep the indoor temperature constant at the desired level. The control box has been calibrated to a curve fitting the required heat energy for different outdoor temperatures.

Fundamentally, the amount  $Q$  [J] of delivered heat is given by

$$Q = \int^V c_p \rho * \Delta T dV \quad (1)$$

where

$V$  is the volume passing from one reading to the next [ $m^3$ ]

$c_p$  is the specific heat capacity of the fluid [ $J/kgK$ ]

$\rho$  is the density of the fluid [ $kg/m^3$ ]

$\Delta T$  is the temperature difference [ $K$ ].

It follows from equation 1 that an increased heat demand can be met either by increasing the volume flow or by increasing the difference between supply and

return temperature (delta-T). The substation works optimally when delta-T is kept as large as possible. On the system level it is also important to keep the overall temperatures as low as possible for maximum efficiency of the production facilities and to keep distribution losses low.

For billing purposes heat load values are stored per month, but in many substations heat load values per hour can also be extracted. When working with hourly values, the heat load is commonly referred to as Power with unit kilowatt although it is actually the energy that has been calculated. This works naturally as the total heat energy during one our, in kilowatt-hours, is equivalent to the average heat power.

A majority of substations perform sub-optimally. There are many different causes of faults. They can be categorized as construction errors, malfunctions, incorrect settings, errors in the surroundings or statistical errors. Some faults strike suddenly, while others are successively increased over time. Typically a faulty substation will have a lower delta-T, compensated by a larger flow. This is inefficient for the individual building and as the problem is very common the insufficient cooling also affects the return temperatures in the primary system, leading to inefficient system performance.

## 2.2 Heat Load Variations

The heat load consumed by a specific building is a consequence of the building's heat demand. For residential buildings as well as most commercial and public sector buildings this demand consists of space heating and hot water preparation. The space heat demand is the amount of heat required to heat incoming cold ventilation air and compensate for heat transmission through walls and windows, plus minus other gains and losses. For a typical building the space heat demand  $P$  can be expressed as

$$P = C * (T_{in} - T_{out}) + \text{wind chill} - \text{solar gains} - \text{indoor heat gains} \quad (2)$$

where  $T_{in}$  and  $T_{out}$  are the indoor and outdoor temperatures.  $C$  is a constant sum of a heat transmission factor and a ventilation rate factor.

The heat load shows regular variations with yearly, daily and, for some building types, weekly cycles. The seasonal variations are of much larger magnitude than daily variations. The patterns of the daily variations also varies with season. While the seasonal differences mainly originate from weather effects, the daily variations are mostly linked to consumer behaviour. For example, heat demand is lowest during nighttime although it is colder outside, because of low water usage and few door/window openings. The impacts of social behaviour lead to different daily load patterns for different building types. In many commercial and public sector buildings the ventilation is often turned off at nighttime, and sometimes

also during weekends, creating a clear weekly pattern. This weekly pattern is significant also when studying the load for the network as a whole.

Figure 2 shows the average weekly pattern for a whole district heating system, where the year has been divided into four seasons, based on heat load patterns in northern European climate conditions. During Mid-winter(December-February) the differences between day and night are significant, with two daily peaks at morning and evening. In Early spring / Late autumn (March-April and October-November) the solar gains around noon eliminates most heat demand during mid-day, but clear peaks at morning and evening remain. In Late spring / Early autumn (May and September) there are pronounced peaks only in the morning and for the Summer (June-August) almost no daily variations in the heat load can be seen.

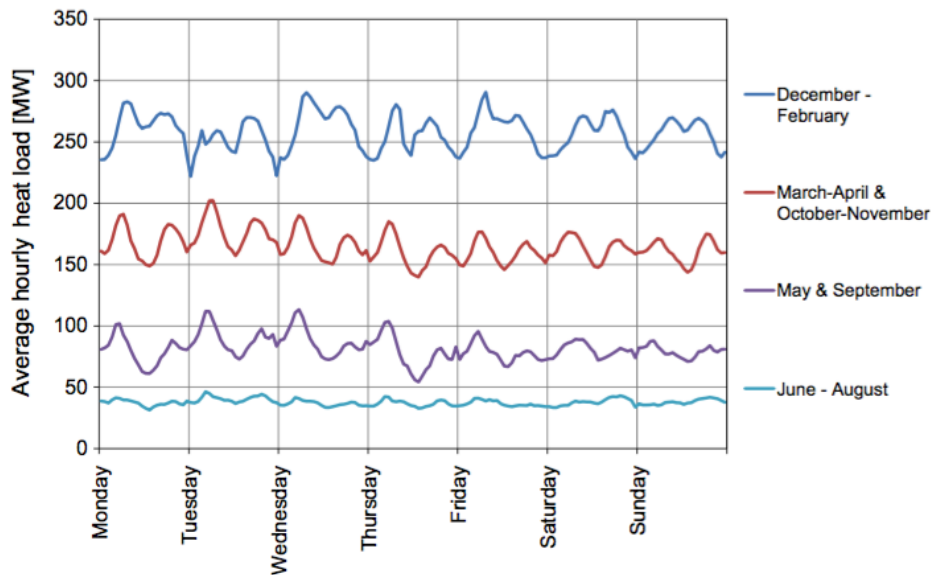


Figure 2: Average weekly heat load patterns for a district heating system in Sweden, categorized by season. (Gadd & Werner, 2013a)

For individual multi-family houses, the heat load for space heating demands largely follow that of the system as a whole. The heat power demand for hot water preparation depends mainly on the desired hot water flow for uses like tap water and showers, that use a lot of heat during a short period of time. A typical resident only uses hot water for 15 minutes per day. Therefore the pattern for a single building is irregular with sudden spikes. Looking at the network as a whole these peaks are evened out due to diversity effects, but for individual customers this adds to the social effects. As most water usage in residential houses occur during morning and evenings and typically no water is used during nighttime this amplifies the pattern of the space heating and also gives some daily cycle during the summer. The hot water heat demand is less seasonally dependent than the space heating demand but a bit lower during summer than winter.

The weather effects affecting the heat demand naturally affects the load patterns. Most importantly, the load varies according to the energy signature of the network

or building. The heat load typically depends linearly on the outdoor temperature up to some break point where space heating is no longer demanded and from there the signature profile is flat. Solar gains and wind chills are also significant, although they have much smaller effect than temperature. Solar gains lowers the heat demand especially at temperatures between 8 and 16 °C, i.e. during Early spring / Late Autumn as seen in the weekly pattern. The wind chill, on the other hand, adds to the heat demand especially at temperatures between -4 °C and 8 °C.

Over the year, for a typical example in nordic conditions, space heating makes up around 60 % of the heat load and hot water preparation makes up around 30 %. In the summer however space heating is little used or completely turned off so heat is used only for hot water preparation. The annual heat load reduction from solar gains is 5 % of the total and the corresponding heat loss from wind chill is 4 %.





## 3 Mathematical Tools

This section explains the mathematics behind the models developed in this thesis. The two main tools are Linear Regression and ARMA-filters. Before introducing the latter, brief introductions of Stochastic variables and processes are given as a framework. Sections 3.1 and 3.2 are based on (Blom, Enger, Englund, Grandell & Holst, 2005). Section 3.3 is based on (Lindgren, Rootzén & Sandsten, 2014). Section 3.4 is based on (Jakobsson, 2015).

### 3.1 Linear Regression

In statistical modelling two types of variables are used. Endogenous variables whose values are determined by other variables in the modelled system, and exogenous variables that affect the modelled system without being affected by it. Modeling a district heating system, for example, the energy consumption could be selected as an endogenous variable and the outdoor temperature would be an exogenous variables.

A simple model to describe the relation between two series of values is the linear regression model. In such a model the relationship between the exogenous variable  $x$  and the endogenous variable  $y$  is given by

$$y = \alpha + \beta x \tag{3}$$

where  $\alpha$  and  $\beta$  are constants. If  $x$  and  $y$  are plotted against each other, the points  $(x_i, y_i)$  in the diagram are expected to follow a straight line with slope  $\beta$  and offset  $\alpha$ . Creating such a model, the target is to estimate  $\alpha$  and  $\beta$  so that all measured values are as close as possible to the straight line. One method for this (that has been proved to be optimal in conditions true for most applications) is the Least Squares method. This method aims to minimize the function  $Q$ , defined as

$$Q(\alpha, \beta) = \sum_{i=1}^n (y_i - \mu_i)^2 \tag{4}$$

where  $\mu_i = \alpha + \beta x_i$ . There are many ways to numerically find the minima.

Linear regression models can naturally be extended to an arbitrary number of exogenous variables.

## 3.2 Stochastic Variables

A stochastic variable is a random variable. Instead of having a set value like a classic variable, it is described by its statistical properties, giving the probability of different realizations (outcomes). These properties are expected value, variance and distribution. The expected value  $E$  of a stochastic variable  $X$  is defined as

$$E(X) = \int_{-\infty}^{\infty} x f_X(x) dx \quad (5)$$

where  $f_X(x)$  is the probability density function that for every value  $x$  describes the probability that  $X$  takes this value when realized. The expected value can be interpreted as the centre of gravity of the distribution function. For a large number of realizations, the expected value coincides with the mean value of these, denoted  $\mu$ . The variance  $V$  is a measure of the spread of the realizations, defined as the expected value of the deviation from the mean, squared, according to

$$V(X) = E[(X - \mu)^2]. \quad (6)$$

Often the standard deviation,  $\sigma$ , is used instead of the variance. This is the square root of the variance, and is convenient to use as it has the same dimension as the original variable. Because of its relationship to the standard deviation the variance is also denoted  $\sigma^2$ .

Stochastic variables with the same expected value and variance, can still differ in the shape of their probability distribution. Two examples of common distribution types are the rectangle distribution that gives equal probability for all values between two limits, and the normal distribution which is often encountered in practical applications and natural phenomena. The probability density function of the normal distribution is given by

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2} \quad (7)$$

where again  $\mu$  is the expected value and  $\sigma^2$  the variance.

