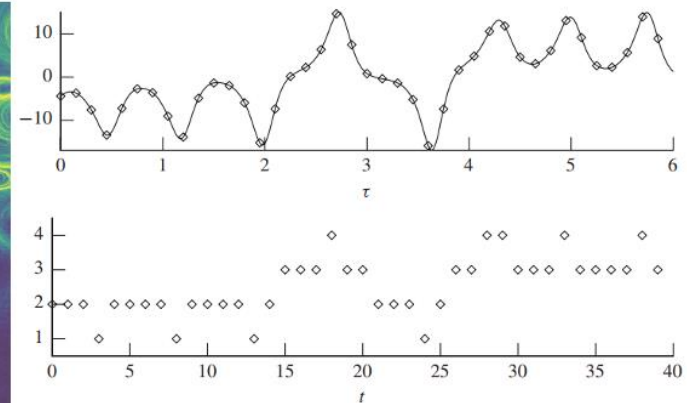
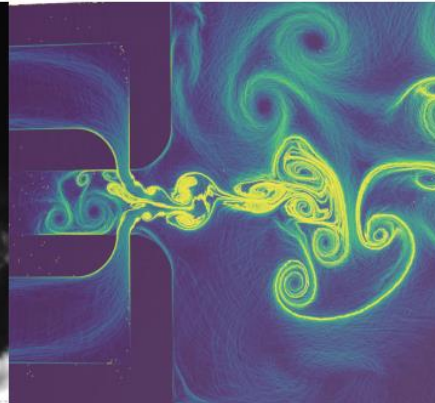



# Data Driven Engineering II: Advanced Topics

## Feature Engineering I

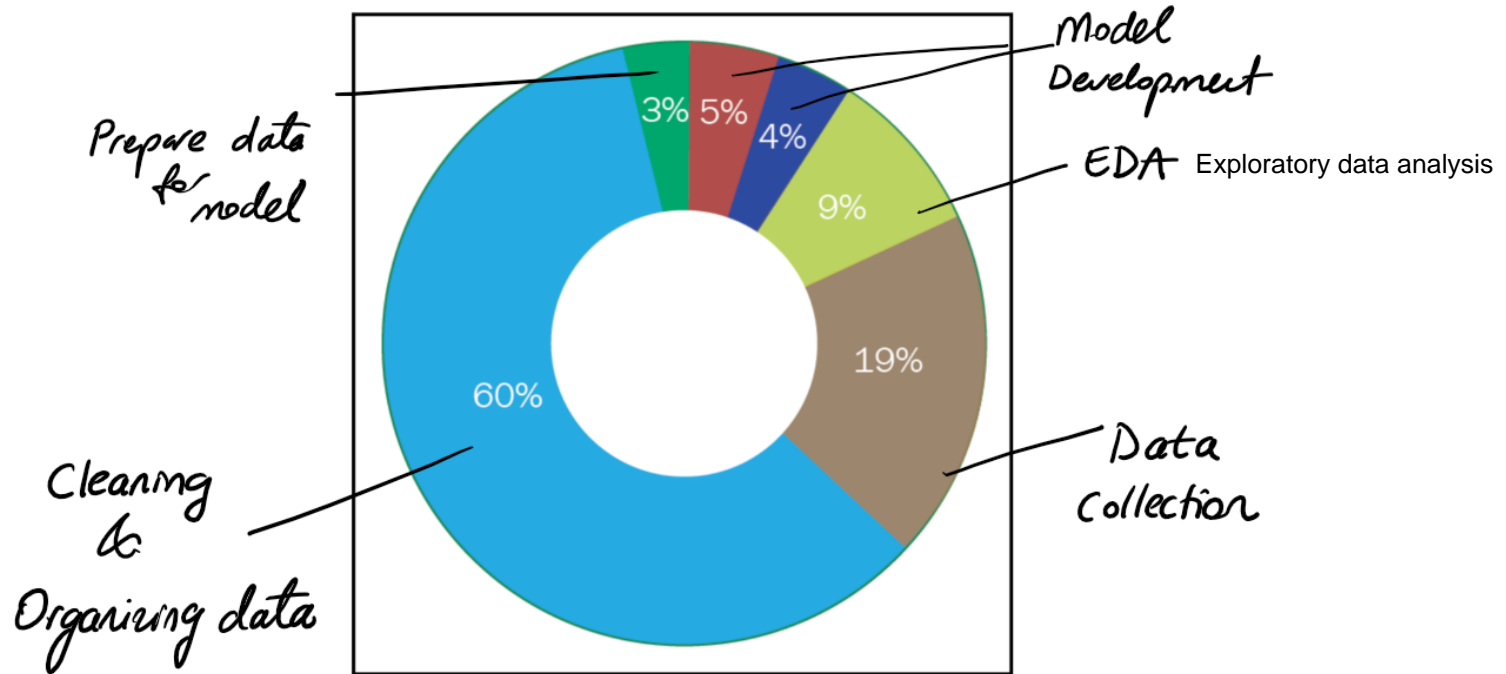
Institute of Thermal Turbomachinery  
Prof. Dr.-Ing. Hans-Jörg Bauer



## To Do List

- ☐ Check & think about projects
- ☐ Register via Ilias ~~⇒~~ latest 13<sup>th</sup> May
- ☐ Check dataset & materials
-  Involve in projects ~~⇒~~ Register for HPC Access  
(ILIAS)

How do you spend your time?

