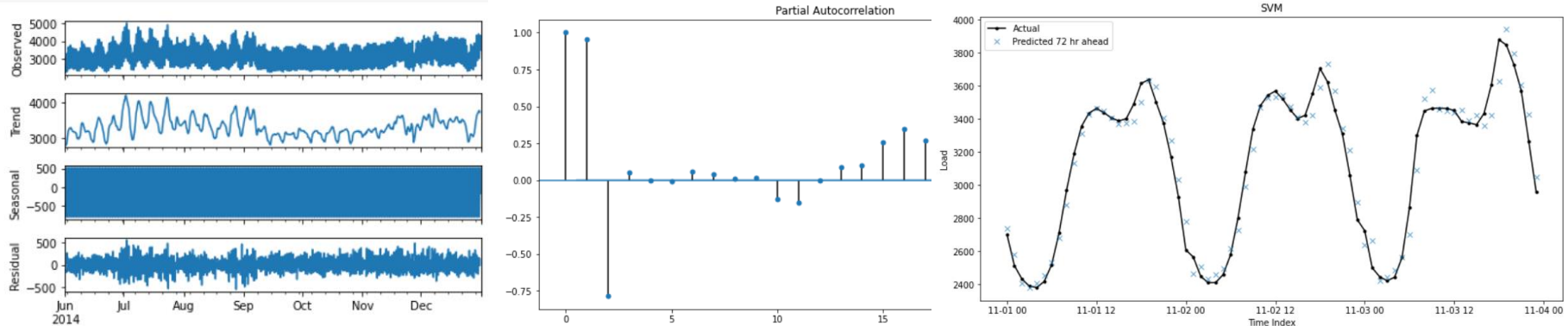


# Data Driven Engineering I: Machine Learning for Dynamical Systems

## Analysis of Dynamical Datasets I: Time Series

Institute of Thermal Turbomachinery  
Prof. Dr.-Ing. Hans-Jörg Bauer

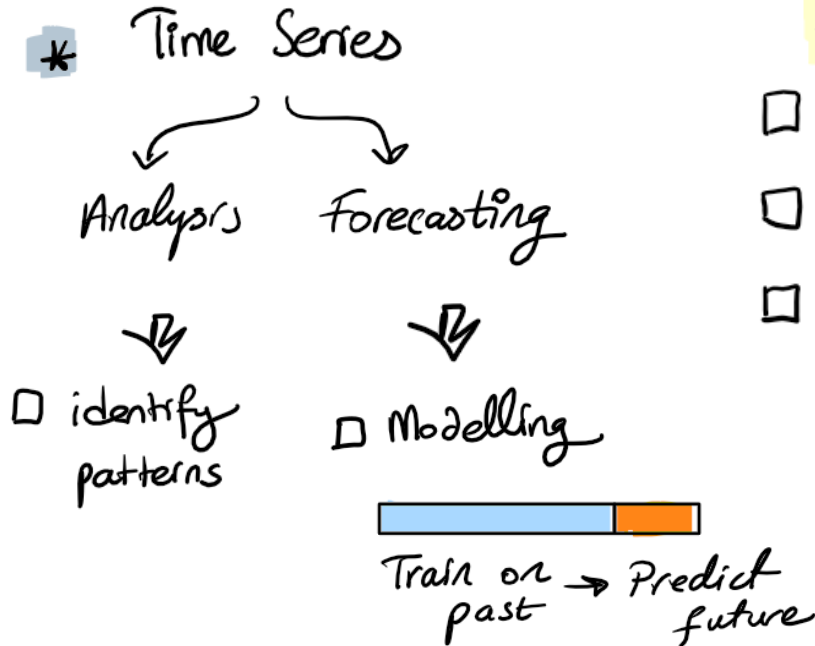


# Dynamical Datasets I: Time Series

## Outline

- \* Time Series = Overview
- \* Statistical Models for time series
- \* State space models  $\Rightarrow$  DDE II
- \* Machine Learning Part I
- \* Machine Learning Part II

# Time Series Analysis



Relatively new field:

- Forecasting ~ old as humankind
- Autoregressive model ~ 1920s
- Box-Jenkins Model ~ 1970