Given two stochastic variables  $X$  and  $Y$ , it is sometimes of interest to study their dependence on each other. One measure of this is the Covariance,  $C$ , defined as

$$C(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]. \quad (8)$$

Note that if  $X = Y$  this is the expression for variance. Another measure is the Correlation,  $\rho(X, Y)$  which is the covariance normalized to values between 0 and 1.

### 3.3 Stationary Stochastic Processes

A stochastic process  $\mathbf{Y} = Y_1, Y_2, Y_3, \dots$  is a sequence of stochastic variables. Like the stochastic variable the process only takes the form of specific values when the process is realized. Thus one stochastic process can lead to many or even infinite combinations of different realizations.

A stochastic process is called stationary if the expected value and variance is constant over the sequence and the correlation between any two variable  $Y_t, Y_s$  only depends on the distance  $\tau = |t - s|$ .

A special example of a stationary stochastic process is the White Noise process. This is a completely random, normally distributed process where there is no correlation between any variables in the sequence.

Stationary stochastic processes occur in many different fields of science and technology, where some examples are ocean waves, speech recordings, financial markets and, indeed, energy consumption. In many applications, including this thesis, the stochastic process is a series of measurements of some physical quantity at different moments in time. In the following this will be referred to as a Time Series, and the distance  $\tau$  is equal to the time lag between two data points.

### 3.4 ARMA-processes

An important group of stationary stochastic process are Auto-Regressive (AR) or Moving Average (MA) processes, and combinations and variations of these such as ARMA, ARIMA and SARIMA processes. In fact, any stochastic process can be expressed on this form, given enough parameters. Especially in the field of Time Series it is common to use these processes as models for any real process.

For an Auto-Regressive process every value  $y_t$  in the series can be expressed as a linear regression against the earlier values in the series, as in

$$y_t = a_1 y_{t-1} + a_2 y_{t-2} + \dots + a_p y_{t-p} + e_t \quad (9)$$

where  $e_t$  is an innovation term, randomly picked from a normal distribution.  $p$  is called the order of the model.

The Moving Average process, on the other hand, is a process where every value can be expressed as a linear regression against the innovation terms from earlier points in the series, according to

$$y_t = e_t + c_1 e_{t-1} + a_2 c_{t-2} + \dots + c_q y_{t-q}. \quad (10)$$

Here  $q$  is the model order. Combining these two gives a so called ARMA process of order  $p, q$ , expressed as

$$y_t = a_1 y_{t-1} + a_2 y_{t-2} + \dots + a_p y_{t-p} + e_t + c_1 e_{t-1} + c_2 e_{t-2} + \dots + c_q e_{t-q}. \quad (11)$$

Now, introducing the Z-transform, the operator  $z^{-k}$  is defined by the expression

$$z^{-k} y_t = y_{t-k}. \quad (12)$$

In this notation, the AR-polynomial  $A(z)$  and the MA-polynomial  $C(z)$  are defined by rewriting equations 9, 10 as

$$(1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_p z^{-p}) y_t = e_t \quad (13)$$

$$A(z) y_t = e_t \quad (14)$$

and

$$y_t = (1 + c_1 z^{-1} + c_2 z^{-2} + \dots + c_q z^{-q}) e_t \quad (15)$$

$$y_t = C(z) e_t \quad (16)$$

respectively. AR and MA can be regarded as special cases of the ARMA model, expressed on this form as

$$y_t = \frac{C(z)}{A(z)} e_t. \quad (17)$$

In this notation, the ARMA process can be regarded as a white noise sequence filtered through a black-box system, completely described by the  $A(z)$  and  $C(z)$  polynomials. The filtering is mathematically performed as a convolution.

When modelling a Time Series as an AR, MA or ARMA process, the fit of the model can be evaluated by filtering the series through the inverse of the process and study how well the residual resembles white noise. If the model residual (the difference between the modelled and real data) is completely random, then it is impossible to create a better model.

For the process to be stationary and invertible, all roots of the  $A(z)$  and  $C(z)$  polynomials must be inside the unit circle.

Sometimes it is of interest to make a model of a non-stationary process. For this it is necessary to add some extensions to the ARMA-model, creating the integrated and seasonally integrated ARMA models, ARIMA and SARIMA. An ARIMA model utilizes the fact that many non-stationary processes have a stationary derivative. Instead of modelling the process itself as an ARMA, its derivative is modelled as an ARMA-process. For a discrete process the analogy of a derivative is simply the difference between adjacent values. To differentiate a process, the process is convoluted with the difference operator,  $\nabla$ , defined in z-notation as

$$\nabla = (1 - z^{-1}). \quad (18)$$

It is also possible to differentiate the process more than once, if necessary to make it stationary. The ARIMA process of order (p,d,q) is accordingly defined as

$$y_t = \frac{C(z)}{\nabla^d A(z)} e_t \quad (19)$$

where d is the number of differentiations. Sometimes the process is non-stationary in the way that it has some seasonal cycle with period  $S$  (typically a diurnal, weekly or yearly cycle for practical applications). Then another difference operator,  $\nabla_S$ , is applied so that every value in the transformed series is the difference between the original value and the value one cycle back.  $\nabla_S$  is defined as

$$\nabla_S = (1 - z^{-S}). \quad (20)$$

In its complete form, the SARIMA model also includes seasonal AR and MA polynomials, enabling regression on values and innovations an integer number of cycles back. These are denoted  $\mathcal{A}(z^S)$  and  $\mathcal{C}(z^S)$  and their orders are  $P$  and  $Q$  respectively. The complete SARIMA model of order  $(p, d, q) \times (P, D, Q)_S$  is given by

$$y_t = \frac{\mathcal{C}(z^S)C(z)}{\nabla_S^D \nabla^d \mathcal{A}(z^S)A(z)} e_t. \quad (21)$$

The various kinds of ARMA-models can be combined with linear regression against exogenous variables. In its simplest form, sometimes referred to as an ARMAX, the exogenous regression is added as an extra term to equation 17. A more advanced approach is the Box-Jenkins method that has a separate quote of AR and MA polynomials in the extra term.

### 3.4.1 Model order selection and parameter estimation

Modelling a real sequence of measurements as an ARMA process, one of the challenges is to find good model orders. Two tools to help with this are the Auto-Correlation Function (ACF) and the Partial Auto-Correlation Function (PACF).

As the correlation of any two variables in a stationary stochastic process only depends on the distance between them, it is possible to define the auto-correlation function as

$$r_Y(\tau) = \rho(Y_t, Y_{t-\tau}). \quad (22)$$

The ACF is a symmetric function and  $r(0) = 1$  always holds as any variable is totally correlated with itself. For a white noise process the ACF will be zero for all other lags. For an ideal MA-process of order  $q$ , the ACF will be non-zero for all lags up to  $q$  and then zero. This makes the ACF a very useful tool to identify a process as an MA-process of a certain order.

For an AR-process the ACF will decay as a damped exponential or sine function, regardless of model order. It would be useful to have another tool that helped model order selection of an AR-process in the same way the ACF does for MA-processes. The Partial Auto-correlation (PACF) is defined for this purpose. This is a function that for a given process, for every lag  $k$ , gives the estimated value of the last coefficient in an AR representation of this process. For example given an ideal AR(2)-process, the first two values will be non-zero, but when modelling this as an AR(3)-process the third coefficient will be estimated to zero. For an MA-process the PACF will decay as a damped exponential or sine function.

For an ARMA( $p,q$ )-process, both the ACF and the PACF will decay as a damped exponential or sine function after lag  $|p - q|$ .

After deciding on a model type and order, the coefficients of the AR- and MA-polynomials need to be estimated to fit the data. One option is to use the least squares method, but in the general case this is not the best method. Another option is the Maximum Likelihood estimation. For this method, the Likelihood function,  $L(\Theta)$ , where  $(\Theta)$  is a vector containing all the coefficients, is defined as

$$L(\Theta) = f_{\mathbf{X}}(\mathbf{x}; \Theta). \quad (23)$$

Here,  $\mathbf{x}$  is the observed realization.  $f_{\mathbf{X}(\mathbf{x};\Theta)}$  is the probability of outcome  $\mathbf{x}$ , given  $\Theta$ . The  $\Theta$  that maximizes the Likelihood function contains the combination of parameter values that are most probable to result in the observed realization. In practical situations the probability density function is usually not known, but there are various algorithms designed to approximately find the Maximum Likelihood estimation.

## 4 Earlier Work

So far little work has been done on modelling of individual buildings in a district heating system. However, many papers have been written on modelling of the entire system for forecasting of heat loads in order to optimize plant operation.

An advanced model of the entire district heating system was presented in a doctoral thesis at LTH by Arvastsson (2001). He used a grey-box approach, meaning the model included both physical and stochastic components. The model was developed to work both for prediction and simulation. Many parts of the network and production facilities were modelled on a detailed level, while the consumer buildings were given less emphasis.

Dotzauer (2002) claimed that Arvastssons models were good for simulation but argued that when predicting loads in a real system, uncertainties in weather forecasts and lack of well measured data would limit the usefulness of such a complex model. He suggested a simpler approach where the heat load was modelled as one temperature-dependent part and one social part. The temperature dependence was subtracted based on a rough energy signature and the residual was expected to follow set historical weekly patterns. The model was not very accurate, but he argued that it was of greater significance to improve weather forecast predictions than to make a more precise model.

For some time the debate seemed settled, but as more detailed data became available and computer performance improved over the next decade more advanced models were again suggested. Madsen and Aalborg Nielsen (2006), like Arvastsson, use a grey-box model. Unlike Arvastsson they focus on properties of buildings and include relations from heat transfer theory (ventilation, solar radiation etc).

Madsen and Aalborg Nielsen argue that it would be hard to find a good model only using stochastic components. Some others have tried though, and claimed to have managed well. Grosswindhager (2011) uses the same method as Dotzauer to subtract the temperature dependent part using the energy signature, and then models the residual as a SARIMA-process. Chamchov (2010) uses the more advanced Box-Jenkins method.