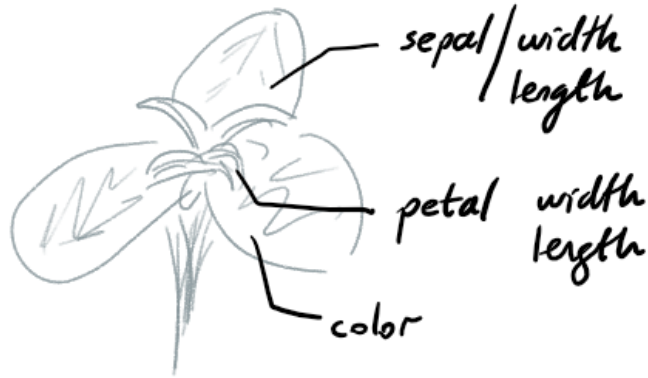


# Today's Agenda:

- 1) Feature Engineering (FE)
- 2) Data Preparation  $\Leftrightarrow$  FE
- 3) Continuous  $X \Rightarrow y$

# Introduction to features

□ Feature := attributes defining properties of objects

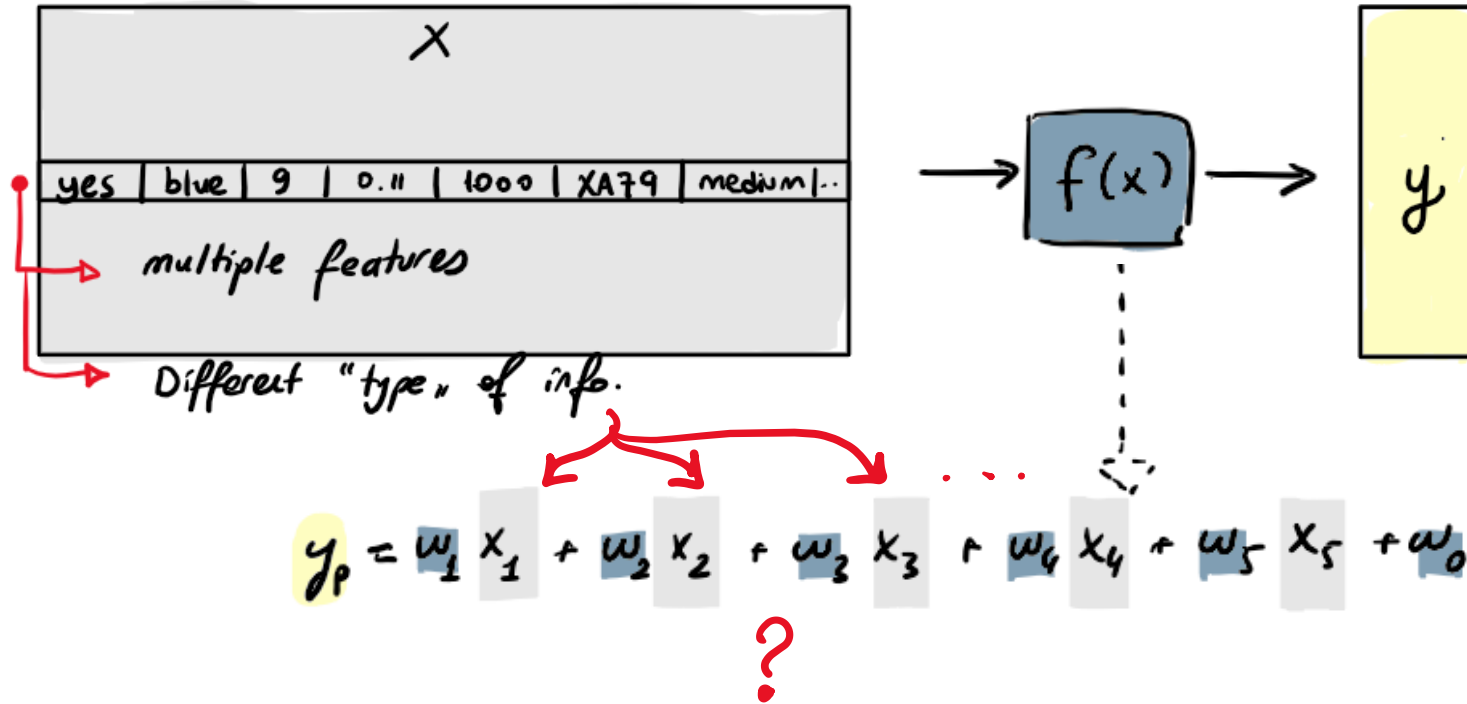


□ Goal := how features are distributed among objects  
:= how they are assoc. with a dep. var. (label)



Success  $\Rightarrow$  Depends on how informative features are  
-----> Data Collection >> -----

# Features in real life problems



# Features in real life problems

yes	blue	9	0.11	1000	XA79	medium
-----	------	---	------	------	------	--------

$\Rightarrow x_{\perp} :=$  vector description of  $y_{\perp}$   
 $\hookrightarrow$  converted to numbers

- Homogenous / Heterogenous in type

**Categorical** : feature value has no order in it.

e.g. // "color",  $=$  [blue, white, black, pink]

$\Rightarrow$  [0, 1, 2, 3]  $\Rightarrow$  vector space  $\gg$  ordered !

**Ordinal** : Set of ordered values  $\Rightarrow$  different than cat. !

e.g. // "Doneness"  $\rightarrow$  [Rare  $\rightarrow$  Medium  $\rightarrow$  well-done]  $\Rightarrow$  [0, 1, 2]

Categorical and ordinal features are two types of variables used in statistics and data analysis. They both describe non-numeric data, but they differ in terms of the information they convey and how they can be analyzed.

Categorical features:

- Categorical variables, also known as qualitative or nominal variables, represent distinct categories or groups. These categories have no inherent order or ranking.
- Examples include gender (male, female), colors (red, blue, green), or types of animals (cat, dog, bird).
- Categorical variables are typically used to describe the qualities or characteristics of a dataset, and they are often analyzed using techniques like frequency tables, bar charts, or chi-square tests.

Ordinal features:

- Ordinal variables, also known as ordinal data, represent categories with a specific order or ranking. These variables can be arranged in a meaningful sequence, but the differences between the levels are not quantifiable.
- Examples include survey responses (strongly disagree, disagree, neutral, agree, strongly agree), educational levels (high school, undergraduate, graduate), or customer satisfaction ratings (poor, average, good, excellent).
- Ordinal variables can be analyzed using methods similar to those used for categorical variables, but they can also be analyzed with techniques that account for the inherent order, such as the median, percentiles, or non-parametric statistical tests.