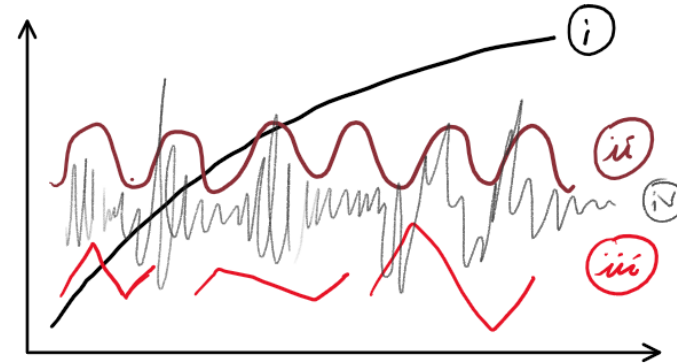
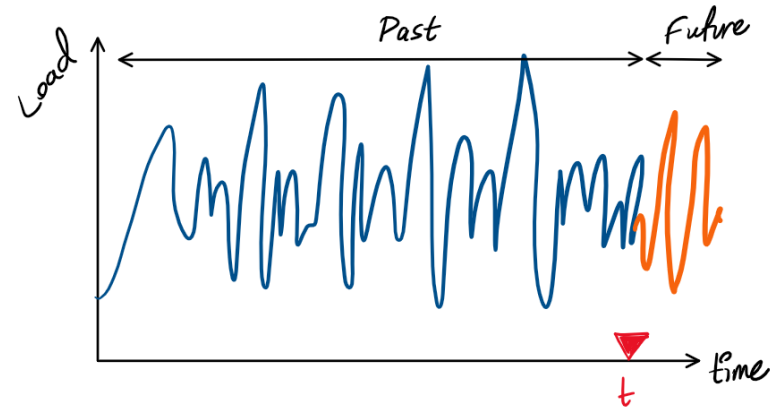


"All models are wrong,  
but some are useful."

G. Box

# Time Series Analysis

## \* Components of time series



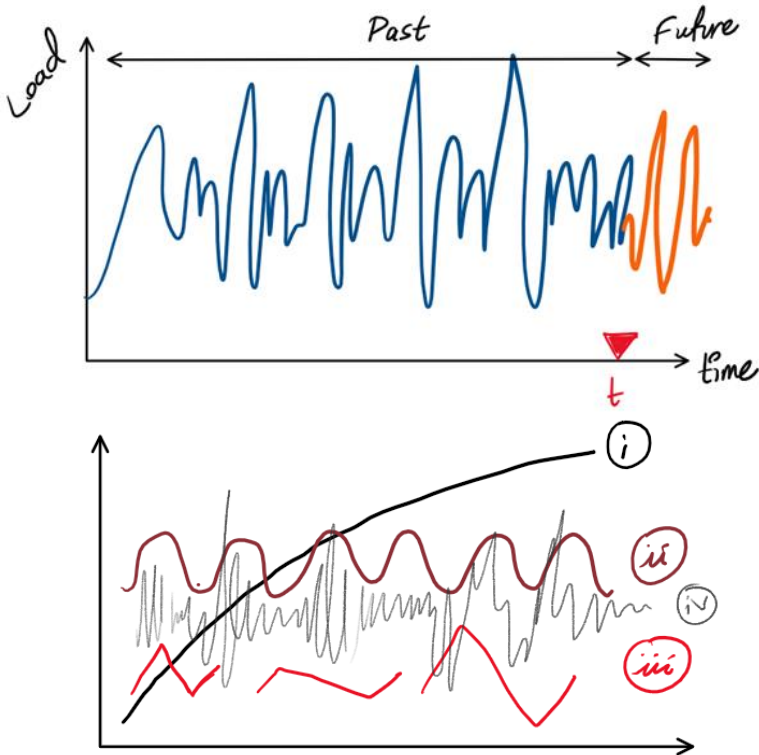
i) Long term trends

iii) Cyclic variations

ii) ST seasonal variations

iv) Random fluctuations

# Time Series Analysis



\* Components analysis



"Stationarity"

$$[\bar{X} \text{ \& } \sigma \neq f(t)]$$

~ Strong ~

$$[\bar{X} \text{ \& } \text{auto covariance} \neq f(t)]$$


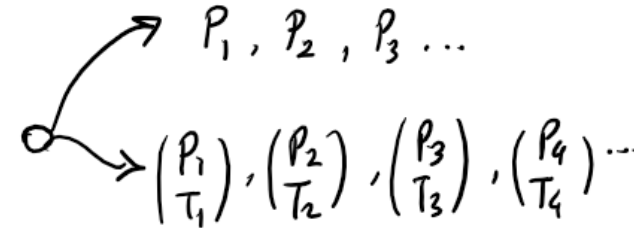
~ weak ~

\* Forecasting:

☑ There is an ordered relationship between observations

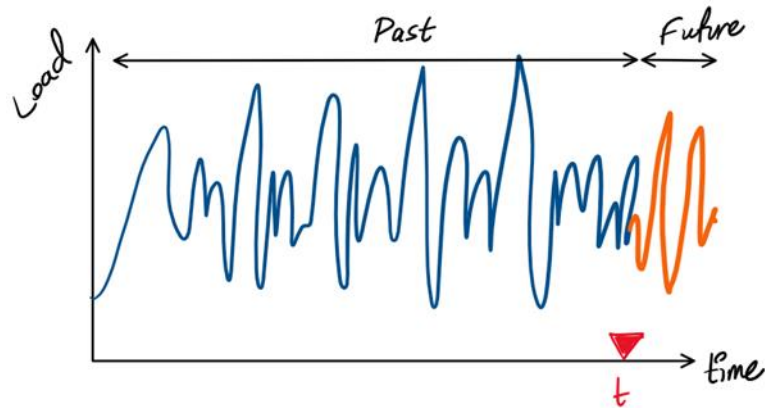
$$t_{i-j} \rightarrow t_i$$

## Before we begin:

- \* horizon of your model (short term vs. long term)
- \* level of granularity you need ( $\Delta t_i$ ) 
- \* univariate or multivariate models 

Before we begin:

\* Sliding Window



time	Load
0	321
1	316
2	314
3	314
4	318
...	

}  $y(t)$

$$y = f(x)$$

Before we begin:

\* Sliding Window

$W=1$

$$y = \begin{bmatrix} 101 \\ 14 \\ 46 \\ 84 \\ 72 \\ 13 \end{bmatrix}$$

→

x

y

NaN

101

→

101

14

→

14

46

→

46

84

→

84

72

→

72

13


$W=2$


	$x_1$	$x_2$	$y$
$\rightarrow$	NaN	NaN	101
$\rightarrow$	NaN	101	14
	101	14	46
$\rightarrow$	14	46	84



Before we begin:

\* Single/multistep forecasting

① Direct multi-step:   $N$  models  $y = f(x)$

② Recursive multi-step:  1 model

$1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4, \dots, N-1 \rightarrow N$

$1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4, \dots, N-1 \rightarrow N$

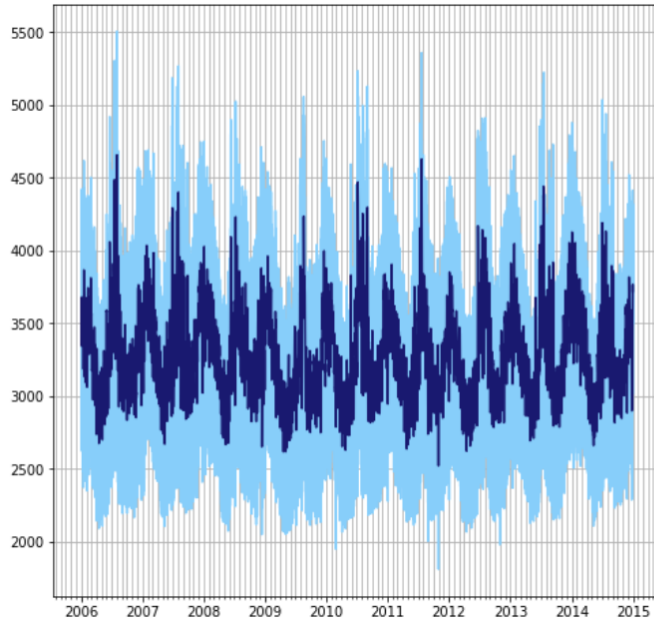
③ Multiple output:  $[ \text{history} ] \rightarrow [ \text{future} ]$

# Time Series Analysis

## Work flow template:

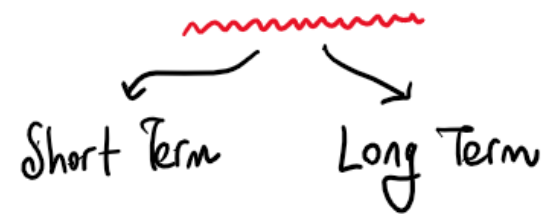
- 1) Understand the problem/business
- 2) Data exploration
- 3) Data preprocessing // feature eng.
- 4) Shortlist the models/algorithms
- 5) Train your model
- 6) Evaluation phase