Fang and Lahdelma (2016) developed and compared several models for the district heating network in Espoo, Finland. They used models similar to Dotzauer but included effects from wind speed to complement the temperature and improved the social pattern by accounting for midweek holidays. They then compared these to SARIMA models, and concluded that on the system level the first type of models performed slightly better.

Gadd and Werner (2013a, 2013b) defined quantities describing the magnitude of daily and seasonal changes and used these to examine the conditions for central heat storages in order to eliminate daily variations from the system.

All mentioned studies have focused on the entire system from the plant perspective, but several of them have commented on the possibility to study individual substations. Arvastsson writes (2001, p.3) "Detailed models of all consumers connected to a district heating network are hardly meaningful and not possible without measurements close to the different consumers." Grosswindhager (2011) writes "Consumer load forecasting were not treated, due to the highly stochastic nature of the consumer data, which would make it necessary to build several individual models." As we now have measurements available from substations Arvastsson's comment is outdated. Grosswindhagers objection remains to be proved or disproved by scientific studies.

Gadd and Werner (2013a) also discuss the difficulties in modelling individual substations stating that "A large variation of heat load patterns among various buildings implies that a standard heat load pattern for customer substations does not exist."

Modelling has been done on the heat demand for individual buildings without a strict connection to district heating. One such model was made by Bacher, Madsen, Aalborg Nielsen and Perers (2013). They created a model for single-family houses, based on Madsen and Aalborg Niensens earlier work on district heating systems. A similar grey-box model is used with focus on heat dynamics of buildings, where the building temperature is seen as low-pass filtered outdoor temperature. Some harmonic curves are fitted to the model and then the remaining noise is modelled as an AR process. Here, the hot water preparation is filtered away so that only space heating is considered. The model works for all buildings in the study (with individual parameter fitting), but it is suggested to study social behaviour patterns and solar radiation effects further to improve the model.

Zhou, Wang, Xu and Xiao (2008) made an advanced model and applied it to a 50-floor multipurpose building in Hong Kong. The approach was similar to the study by Bacher, Madsen, Aalborg Nielsen and Perers, but more detailed. They also included air humidity as a factor, which they argued was relevant in this geographical region.

Sandin, Gustafsson and Delsing (2013) studied fault detection in substations. They created linear regression models with the outdoor temperature as input and different variables available in metering data as output and studied outliers and drifts that deviated from these models, also taking account daily and weekly cycles. Their methods were not tested in practice, but similar work done by master student Lindquist (2010) at Chalmers University has been implemented at Göteborgs Energi.



## 5 Method

In this section the methods used to develop the models and answer the research questions are presented. The type of models used is motivated and the available data is described. The modelling procedure is explained, as well as the tests used to evaluate the models performance, general applicability and usefulness in fault detection applications.

### 5.1 Model Creation

In the choice of output variables for the model there were several alternatives. As the aim is not to create a cost forecasting tool it is not necessary to choose power as output variable. However, such a model is interesting as it can be compared to models created for entire systems, and the relationship of power versus outdoor temperature is well studied. Good alternatives are flow and delta-T, as these values say more about the characteristics of the substation. It can be argued that changes in flow values are just consequences of delta-T characteristics and the heat demand. Therefore heat power and delta-T were chosen as the two endogenous variables to study.

Selecting the exogenous variable, one option was to use the outdoor temperature, as has been done in most studies on the system level. Another possibility that exists when studying individual buildings is to use the supply temperature as an exogenous variable, as it is not controlled by the substation and continuously changes both daily and seasonally to meet the system demand. To the authors knowledge, this has not been done before. It was decided to try linear regression against both input variables and also to combine the two.

Looking at the earlier studies there were basically two categories of methods used to model the patterns caused by social behaviour. Either regression on some deterministic pattern based on historical data was used, or ARMA-type models were created. Fang Lahdelma (2016) made a comparison of the two approaches and concluded that on the system level the first category performed slightly better. However the idea is that the developed model should be applicable to a large number of substations and according to Gadd & Werner (2013a), a standard heat load pattern for all substations does not exist. Therefore it was decided to use a model of ARMA-type.

An option was to filter out the hot water use and only study space heating, as done in the study by Bacher et al (2013). This approach was not chosen though, because it was of interest to see what patterns could be modelled also in the water usage. Random noise caused by irregular water usage could indeed raise the magnitude of the model errors, but should not affect the choice of model.

### 5.1.1 Data

Data from the district heating network in Karlshamn, Sweden was provided by NODA Intelligent Systems. Data from five different substations were used in the study. One of them was known to be a well-performing substation and therefore was chosen as the reference substation used in the modelling. For the reference substation, data was available for the years 2015-2017. The others had unknown characteristics and different time periods.

For all substations hourly measurements of outdoor temperature, supply temperature, return temperature, flow, volume (accumulated flow), heat power and accumulated heat energy was available, with exception for some missing values.

To enable using the same model for all substations, regardless of building size, all measurement sequences  $Y$  were normalized to a dimensionless scale between 0 and 1, according to

$$Y_{normalized} = \frac{Y - Y_{min}}{Y_{max} - Y_{min}}. \quad (24)$$

### 5.1.2 Modelling Procedure

For the modelling the data was split into three sets, following established principles of cross-validation. The first set is for modelling, the second for validation and the third for testing. The modelling set is examined to decide on model orders and estimate parameter values. The validation set is used to confirm or dismiss the suggested models. These two are used in an iterative manner to find the best possible match. The idea behind this approach is to avoid over-modelling, i.e. that the parameters are unintentionally fitted to describe the unpredictable random noise in the specific sample.

The third set is put aside and not involved at all in the creation of the model. This is only used in the end to see how well the model actually performs when it is done. As data was available for three years it was intuitive to split the data by calendar year.

This study started from the standard case where the heat power was modelled with the outdoor temperature as the only input variable, like several earlier studies had done on a system level. The program code to evaluate this case was written in a general format so that it could then easily be modified to run with other combinations of input and output variables. The following procedure was used to find an optimal model.

1. Prepare data by deleting obvious outliers and fill gaps.

2. Perform linear regression on modelling data. Try modifying the exogenous data and see if regression results improve.
3. Study the residual and its corresponding ACF and PACF to find trends and cycles.
4. De-trend modelling data by differentiation.
5. Study the residual and its corresponding ACF and PACF to find appropriate orders for the AR and MA polynomials. Make an educated guess.
6. Apply the model to the validation data and let the computer estimate all coefficients together by maximum likelihood estimation. This includes re-estimating the regression parameter. Study the ACF and PACF of the residual to examine if there is any remaining structure.
7. Repeat steps 5 and 6 until the model has been identified which has the least remaining structure. Ideally the residual is pure white noise.

For this work Python was used, with standard packages as well as the Pandas and Statsmodels toolboxes.

## 5.2 Model Evaluation

To evaluate the performance of the different models that had been created for different combinations of input and output variables, they were applied to the test set. The models were compared by their fit to the measured data as well as the whiteness (lack of structure) and variance of the residuals. The best performing models were applied to data from the other four substations to test their general applicability. It was tried both to apply the models directly and to re-estimate the coefficients.

To investigate if the models could be useful in fault detection applications, two tests were made to cover both slow changes in the substation characteristics and a sudden fault. To study any slow changes, the model coefficients were repeatedly estimated using data from different time periods within the available data set. This was done both for the reference substation and for one of the other substations, that had data available for a longer time period.

To study the effects of a sudden fault the model was applied to a data set containing such a fault and the deviation from the model predictions were studied. Since no data was available from substations with registered faults, a typical fault was simulated on the reference data set. This was done by suddenly raising the return temperature by 10 percent (and adjust delta-T accordingly) at one moment in time and then apply a selected model to this new data set. Apart from analyzing the residual it was also investigated if the fault made any impact on the model parameters, as they were re-estimated based on the faulty data set.



## 6 Modelling Results

In this section it is described how the models were created, based on the method motivated in the previous section. Six different models were made with different combinations of endogenous and exogenous variables. The creation of the first model is described in detail, while the others are described with focus on similarities and differences to the first.

### 6.1 Modelling Heat Power

Using heat power as the endogenous variable, three models were developed with the outdoor temperature, the supply temperature and both as exogenous variables. Figure 3 shows the heat power sequence to be modelled. As explained, the sequence was split in three by calendar year for modelling, validation and testing.

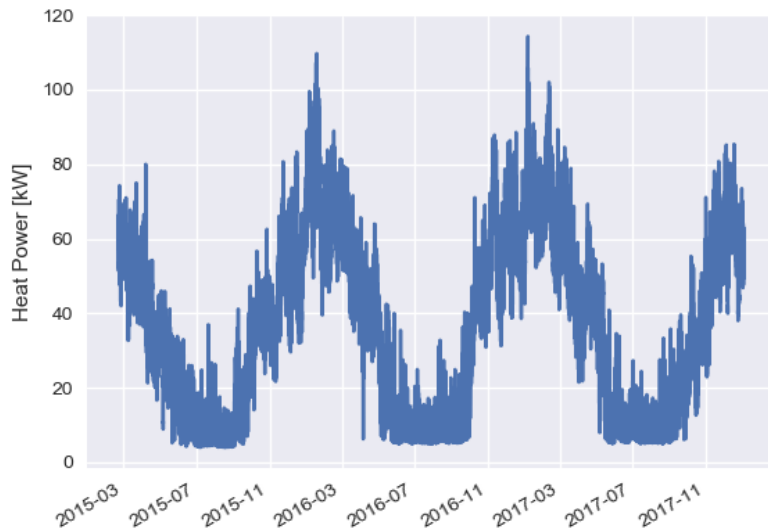


Figure 3: Hourly Heat Power values over the whole data set available for the reference substation.

For linear regression to be a useful tool there needs to be a close to linear dependency between the endogenous and exogenous variables. The plot in figure 4, top left, reveals that there exists such a relationship between the outdoor temperature and the power, at least for temperatures up to some breakpoint around 17 °C. Also between supply temperature and power (same figure, top right) the values are clustered around a straight line across the diagonal, although a more complicated structure is also visible. Even though the relationship might not be as clear for hourly values as for the daily averages plotted in the figure, it was sensible to try both inputs.