# Features in real life problems

yes	blue	9	0.11	1000	XA79	medium	...
-----	------	---	------	------	------	--------	-----

⇒  $x_1$  := vector description of  $y_1$   
↳ converted to numbers

- Homogenous / Heterogenous in type

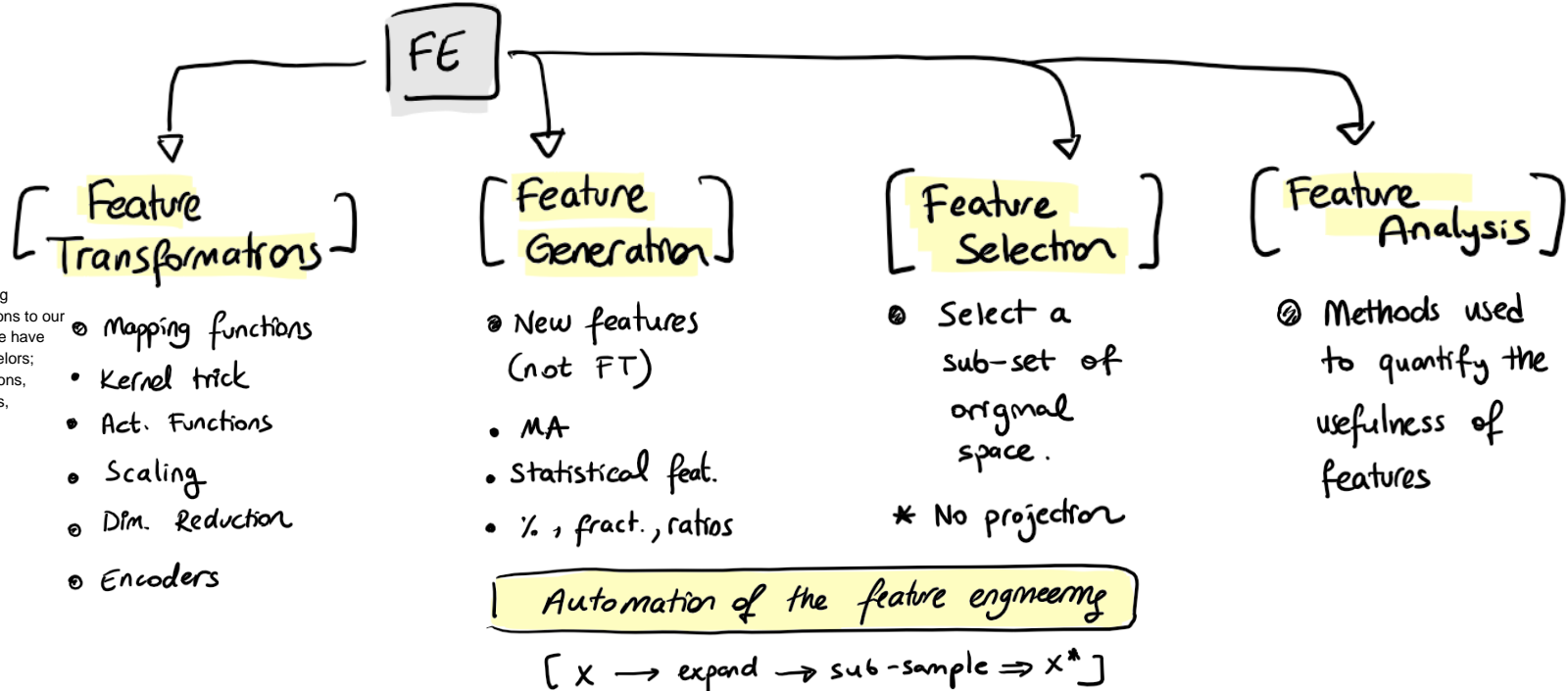
Continuous: quantify the relationship of "nearbyness"

--- ⚠ Ratio Scaled:  $x_1 = 2$ ;  $2 \times 2 \rightarrow x_1$

↳ Object property assoc. with feature  $x_1$  (length) doubled.

$$y_p = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5 + w_0 \quad \Rightarrow \text{How to handle it?}$$

□ Operations we perform on object vectors during the whole pipeline.



Here we are applying mathematical functions to our data. This is what we have learned in our Bachelors; Fourier transformations, differential equations, etc.

# KITHealthTech

KIT Center of Health Technologies

## KARE

The vision of KITHealthTech is to create a unique interaction with patients and citizens, physicians and clinics →

## Thematic Fields

To enable and accelerate cutting edge research in Health Technologies, KITHealthTech will focus on three main overarching Thematic Fields →

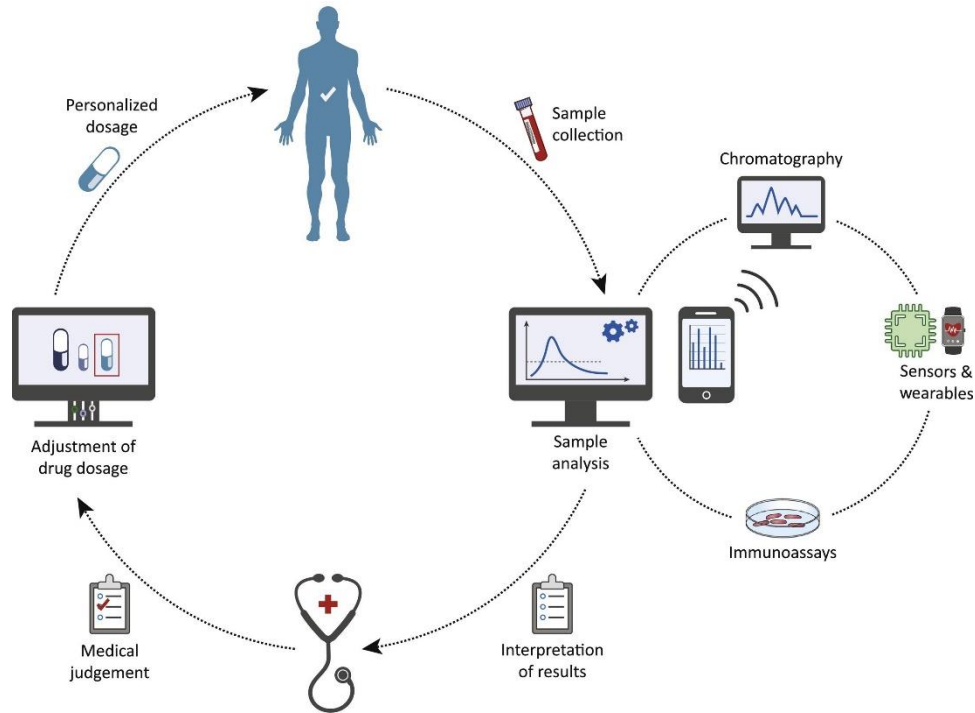
## Research Fields

12 Main Research Topics of KITHealthTech →

## Focus Fields

Focus Fields are dynamic structures giving the FOCUS to highly relevant scientific challenges →