# Case: Energy Demand Forecasting



\* 8 years data of Temp & Load ( $\Delta t = \text{hr}$ )

? Power Demand forecasting



# Case: Energy Demand Forecasting

\* Short term load forecasting :  $\sim 1$  hr to 24 hr  
 $\sim$  demand / supply

← near past is used

← Temperature is an important feature

\* Long term LF :  $\sim 1$  week to months  
 $\sim$  years

} Planning & investment

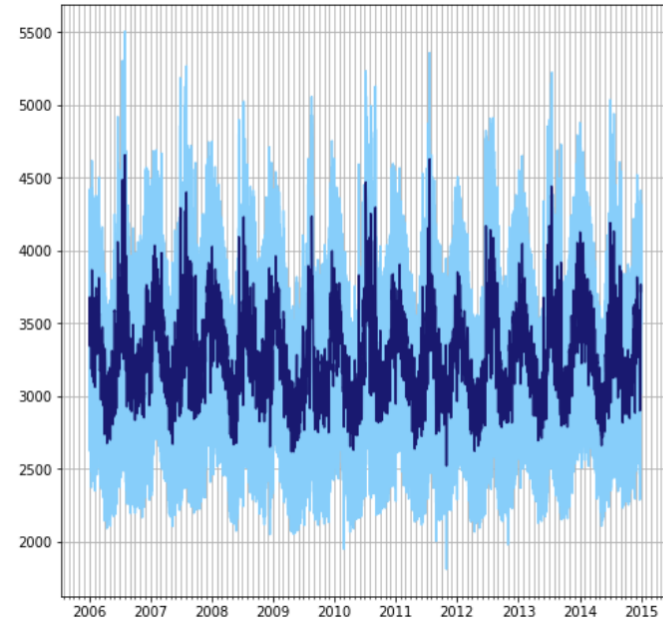
← Seasonal patterns

← Long term trends

← Climate Models

# Case: Energy Demand Forecasting

Typical	STLF	LTLF
Horizon	1 hr - 2 days	$\geq 1$ months
Granularity	$\sim$ hr	$\sim$ hr - day
History Range	$\sim 2$ years	$\sim \geq 5$ years
Accuracy	$\leq 5\%$ error	$\leq 25\%$ error
Forecasting freq.	$\sim$ hr to day	$\geq$ month

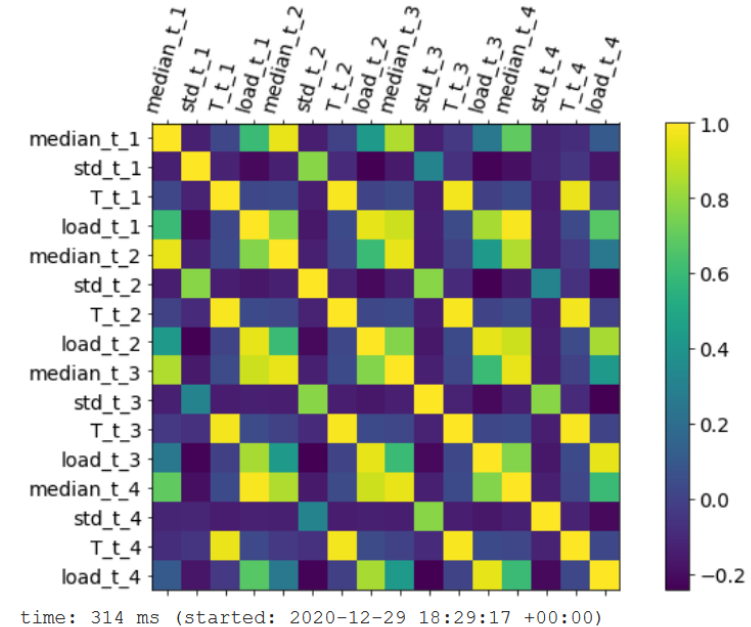


# Data Exploration: What we already know

✓ Basic statistics (mean, median, STD...)