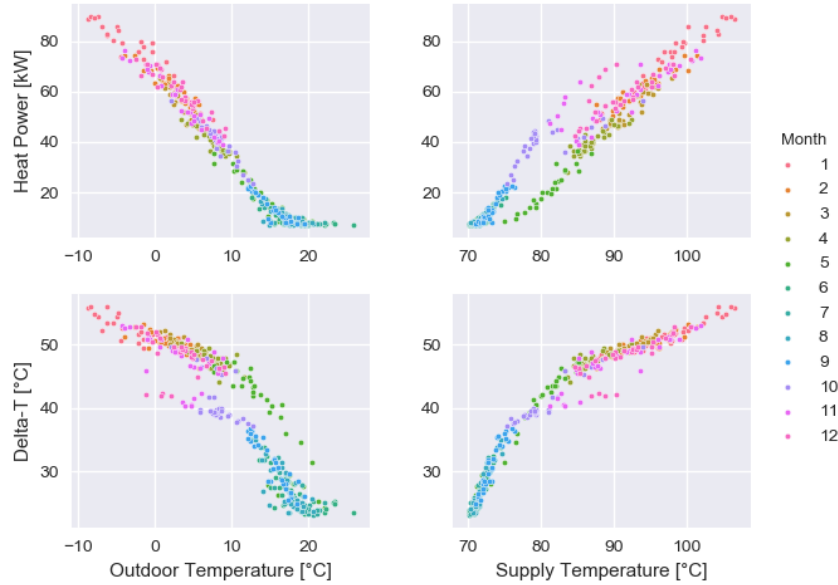


Figure 4: Daily average values of the four quantities used in the modelling. The endogenous variables are on the x axis and the exogenous variables on the y-axis, to visualize their mutual dependencies.

In the first model created, the outdoor temperature was used as exogenous variable. To enable a good linear regression fit the temperature data was modified by setting all temperatures greater than 17 °C equal to 17.

Outdoor temperature and heat power follow each other over the diurnal cycle, while they have an opposite relationship over longer time spans. Thereby it was motivated to try the average over the last 24 hours instead of the instantaneous value. Linear regression models of the heat power were made with one and both of these modifications on the outdoor temperature, and the residual variances were compared. It was concluded that both modifications to the input data increased model performance.

The part of the process predicted by linear regression was subtracted and the residual was studied to decide whether it was stationary or if further transforms, like differencing, would be necessary. The mean and variance weren't perfectly stationary, but it was considerably close. Because of this borderline case models both with and without the difference operator were tried, but it turned out during the process that the process wasn't stationary enough to be modelled as an ARMA-process if the difference operator was omitted.

Seasonal trends were identified by looking at the ACF and PACF in figure 5. As expected a 24 hour cycle is clearly present. A small peak at 168 hours was also seen in the PACF, when plotted for larger lags, suggesting that there could also be some weekly cycle. However it was very weak compared to the 24 hour pattern and therefore neglected.

By filtering the process by convolution with both a standard and a 24 hour differ-

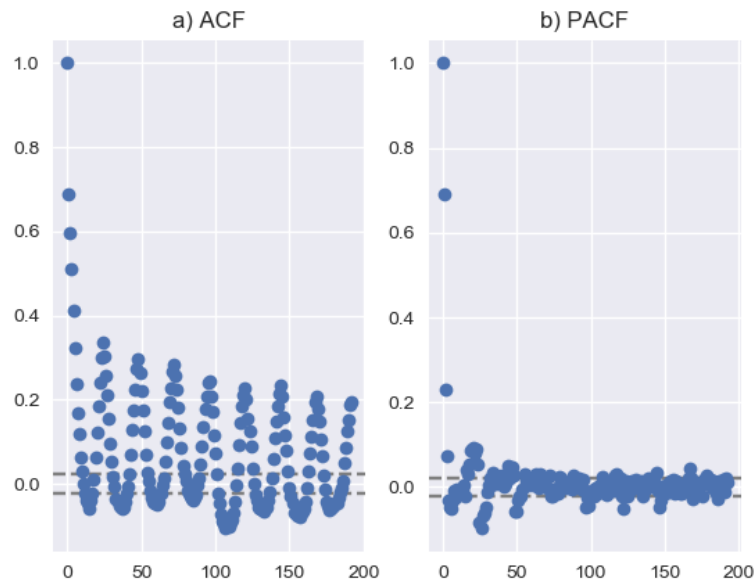


Figure 5: The estimated Auto-Correlation Function (a) and Partial Auto-Correlation Function (b) of the heat power process, after regression against the modified outdoor temperature.

ence operator the trends were removed, and another structure was revealed. The ACF and PACF of the remaining process is seen in figure 6. The single spike in the ACF together with the exponentially decaying PACF agree very well with an MA(1)-process, so it was decided to include such a parameter in the model.

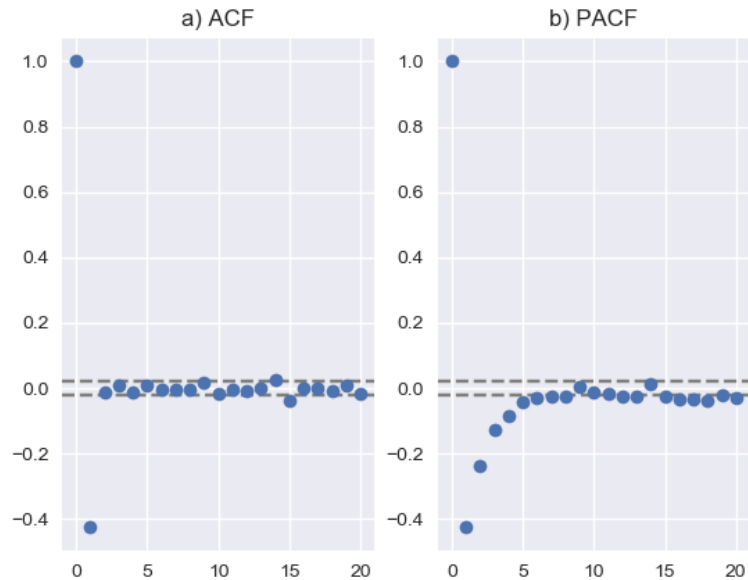


Figure 6: The estimated Auto-Correlation Function (a) and Partial Auto-Correlation Function (b) of the heat power process, after applying 1 hour and 24 hour differencings.

The process was filtered through this SARIMA  $(0, 1, 1) \times (0, 1, 0)_{24}$  model and the

ACF and PACF of the remaining process were studied over a larger number of lags to find if any seasonal AR or MA parameters needed to be included. As seen in figure 7 there is an almost perfect MA(1)-process on multiples of 24, as there is a single spike at lag 24 in the ACF and exponentially decaying spikes on multiples of 24 in the PACF. Hence, a SARIMA  $(0, 1, 1) \times (0, 1, 1)_{24}$  model seemed like a good choice. This was tried on the modelling data with a pleasing result and then confirmed on the validation data. The ACF and PACF of the residual did not agree completely with those of a white noise process, but it appeared to be the closest possible. Some experimenting was made with including higher order parameters, but as expected these coefficients were estimated close to zero and did not improve the model performance when applied to the validation data.

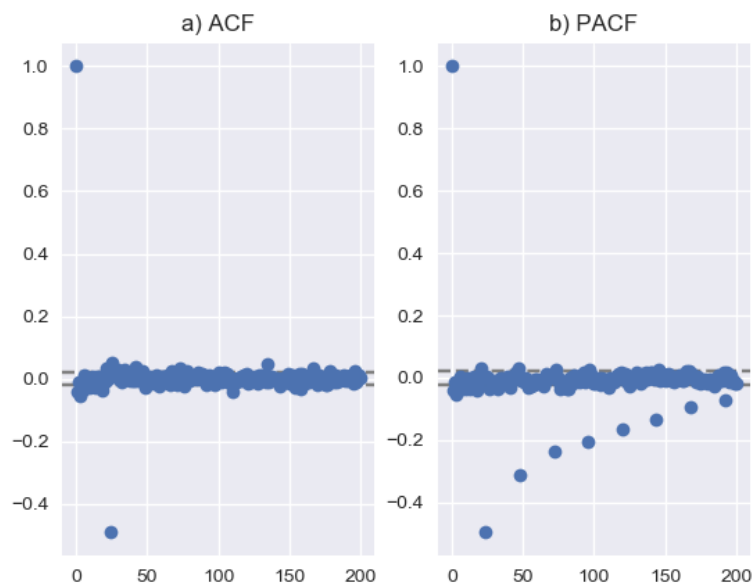


Figure 7: The estimated Auto-Correlation Function (a) and Partial Auto-Correlation Function (b) of the heat power process, after applying a SARIMA  $(0, 1, 1) \times (0, 1, 0)_{24}$  filter.



A second model was made using the supply temperature as exogenous variable. As before it was tried to use the 24 hour mean as regressor, but the best regression results were found using the original sequence. The regression to the exogenous variable was again not enough to make the process stationary, so in order to proceed the process needed to be differentiated. Neither was the diurnal cycle in the supply temperature enough to eliminate the 24 hour season in the data, so a 24 hours differentiation was also applied. The remaining parameters were identified using the ACF and the PACF, again iterating between modelling data and validation data. It was found that the structure was nearly identical to the previous case, so also here it was decided to use a SARIMA  $(0, 1, 1) \times (0, 1, 1)_{24}$  process.

A third model was tried where both exogenous variables were used simultaneously. Just like the first case it was hard to see whether a differencing was needed, so both variants were tried. Also here it was found that a SARIMA  $(0, 1, 1) \times (0, 1, 1)_{24}$  was the best possible fit.

## 6.2 Modelling Delta-T

It was also of interest to develop models with delta-T as the endogenous variable. As seen in figure 4, bottom line, neither the outdoor or supply temperature has a direct linear relationship to delta-T. But studying the plots it is possible to imagine a line across the diagonal that fits the data to some extent, for both input choices. It should be noted however, as delta-T is defined as the difference between supply and return temperature, that a model of delta-T with the supply temperature as input is in practice the same as a model of the return temperature alone.