# Case Study: Th. Drug Management at ICU



<https://doi.org/10.1016/j.tibtech.2020.03.001>

- i. Measuring this drug concentration in blood / plasma / non-invasive options
- ii. Medical interpretation
- iii. Dose regimen adjustment

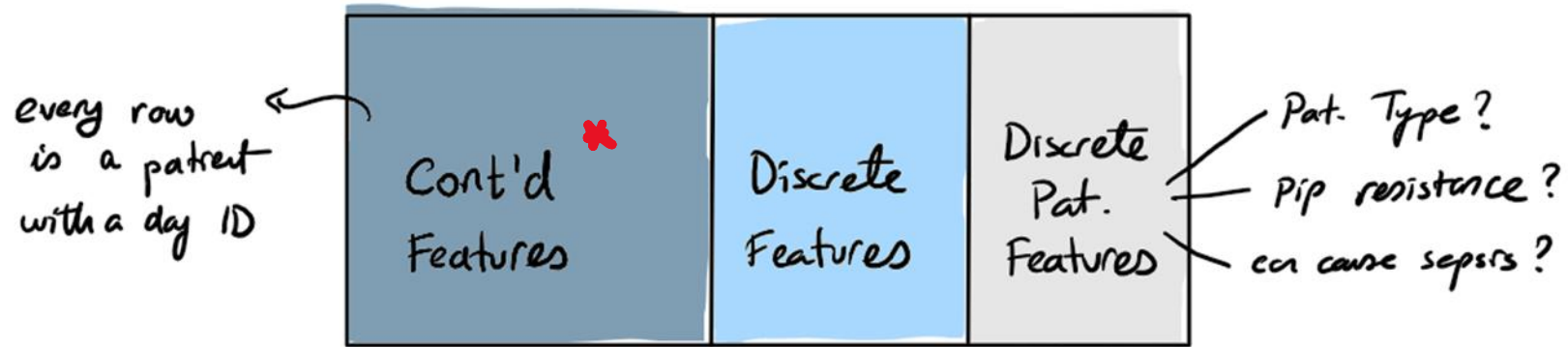
## Objective:

- Personalized dosage regimen

## Challenge (assume measured):

- how to **represent** and **interpret** the **response** of the patient to the drug

## Case Study: Th. Drug Management at ICU



\* 2386 instances  $\Rightarrow$  day info  $\times$  patients

\* [37 + 54 + 108] features

Numerical data : Easy to integrate !

[Q1] Do I know any features that are not related to my objective? (assumption  $\Leftrightarrow$  Domain knowledge)

[Q2] Can I filter uninformative columns?

$\rightarrow \emptyset$  variance

$$y = w_1 x_1 + w_2 \underset{\substack{\downarrow \\ 1}}{x_2} + w_0 \quad \} = w_1 x_1 + w_0'$$

If there is little variance in the column, that column is acting as a variance.

if you have two features with the same or highly correlated information, they may introduce multicollinearity in your model. Multicollinearity occurs when two or more independent variables in a regression model are highly correlated, leading to unreliable and unstable estimates of the regression coefficients. This can make it difficult to determine the individual contribution of each feature to the model and can lead to overfitting.

$\rightarrow$  "little variance" & duplicates

Notebook

Numerical data : Easy to integrate !

[Q3] Do the magnitudes matter for the problem?

•  $(-)/(+)$   $\Rightarrow$  is enough? Is the magnitude important? For example, input of heat transfer is positive means that is going in.

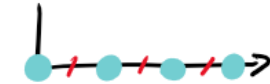
• magnitude @ coarser granularity ?

$\hookrightarrow$  "counts" := # stations, daily visits etc.

We are not interested in the exact number of bus stations as 899 or 901, we need a rough number of 900. That's when quantization comes in.

## Quantization

\* Group counts  $\Rightarrow$  "bins" ☐ Fixed-width binning



$\hookrightarrow$  continuous var.  $\Rightarrow$  Discrete ☐ Quantile binning



Notebook

Numerical data : Easy to integrate !

[Q3] Do the magnitudes matter for the problem?

Scale of features  $\Rightarrow$  models smooth func. of  $x$  Need for Scaling

$\rightarrow$  linear models; clustering; KNN, RBF, SVM ...  
 $\rightarrow$  anything using distance

Logical Functions  $\Rightarrow$  Insensitive to scale Not needed for Scaling  
 $\rightarrow$  step functions; space partitioning (trees)  
 $\rightarrow$  Be careful if domain shifts !



Fixed-width binning (Group counts): This method divides the range of the continuous variable into fixed-width intervals or bins. The width of the intervals is determined by specifying the number of bins or by providing a specific width for each bin. Each data point is then assigned to one of the bins based on its value. The main advantage of fixed-width binning is its simplicity and ease of interpretation. However, it can be sensitive to the choice of bin width and may lead to uneven distribution of data points across bins if the data is not uniformly distributed.

Quantile binning (Continuous variable): This method divides the continuous variable into bins such that each bin contains approximately the same number of data points. To achieve this, the data is first sorted, and then the range is divided into intervals based on the quantiles (e.g., quartiles, deciles, percentiles). This ensures that each bin has an equal number of data points, leading to a more balanced distribution across bins. Quantile binning is particularly useful when dealing with skewed data or when you want to ensure that each bin has an equal number of observations, which can help with some statistical analyses.