✓ Plots  $\Rightarrow$  1D : Temporal data  
 $\Rightarrow$  2D : Scatter plots  
 $\Rightarrow$  Histograms  
 $\Rightarrow$  Box plots, violin plots

✓ Correlation matrix



# Data Exploration: Temporal Nature of data

## ① How to handle "time stamps",

	Date	Hour	load	T
0	01/01/2004	1	NaN	37.33
1	01/01/2004	2	NaN	37.67
2	01/01/2004	3	NaN	37.00
3	01/01/2004	4	NaN	36.33
4	01/01/2004	5	NaN	36.00



	load	T
2012-01-05 00:00:00	3167.0	19.00
2012-01-05 01:00:00	3014.0	22.33
2012-01-05 02:00:00	2921.0	22.33
2012-01-05 03:00:00	2874.0	22.00
2012-01-05 04:00:00	2876.0	21.67



# colab



# Data Exploration: Temporal Nature of data

## ② Temporal data decomposition

Stationarity



how 'stable' your system ☒ Intuition



how much we should expect

the past reflects itself on future ?



"Self Correlations"





# colab

## Data Exploration: Temporal Nature of data

### ③ Feature Eng. for Time Series

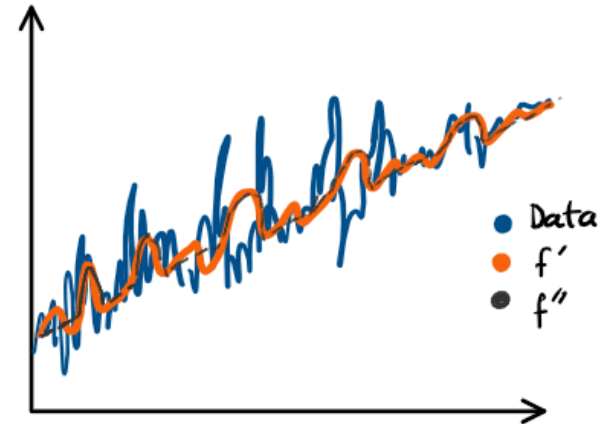
□ Date/time information



# Data Exploration: Temporal Nature of data

## ③ Feature Eng. for Time Series

- Date/time information
- Window functions





# colab

## Data Exploration: Temporal Nature of data

### ④ Self / Auto Correlations in temporal data

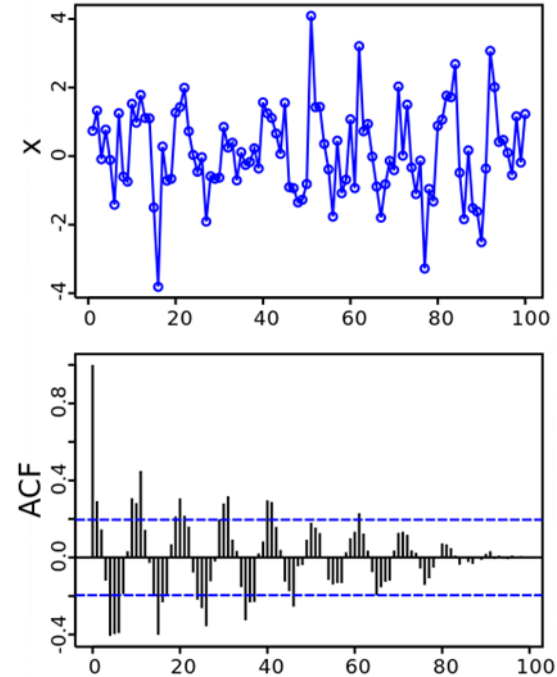
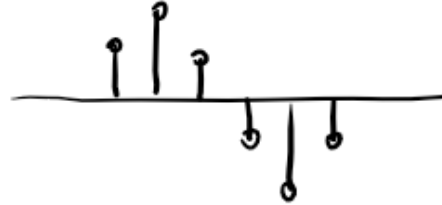
- Autocorrelation function (acf)

- Partial ACF (pacf)

☑ How data points are linearly related as a function of time difference.

# Data Exploration: Temporal Nature of data

- \* ACF  $\Rightarrow$  it preserves the periodicity
- \*  $ACF = 1$  @  $lag = 0$  [self correlated]
- \*  $ACF(\text{white noise}) \rightarrow 0$
- \* ACF is symmetric