At first, the outdoor temperature was used as input. Comparing results from regression tests it was again decided to transform the outdoor temperature to its 24 hour mean, before applying it to the model. Again both 1 lag and 24 lag difference operators were applied to compensate for trends and cycles, and yet again the residual turned out to be best modelled as an MA process with  $q = Q = 1$ . All parameters were estimated forming another SARIMA  $(0, 1, 1) \times (0, 1, 1)_{24}$  process.

Also in the case when supply temperature was used as exogenous variable and when both variables were used together the SARIMA  $(0, 1, 1) \times (0, 1, 1)_{24}$  model turned out to be the best fitting structure.

### 6.3 Summary of Models

Six different models were created, all using a SARIMA  $(0, 1, 1) \times (0, 1, 1)_{24}$  structure combined with linear regression against one or two exogenous variables. They will hereafter be referred to by a capital letter each, according to table 1 that presents a summary of the model coefficients. When the outdoor temperature was used as exogenous input, it was first modified to contain the average over the last 24 hours instead of the instantaneous value. For model A it was also modified according to the energy signature by setting all values greater than 17 equal to 17, before normalizing the data.

Table 1: Summary of the six models created with different combinations of input and output variables. The different parameters are: Regression parameter vs. outdoor temperature, Regression parameter vs. supply temperature, MA-parameter for lag 1 and seasonal MA-parameter for lag 24.

Model	Variable	Outdoor T.	Supply T.	MA.L1	MA.S.L24
A	Power	-0.5001		-0.5459	-0.8921
B	Power		0.3074	-0.6157	-0.9044
C	Power	-0.2956	0.2883	-0.6292	-0.9003
D	delta-T	-0.5464		-0.5022	-0.9066
E	delta-T		0.9559	-0.7055	-0.8704
F	delta-T	0.2375	0.9713	-0.7106	-0.8688

## 7 Model Performance

In this section the results from the tests on model performance are presented, as all six models were tested on the reserved data set. Plots of 1-step predictions vs real data are shown and observations from the residual analysis are presented. The models with the best performance were tested on other buildings and these results are also presented here.

### 7.1 Heat Power Models

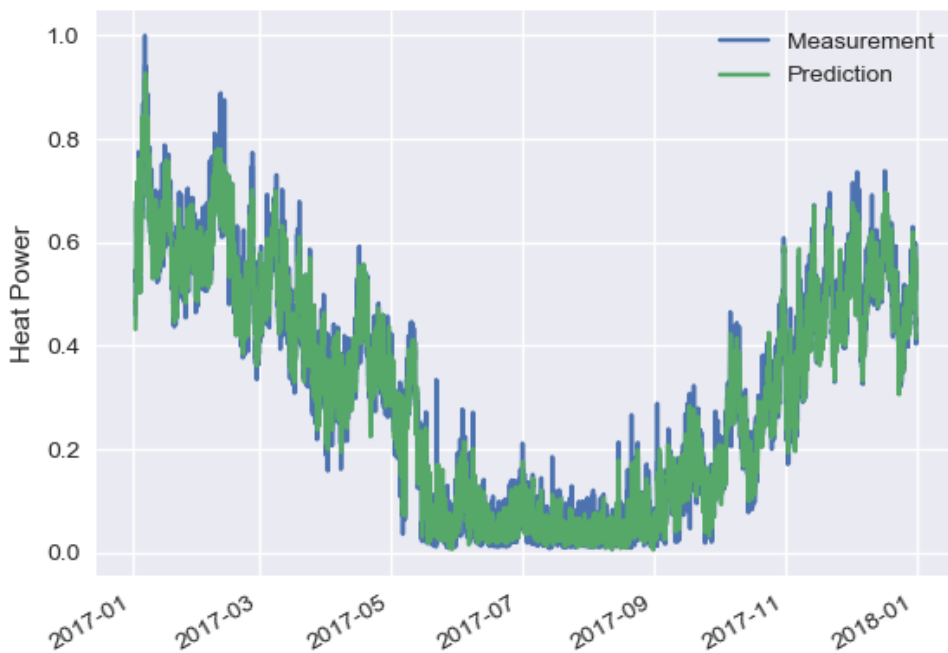


Figure 8: 1-step predictions of the Heat Power (normalized) based on model A, compared to the actual data in the test set.

Models A, B and C were all applied to the test data from year 2017 and the 1-step prediction as well as properties of the residuals were studied. Figure 8 shows the prediction on top of the actual data for model A. All models fitted good enough to look good on this scale, so a closer look of a typical week in early April was also studied. See figure 9. Even from this figure it is hard to draw any conclusions of the performance of the respective model, other than that they all generate similar results. Overall, the difference between the modelled data and the real data was larger than the difference between the different models.

In figure 10 some properties of the model residual for model A are shown. Subfigure a) shows the model residual and its mean value. Subfigure b) shows a cumulative



Figure 9: 1-step predictions of the Heat Power (normalized) based on model A, B and C for a typical week in the test set, compared to the actual data.

sum, meaning that the errors are sequentially added to each other. Subfigure c and d shows the estimated ACF and PACF of the residual. Subfigure e) shows how the variance of the residual varies over time, calculated from a rolling window of 168 hours (1 week). Subfigure f) shows the distribution of the errors, compared to an ideal Gaussian distribution with the same mean and variance. If the residual was pure white noise the distribution would fit under the red curve. The residuals of models A, B and C all show similar characteristics. The mean and cumulative sum is stable, supporting that all trends have been properly taken care of. The ACF and PACF, distribution plot and rolling variance are not exactly those of a white noise process, which would have been the ideal result, but they are not very far off.

In the cumulative plot it is seen that it takes a few weeks for the residual to reach a steady level and after that it fluctuates around an equilibrium. The time it took for the residual to stabilize was examined to compare the models. In Figure 11 the first two weeks of the cumulative sums have been zoomed in. It is shown that the residual of model B takes longer to stabilize than the other two. After 14 days all three models have stabilized. The error variances were 0.0010 for all models. In the calculation of these, the first 14 days have been omitted.

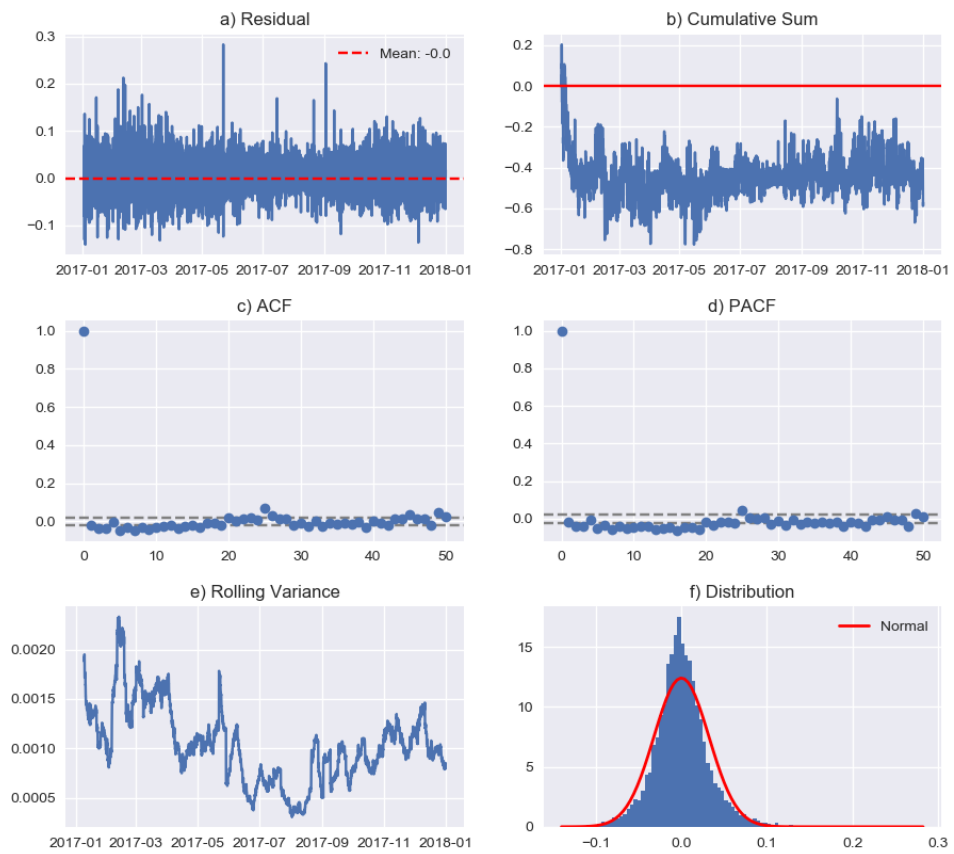


Figure 10: Properties of the model residual for model A. The subplots are described in the text.



Figure 11: The cumulative error for the three models A,B and C over the first 14 days of the test set.

## 7.2 Delta-T Models

The same tools were used to evaluate the performance of models D, E and F. Prediction results are presented in figure 12 and 13. Like the earlier case, all models managed to fit the data about equally well, in a large scale perspective.