# Data Scaling

## □ min-max scaling

$$x^* = \frac{x - \min(x)}{\max(x) - \min(x)} \Rightarrow [0, 1]$$

## □ Standard Scaling

$$x^* = \frac{x - \text{mean}(x)}{\text{sqrt}(\text{var}(x))}$$

## □ l-norm Scaling

$$l_2 \Rightarrow x^* = x / \|x\|_2 \Rightarrow x / \sqrt{x_1^2 + x_2^2 + \dots + x_m^2} \text{ for } m \text{ examples}$$

If X is sparse:

[ 0 ... 0 ... 1 ... 0 ]

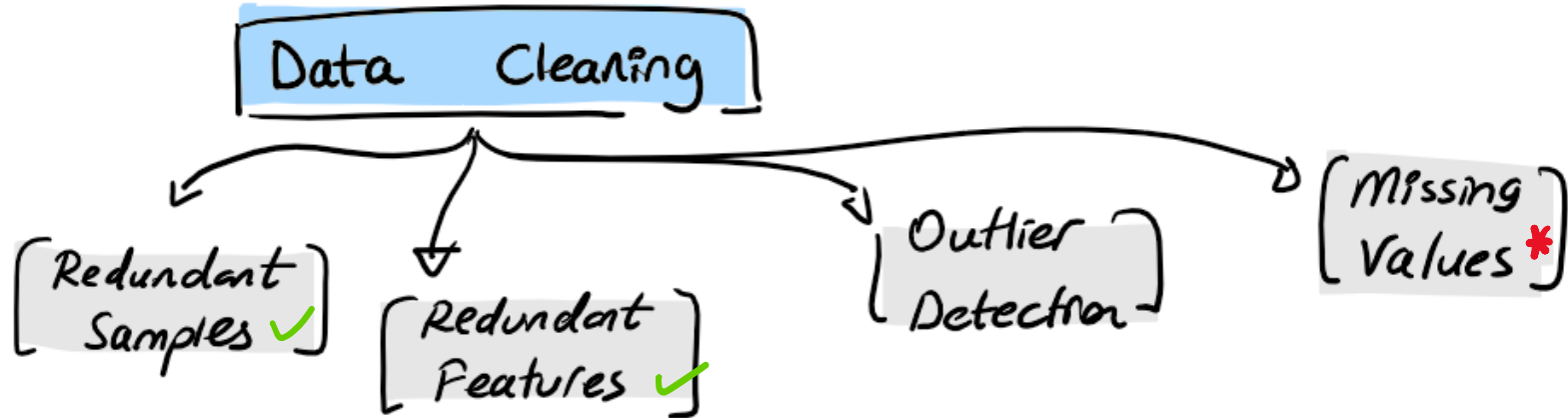


[ 0.1 ... 0.4 ... 0.1 ]

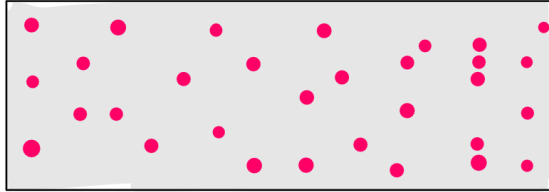
"Dense"

Notebook

Numerical data : Easy to integrate !

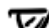



# Handling the missing values: imputation



Fill up the column, (Which are the closest neighbors in different rows?  
And then compute the average and impute the missing value.


In every column you just train a bunch of regression models in each column.

- Simple Imputer  $\Rightarrow$  mean, median, most-frequent, constant
- KNN Imputer  $\Rightarrow$  use "n" neighbours to fill missing cells  
 Distance-based
- Iterative Imputer  $\Rightarrow$  Loop over to solve multivar. reg. problem  
 Distance based

Notebook

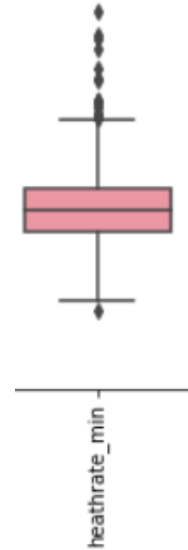
# Outlier Detection & Removal

□ Distribution  $\sim$  Gaussian-like ; use SD as a cut-off  
 $\hookrightarrow$  2SD for small dataset ; 4SD for large datasets

□ IQR based  $\gg$    $\gg$   $IQR = Q3 - Q1$   
 $\hookrightarrow$  character. length scale

□ Automated Outlier Detection

- ⊙ Local outlier factor
- ⊙ Isolation forest



## Local Outlier Factor

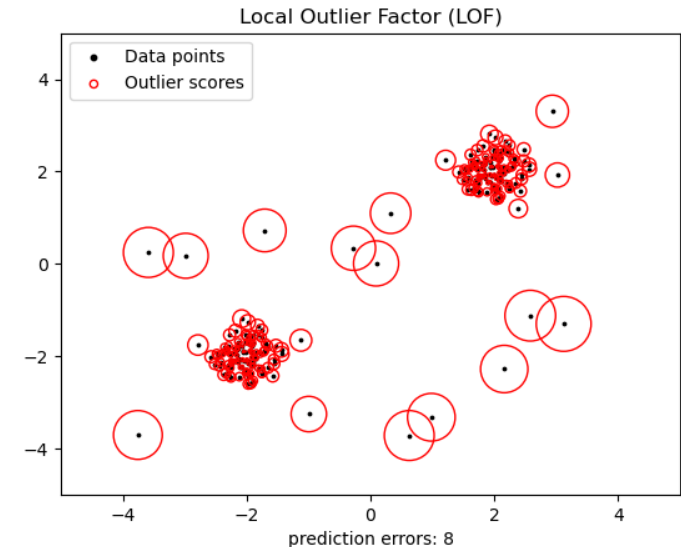
□ Idea is similar to DBSCAN := Local density

□ Compare the density wrt.  $k$  nearest neighbor

→  $LOF \ll 1 \rightarrow$  inlier;

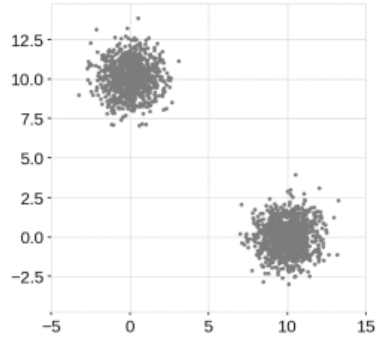
→  $LOF \sim 1 \rightarrow$  similar density

→  $LOF \gg 1 \rightarrow$  Outlier

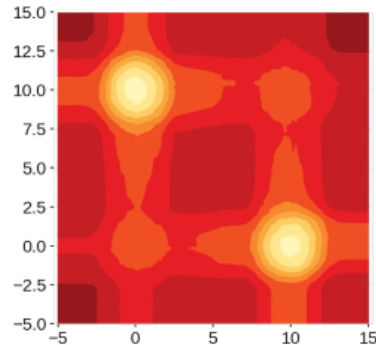


# Isolation forest

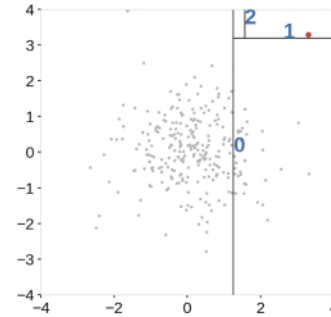
- Anomalies are few in number.  
↳ very different in feature space.