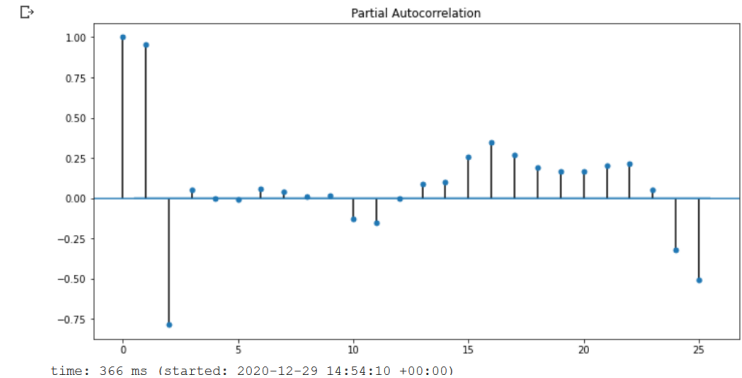
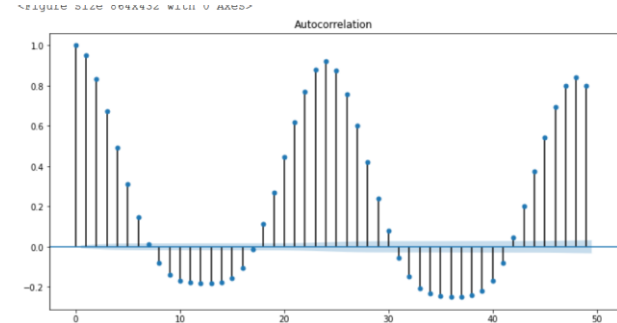


# Data Exploration: Temporal Nature of data

- \* PACF  $\rightarrow$  which time lag is "informative",  
 $\sim$  filters periodic behavior



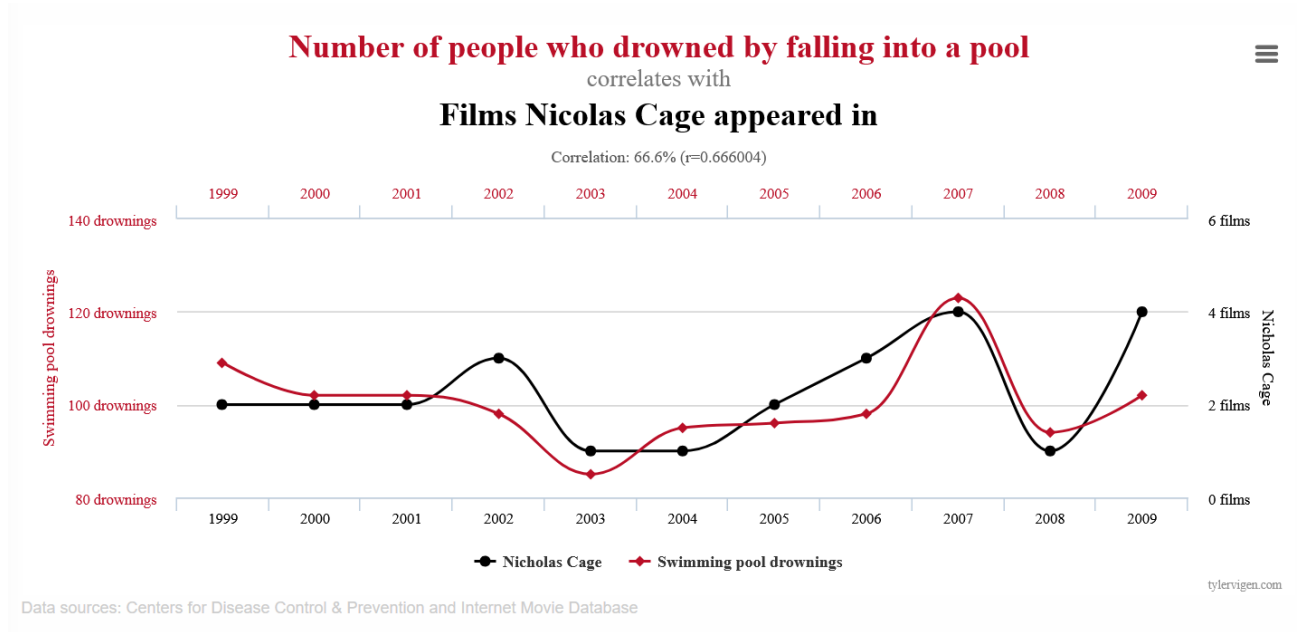
- \* PACF  $\rightarrow$  determine the "order" of a model



time: 366 ms (started: 2020-12-29 14:54:10 +00:00)



# Spurious Correlations





# colab

# Overview of Statistical Models

- \* Obj: (i) find time-related trends
- (ii) find seasonality
- (iii) find auto-corr. (corr. wrt. time)