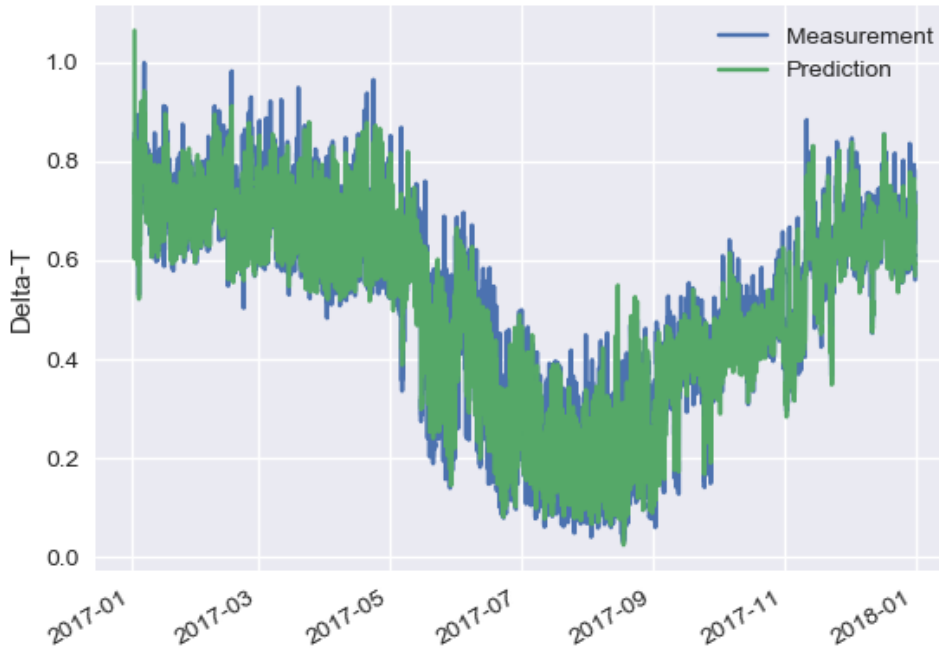


Figure 12: 1-step predictions of delta-T (normalized) based on model D, compared to the actual data in the test set.

The residual properties for model D are presented in figure 14. All three models show similar characteristics. After stabilizing, the errors are evenly distributed and quite close to fit the normal distribution profile. The ACFs and PACFs are significantly different from zero for most lags, but follow a flat profile only slightly off-zero. For all three models, both magnitude of the residuals and their variance is clearly larger during the summer months than the rest of the year. A closer view of the cumulative sum of errors the first 14 days is presented in figure 15. They all seem to stabilize on a steady level after around 10 days. The error variance was 0.0029 for model D and 0.0024 for the other two.

Since the error variance of model A, B and C were of equal magnitude the decisive property was the time it took for the models to stabilize. In fact this was the only part of the test where it made any clear difference what exogenous variable was used for the regression. Accordingly, model A was appointed as the best choice.

As for model D, E and F the first thing that set them apart was that model D had a little larger error variance, meaning it's 1-step predictions were slightly

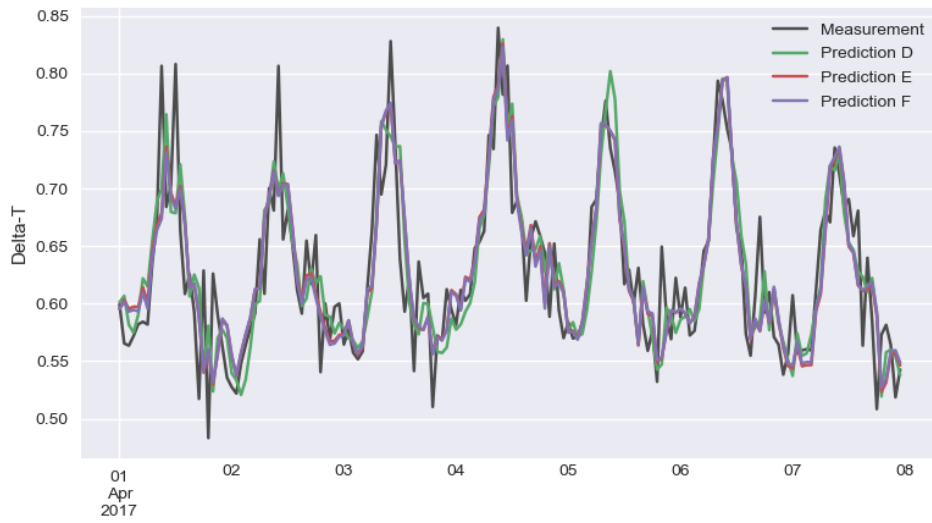


Figure 13: 1-step predictions of delta-T (normalized) based on model D, E and F for a typical week in the test set, compared to the actual data.

less accurate. Model E was appointed as the best choice as it was the fastest to stabilize. Model F also performed well, but not well enough to motivate the use of an extra parameter. Furthermore it seemed unreliable to use a model where the regression coefficient against the outdoor temperature was positive. Model E was kept together with model A for the remaining studies.



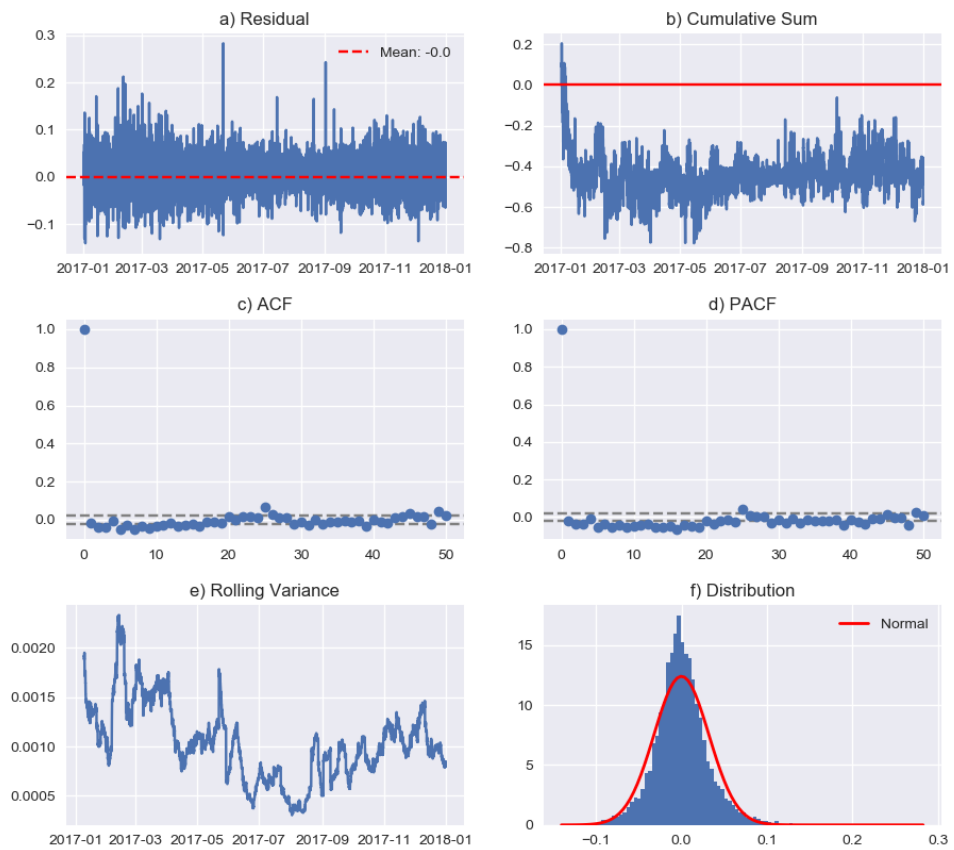


Figure 14: Properties of the model residual for model D. The subplots are described in the text.

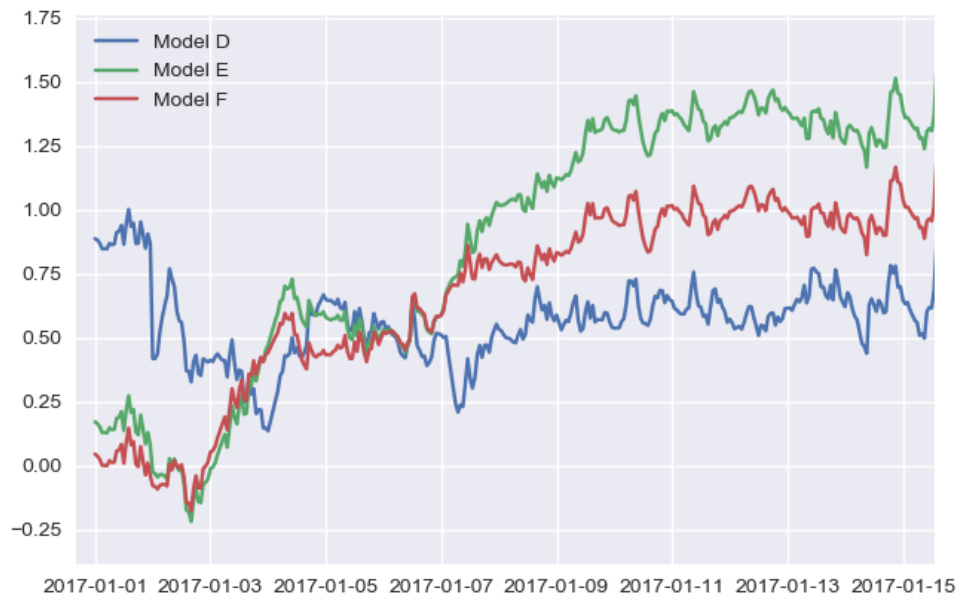


Figure 15: The cumulative error for the three models D, E and F over the first 14 days of the test set.

### 7.3 Other Buildings

The four other buildings were tried with method A and E as these were the best models accounting for error variance, stability and simplicity. It was tried both to apply the models directly with the coefficients estimated for the reference building and to re-estimate the coefficients to fit respective building.

First, the temperature breakpoint for each building's energy signature was identified by visual inspection of the heat signature plot. While building 1 had a breakpoint at 17, like the reference building, the others seemed to have a linear temperature dependence all the way up to 21-22 degrees. Model A was applied to data from one year, and the 1 step predictions were compared to measured data. For all buildings the prediction fitted rather well during the heating season, but rather poorly during the summer. In other words, changing the breakpoint of the heat signature was not enough to make the model fit all buildings. The error variance for each case is listed in table 2, together with re-estimated coefficient values to fit the data from each building and the variance from this test. The numbers show that the overall error variance was reduced only marginally by re-estimating the model coefficients, and for all buildings except number 2 it is at least double that of the reference case. It is also shown that the regression against the outdoor temperature is weighted lower than for the reference case, especially for building 2 and 3.