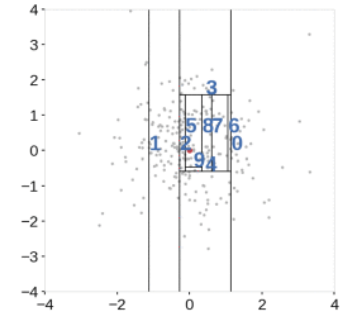
(a) Two normally distributed clusters



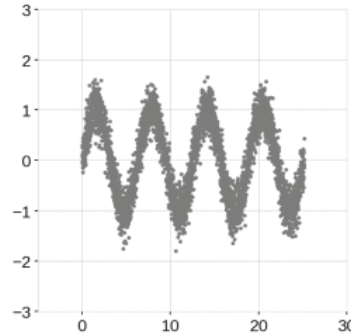
(b) Anomaly Score Map



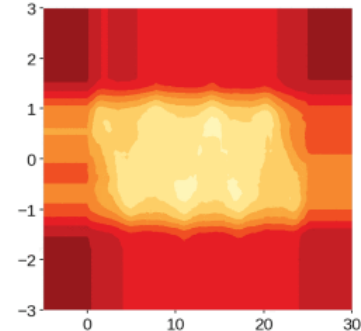
(a) Anomaly point



(b) Nominal point



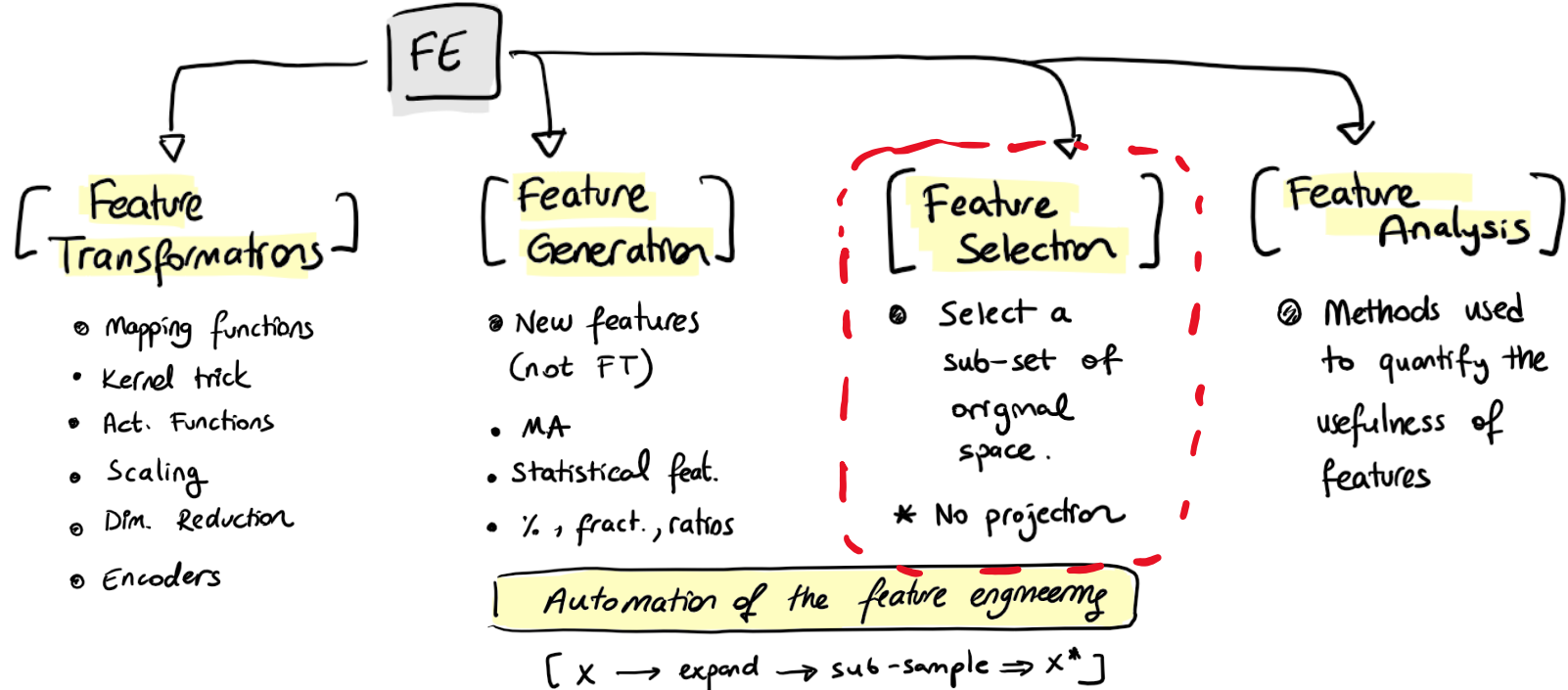
(a) Sinusoidal data points with Gaussian noise.