## AR Model: Auto Regressive

\*  $y_t = a_0 + a_1 y_{t-1} + \text{Err}$  } history := 1 lag

\* Order( $p$ ) := history info;  $p=2$

$$y_t = a_0 + a_1 y_{t-1} + a_2 y_{t-2} + \text{Err}$$

\* Order  $\leftarrow$  "pacf"

# Overview of Statistical Models

## AR-I-MA : AR - Integrated - MA

- \* add differencing  $\Rightarrow$  Remove trends  
"baseline correction"

- \*  $ARIMA = f(p, d, q)$ 

$$\left\{ \begin{array}{l} (0, 0, 0) \rightarrow \text{white noise} \\ (0, 1, 0) \rightarrow \text{random walk} \\ (0, 1, 1) \rightarrow \text{exp. smoothing} \end{array} \right.$$

- \* SARIMA := Seasonal ARIMA

□ Adjacent points in time can have influence on one another

## MA : Moving Average

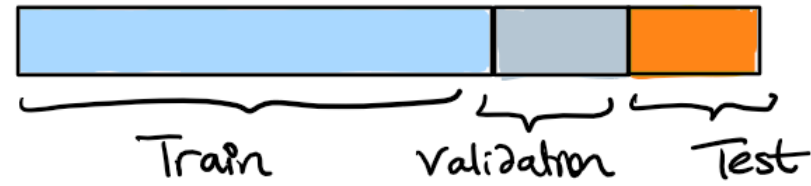
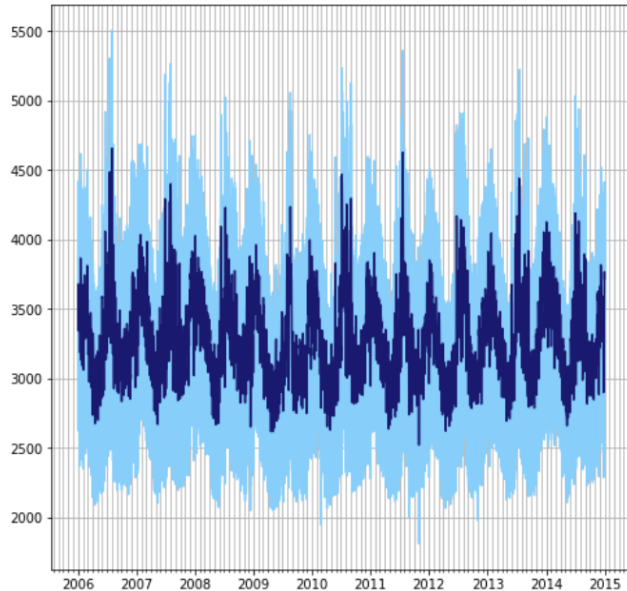
- \*  $y_t = a_0 + E_t + a_1 E_{t-1} + a_2 E_{t-2} + \dots + a_q E_{t-q}$   

$\downarrow$   
Errors dissipate in time

- \*  $q \leftarrow \text{order}$

- \*  $q \leftarrow \text{ACF}$   
"stop error propagation sharply"