Table 2: Results from model A applied to the five different substation. The residual variance is given for the direct application of the model estimated for the reference substation. Both the variance and the parameter values are given for every substation as the parameters were re-estimated for the respective data set.

Building	Variance	Var. Re-fit	Outdoor T.	MA.L1	MA.S.L24
Ref.	0.0010		-0.50	-0.55	-0.89
#1	0.0023	0.0022	-0.39	-0.64	-0.88
#2	0.0014	0.0013	-0.17	-0.38	-0.87
#3	0.0024	0.0023	-0.12	-0.48	-0.82
#4	0.0024	0.0023	-0.34	-0.38	-0.89

The same tests were carried out for model E, and the results are presented in table 3. In this case the original model fitted the other buildings about as good as the reference building, based on the error variances. The improvements when re-fitting the model to each building is again only marginal. Also here building two stands out as it gives a much lower variance than the others. The coefficient values are similar for all models, except for the lag 1 MA coefficient for the reference building and the regression against supply temperature for building 4.

Table 3: Results from model E applied to the five different substation. The residual variance is given for the direct application of the model estimated for the reference substation. Both the variance and the parameter values are given for every substation as the parameters were re-estimated for the respective data set.

Building	Variance	Var. Re-fit	Supply T.	MA.L1	MA.S.L24
Ref.	0.0024		0.96	-0.71	-0.87
#1	0.0029	0.0028	1.06	-0.59	-0.93
#2	0.0010	0.0010	0.92	-0.52	-0.89
#3	0.0024	0.0022	1.00	-0.45	-0.89
#4	0.0019	0.0018	1.27	-0.49	-0.92

## 8 Detecting Drifts and Faults

To further explore possible applications of the models and answer the last of the questions posed in this thesis, it was studied if the selected models could be used to observe differences between a well performing installation and a faulty or poorly performing one. Both slow and sudden changes were studied. The results are presented in this section.

### 8.1 Detection of Slow Changes in Performance

Given the full data set of the well performing substation, the model coefficients were studied as they were repeatedly re-estimated for all 12 month periods available, with one sample per quarter. The parameters for model E were estimated almost identical, independent of the time period, with the regression parameter against supply temperature being virtually 1 all the time. This was expected as a change of one unit in supply temperature will always lead to a one unit change in delta-T. For model A the regression parameter and the lag 1 MA parameter shifted a little over time, in opposing directions so that the sum of all parameters was held fairly constant.

For one of the other four buildings (building 1), with unknown performance rating, there was data available for a longer time period. The test was repeated on that building. The result is presented in figure 16 for model A and figure 17 for model E. For both models there is a slow drift in all parameters.

All in all it seems like the model is more stable over time for the substation that is known to be well performing than the substation for which the performance is undocumented. However the results are not as clear as they may seem, since there was missing data for a year's time (and another year until the next calculated sample). In the plots the interpolation over the gap looks like a perfect slow drift, which might not be true at all. Although it looks like the trend continues on both sides of the gap. Notably, it seems like the parameters actually drifts towards better performance, closer to the characteristics of the well performing substation.

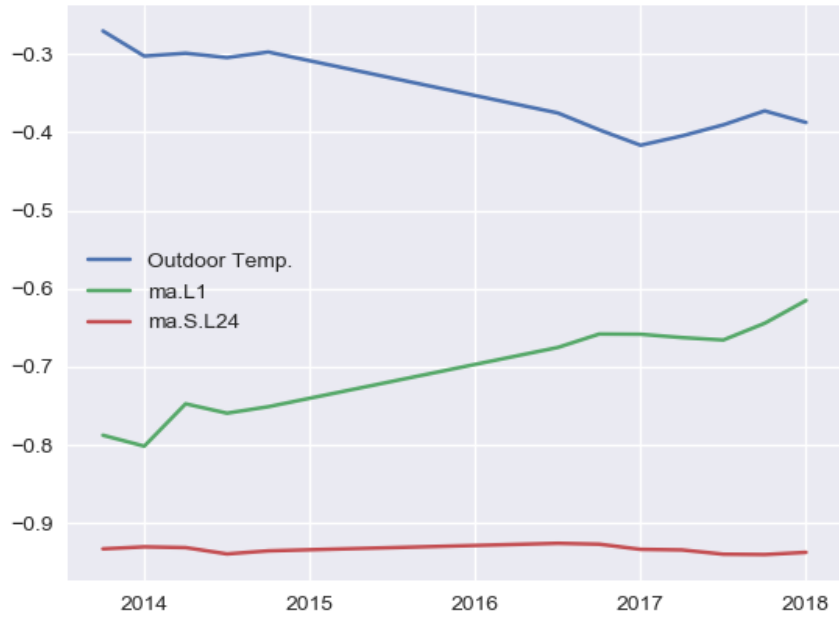


Figure 16: The parameter values for model A applied to building 1 as they are estimated over different yearly windows, with 1 sample per quarter.

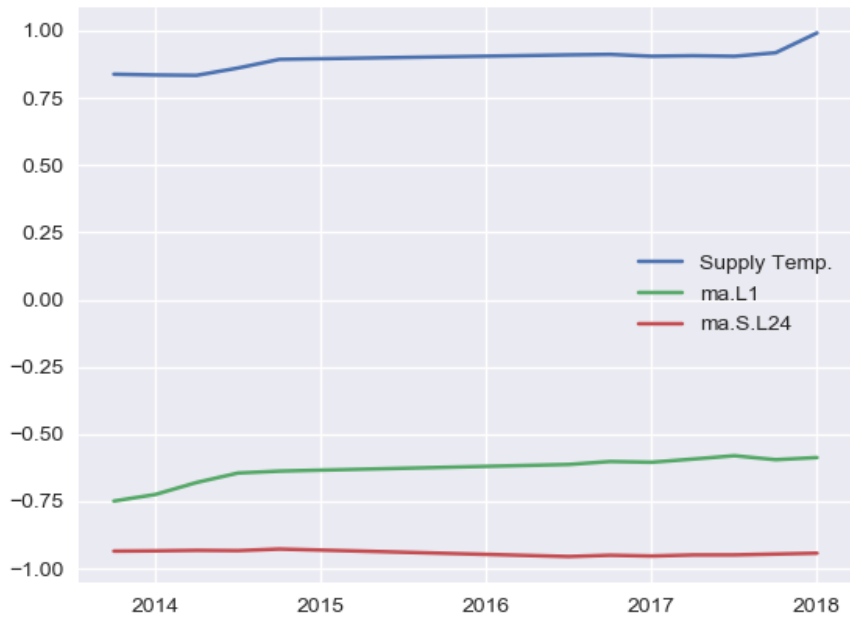


Figure 17: The parameter values for model E applied to building 1 as they are estimated over different yearly windows, with 1 sample per quarter.

## 8.2 Detection of Sudden Faults

To examine how the model dealt with a more sudden kind of fault, such a fault was simulated at January 1st 2017 on the reference substation. The fault raised the return temperature so that  $\Delta T$  decreased by 10 % for all hours after that. The process was filtered through model E and the residual was examined to see if there was any sign of the fault. Figure 18 shows the residual itself and a cumulative sum of its values. The data itself is basically unchanged, but in the cumulative sum a dip can be seen when the error strikes. The dip is large enough to stand out from the general noise. However, when it was tried to add an error of the same magnitude during the summer, when the residual is noisier, it was hard to distinguish the error solely by inspecting the plots.

It was also examined whether such a sudden fault led to any drift in the parameters when these were re-estimated for the faulty installation, but that was not the case.

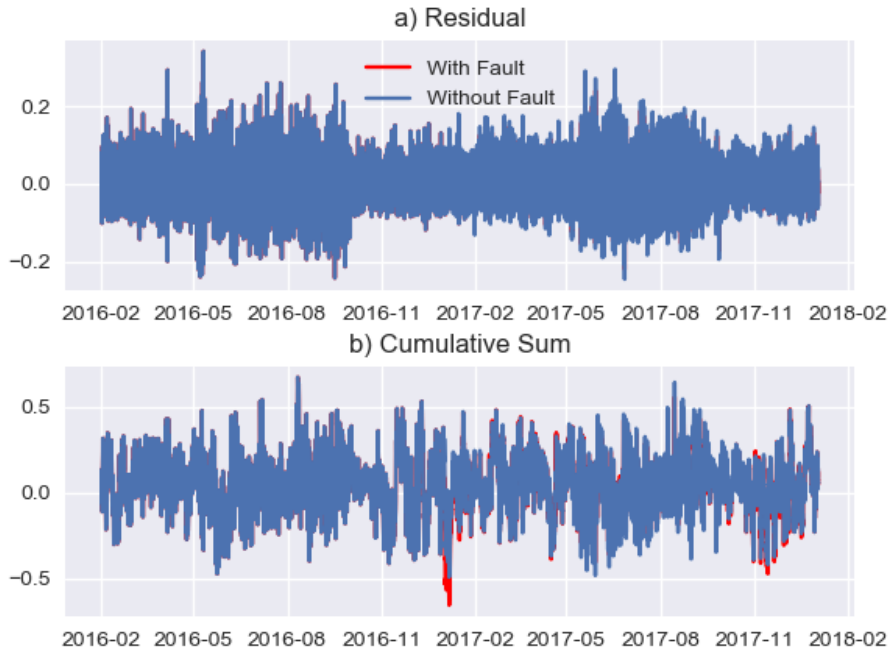


Figure 18: Residual of model E applied to the reference substation with and without a simulated fault induced at January 1, 2017. Subplot a) shows the residual itself and subplot b) a cumulative sum of the residual.





## 9 Discussion

In this section the results are analyzed in relation to the theory. In the light of the analysis the methods used in the work are discussed. Also, some suggestions for future studies are presented.