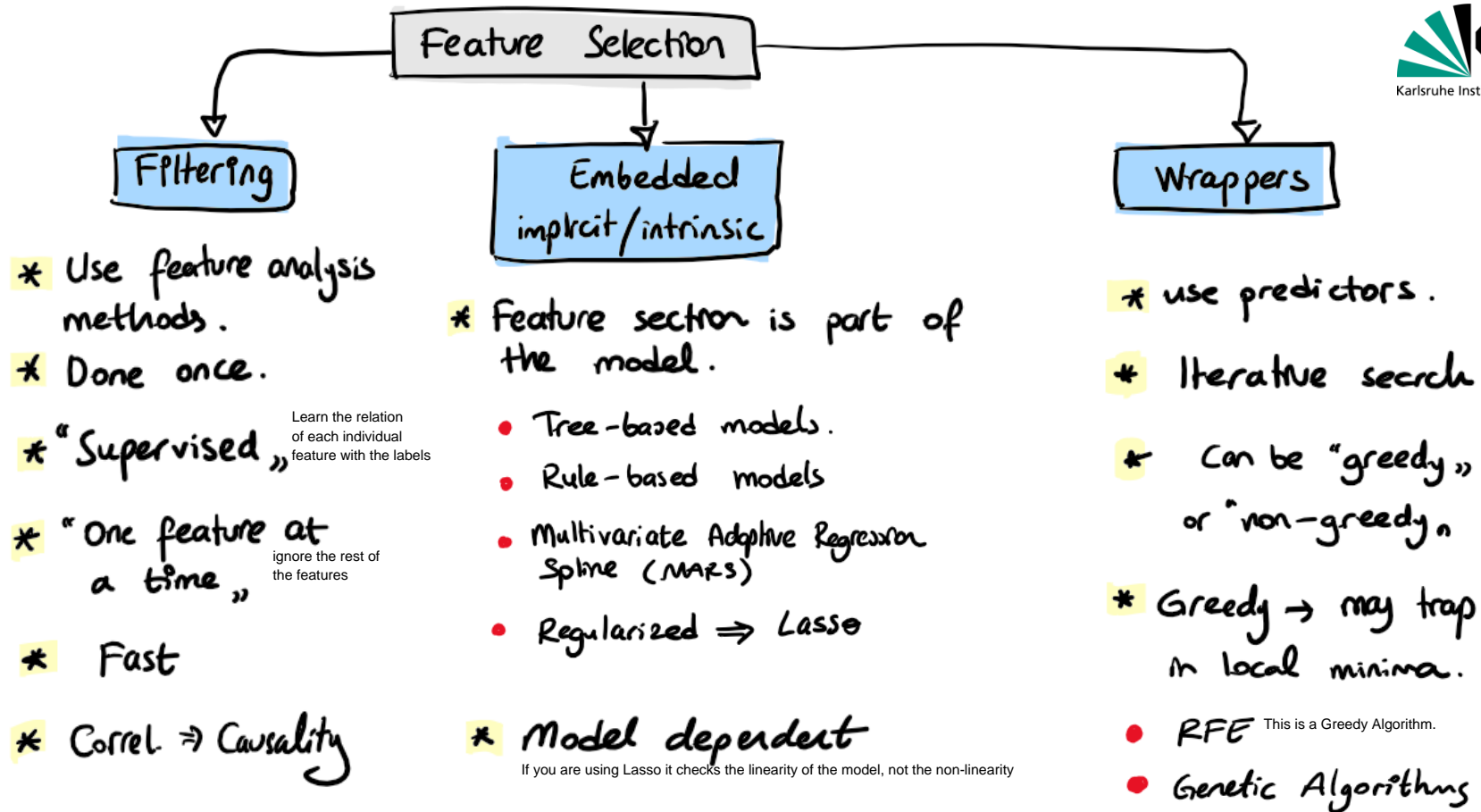
(b) Anomaly Score Map

Notebook

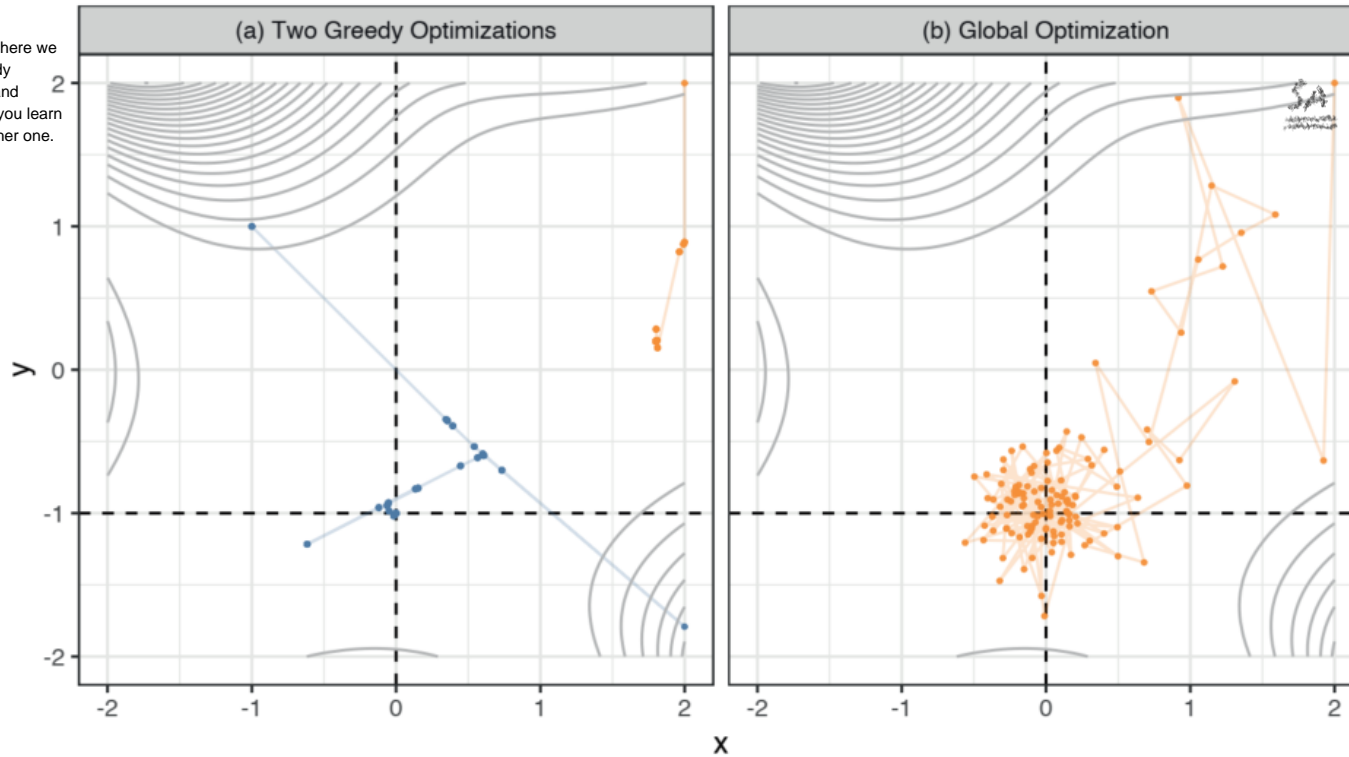
□ Operations we perform on object vectors during the whole pipeline.