# Model Training



Train → model fitting

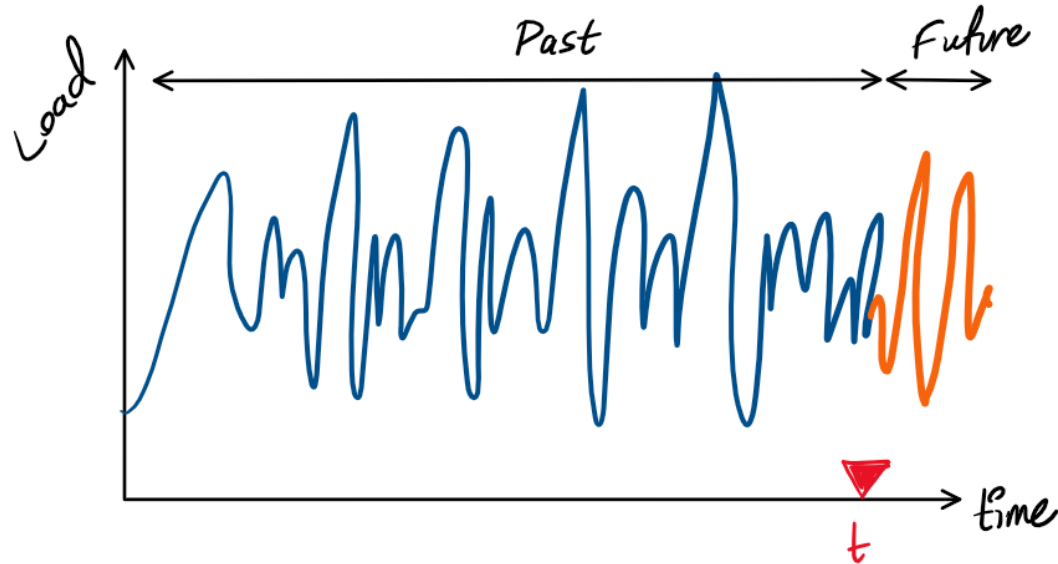
Valid. → hyper-parameter tuning

Test → Performance evaluation



# colab

how can we use ML algorithms?



In M.L.:

$$* [x] \rightarrow [y]$$

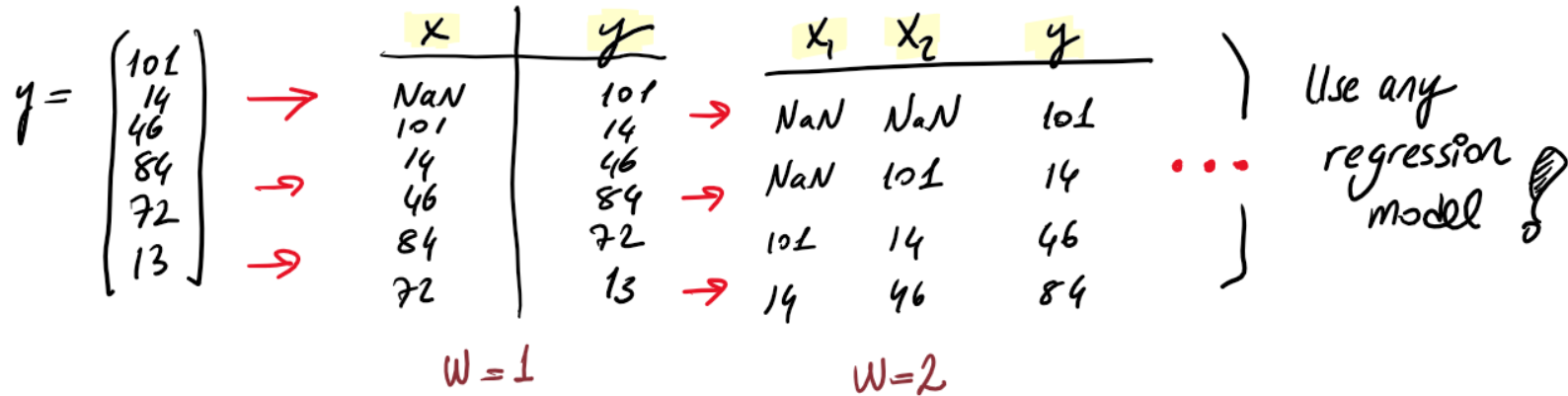
⇒

time	Load
0	321
1	316
2	314
3	
4	318
...	

} y(t)

# how can we use ML algorithms?

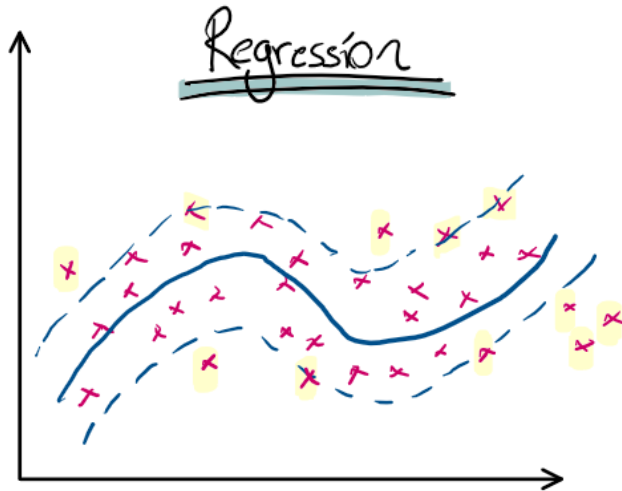
\* Time Series  $\Rightarrow$  "Supervised learning task" [batch // real-time]





# # Model Selection: SVM for Regression

## — SVM —



- \* Fit as many instance as possible
- \* "Street" width is controlled by margin  $\epsilon$ .
- \* Convex optimization problem;
  - ☒  $C$
  - ☒  $\epsilon$
  - ☒ Kernel



# colab

# Additional Notes