### 9.1 Model Structure and Performance

If any of the suggested models had had a residual that was ideal white noise, then there would be no difference between the model and reality. For a model with limited physical support and only three adjustable parameters such results would have been truly remarkable. The resulting residuals for all combinations of input and output variables are close enough to white noise to support that the best model orders possible were chosen for the SARIMA model, especially given that the ACF and PACFs used for the development agreed very well with theory. However the results are disappointing in the way that the optimal model turned out to be a pure moving average model in both the local and seasonal polynomials. By nature, such a model is very flexible and compensates for any deviations from the structure. This means that apart from the already established fact that there is a 24 hour cycle in the data, the models don't reveal much about the underlying process, such as time constants or causal relationships. On the other hand, the flexibility of the model allows it to adapt to the changes in daily variation pattern between seasons. For the reference model it even fitted well during summer, when heating is only used for hot water preparation.

It is not surprising that the same model structure was optimal for the heat power and delta-T models, as delta-T is used in the calculation of the heat energy. However, it was unexpected that the same model structure was found to be optimal for all combinations of exogenous variables, that the parameters and performance was so similar, and that so little weight were given to the exogenous regression parameters. These are reasons to believe that a model without any exogenous regression at all would perform about as well. The models were not compared on their performance as forecasting tools more than one hour ahead of known data since that was not the aim of the study. That could however be a way of distinguishing the models' performance as longer forecasts would converge towards the exogenous regression, or zero if no exogenous variable was used. The comparison on how fast the models are to stabilize at the beginning of the sample tell something about this, but is not sufficient as empirical evidence.

## 9.2 Applicability to Other Buildings

Mixed results were obtained when the models were applied to other buildings. It was pleasing to see that the parameter values were about the same when re-fitted on different buildings. That the delta-T model fitted better than the heat power model is probably not because the behaviour of delta-T is more universal. Rather it is a sign that the regression towards exogenous variables are harder to transfer to other buildings than the MA parameters. A possible source of error is the normalization procedure that could have been done differently. For example the normalization could consider mean and variance of the samples instead of minimum and maximum values. Another option for why the heat power model did not fit the other buildings completely is that they were actually performing sub-optimally and hence had different characteristics. This would not agree with the theory though, as under-performing units rather differ in delta-T patterns than energy signature. It seems that the biggest differences between the buildings happen during summer, where different control strategies are applied. One way to get around this is to simply exclude the summer months from the models. This would probably also increase accuracy for the model in the first place. Possibly this could even uncover a more detailed structure in the process.

So do the results support or oppose Grosswindhagers statement that load forecasting at individual substations need separate models for each customer? Regarding heat power forecasting it depends on how much the above suggested improvements would increase the transferability of the temperature-dependent part of the model. For the stochastic part, as well as for the delta-T model, it seems that the biggest problem is not that each substation needs its own model. The problem is rather that the load pattern from each individual substation is fundamentally so random that the best possible statistical model does not have enough accuracy to be useful in forecasting.

## 9.3 Applicability in Fault Detection

In the study of slow changes in the substation characteristics, some interesting results were obtained, especially when investigating building 1. As discussed, there are a lot of uncertainties regarding the period where data was missing, and the substation performance rather seems to change for better than worse. Nevertheless, it is clear that some change in the characteristics of the substation over time can be discovered by comparing values of the estimated model parameters. To draw any further conclusions it is necessary to study more substations, preferably units that can be cross-validated with documentation of faults, maintenance actions and changes in settings.

The simulated sudden fault could be detected by examining the cumulative plot of the residual. However, if the fault is as clear as the one used in the test, it is probably easier to detect the fault by other methods. Possibly the models

would have been more useful in such applications if the regression parameters had been fixed instead of re-estimated together with the other parameters. It turned out that they were consistently given a lower value compared to the pure regression models used during the modelling process. For a delta-T model with a heavier weight on the exogenous regression to have any acceptable accuracy for model, the regression cannot be purely linear as the relationship between the variables are not. For example a piecewise linear curve with an arbitrary number of breakpoints could have been used and maybe two different fittings should be done for spring and autumn, as it is seen in figure 4 that the pattern is different during different seasons. Possibly such improvements to the heat signature could also have improved the heat power model. If the heat signature had only been a little more accurate the regression would be enough to make the process stationary and thus a differencing of the process would be unnecessary. There is no reason though to believe that such an improvement would have revealed any other structure in the variations than the already discovered MA-process.

## 9.4 Suggestions for Future Work

Apart from the improvements and extensions suggested to the methods of this study, there are a number of other questions related to the topics discussed here that would be recommended for future studies. As it seems impossible to create any detailed model given only the primary side data, it would be interesting to include measurements also from the secondary system, especially indoor temperatures. Another suggestion is to study how the cross-correlation between delta-T and flow varies over time. Such a study would act as a good complement to this one to get a fuller picture of the mathematical structure behind the data. Something else that could be of interest for fault detection applications would be to estimate model parameters for a large group of individual substations within the same network and see if any substations show different characteristics. Slow changes and sudden faults could then be detected by monitoring relations between measurement series from the different units. There is also a lot to study using methods of machine learning and artificial intelligence.

Starting to dig into the subject, it is quickly revealed that even though the technology of district heating has been around for a long time, there is still a lot to discover. The next generation of district heating systems, optimized by methods from mathematical statistics and computer science, does indeed have the potential of becoming an important piece in the energy puzzle to create efficient and sustainable solutions for the cities of the future. In Sweden and around the world.



## 10 Conclusions

With theoretical knowledge and results from earlier studies in mind, a group of mathematical models were developed describing the relationship between different physical quantities and their progress through time. Models were developed for both heat load and delta-T. It turned out that for any combination of selected variables a linear regression combined with a SARIMA  $(0, 1, 1) \times (0, 1, 1)_{24}$  was the best fitted model. One model for each variable was designated as best performing and evaluated by the research questions posed in the introduction.

It was confirmed that the heat load patterns of individual substations are very random in nature and therefore it is not possible to create a model with many parameters. Hence the developed models had a wide applicability but a low detail. The residuals were not perfect white noise, as would have been the case if all predictable effects were captured by the model, but overall the predictions matched the measured data rather well.

Testing the models' applicability to other buildings, the delta-T model fitted just as good as for the reference building. The heat power model had acceptable performance, but did not fit as good as for the reference building. Several suggestions for improvements were suggested to make also the heat power model applicable to other buildings.

A slow change in the performance of a substation could be detected by observing changes in the model parameters over time. However, this needs to be tested empirically and further developed before it could be used in practical applications. A sudden fault simulated on one of the substations could be detected as a deviation from the delta-T model if it was large enough and happened to strike at a favourable time. Suggestions for modifications to the model to perform better at this were suggested.

In summary, the work presented in this thesis gives new insight on the statistical properties of data from individual substations and what can be done using the methods that were used. Hopefully these insights will be helpful to other research projects within the field.



## 11 References

Aalborg Nielsen, H. & Madsen, H. (2006). Modelling the heat consumption in district heating systems using a grey-box approach. *Energy and Buildings*, 38(1), pp. 63-71. /doi.org/10.1016/j.enbuild.2005.05.002.

Arvastson, L. (2001). *Stochastic modelling and Operational Optimization in District Heating Systems*. Diss. Lund: Lund University. <http://lup.lub.lu.se/record/20194>

Bacher, P., Madsen, H., Aalborg Nielsen, H. & Perers, B. (2013). Short-term heat load forecasting for single family houses. *Energy and Buildings*, 65, pp. 101-112. doi.org/10.1016/j.enbuild.2013.04.022.

Blom, G., Enger, J., Englund, G., Grandell, J. Holst, L. (2005). *Sannolikhetsteori och Statistikteori med Tillämpningar*. 5. ed. Lund: Studentlitteratur. ISBN: 9144024428

Chamchov, B. (2010). Heat demand forecasting for concrete district heating system. *International Journal of Mathematical Models and Methods in Applied Sciences*, 4(4), pp. 231-239.

Dotzauer, E. (2002). Simple model for prediction of loads in district-heating systems. *Applied Energy*, 73(3-4), pp. 277-284. doi.org/10.1016/S0306-2619(02)00078-8

Fang, T. & Lahdelma, R. (2016). Evaluation of a multiple linear regression model and SARIMA model in forecasting heat demand for district heating system. *Applied Energy*, 179, pp. 544-552. doi.org/10.1016/j.apenergy.2016.06.133.

Frederiksen, S. Werner, S. (2013). *District Heating and Cooling*. Lund: Studentlitteratur. ISBN: 9789144085302.

Gadd, H. & Werner, S. (2013a). Daily heat load variations in Swedish district heating systems. *Applied Energy*, 106, pp. 47-55. doi.org/10.1016/j.apenergy.2013.01.030.

Gadd, H. & Werner, S. (2013b). Heat load patterns in district heating substations. *Applied Energy*, 108, pp. 176-183. doi.org/10.1016/j.apenergy.2013.02.062.

Grosswindhager, S., Voigt, A. & Kozek, M. (2011). Online Short-Term Forecast of System Heat Load in District Heating Networks. *proceedings of the 31st international symposium on forecasting*. Prague, Czech Republic June 26-29 2011.

Jakobsson, A. (2015). *An Introduction to Time Series modelling*. 2. ed. Lund: Studentlitteratur. ISBN: 9789144108360

Lindgren, G., Rootzén, H. & Sandsten, M. (2014). Boca Raton: CRC Press. ISBN: 9781466586185

Lindquist, P. (2010). *Verktyg för utvärdering av energieffektiviseringar, baserat*

*på effektsignaturanalyser*. Master Thesis, Department of Computer Science and Engineering. Chalmers University of Technology.

Sandin, F., Gustafsson, J. & Delsing, J. (2013). *Fault detection with hourly district energy data: Probabilistic methods and heuristics for automated detection and ranking of anomalies* (Rapport 2013:27). Stockholm: Svensk Fjärrvärme.  
<http://urn.kb.se/resolve?urn=urn:nbn:se:ltu:diva-22021>

Zhou, Q., Wang, S., Xu, X. & Xiao, F. (2008). A grey-box model of next-day building thermal load prediction for energy-efficient control. *International Journal of Energy Research*, 32, pp. 1418-1431. doi:10.1002/er.1458