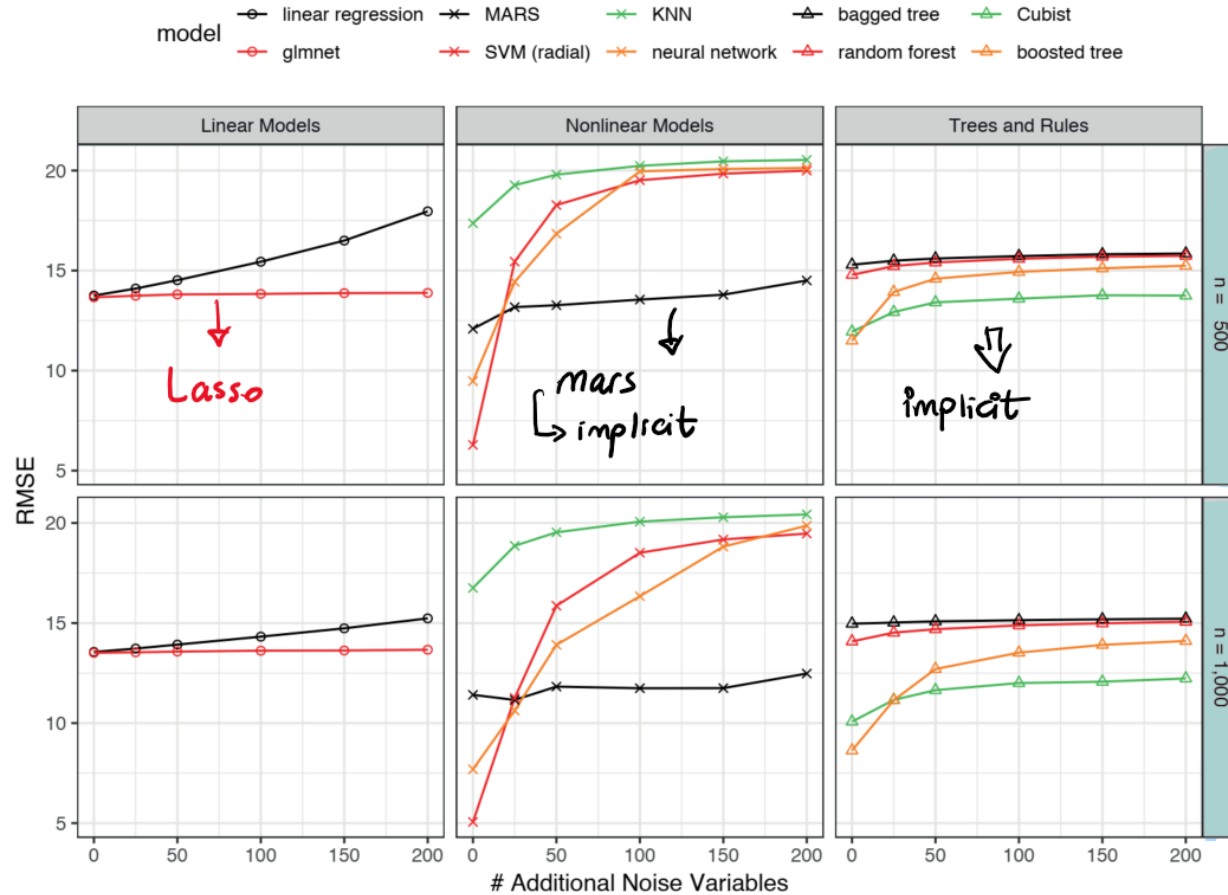




start —•— (-1, 1) —•— (2, 2)



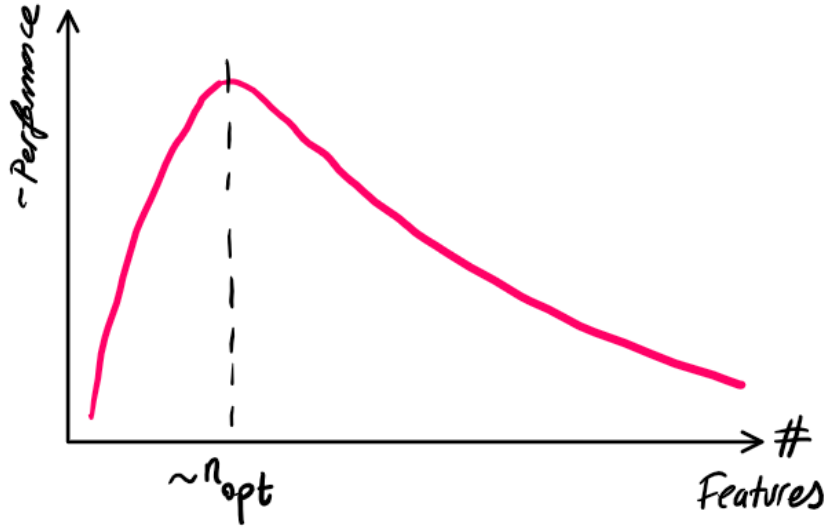
$$f(x, y) = [1 + (x + y + 1)^2(19 - 14x + 3x^2 - 14y + 6xy + 3y^2)] \times [30 + (2x - 3y)^2(18 - 32x + 12x^2 + 48y - 36xy + 27y^2)].$$



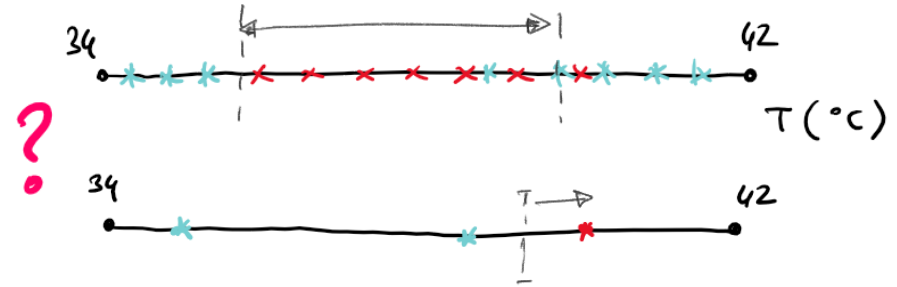
Feature  
Elimination  
helps !

More data may  
or may not  
help !

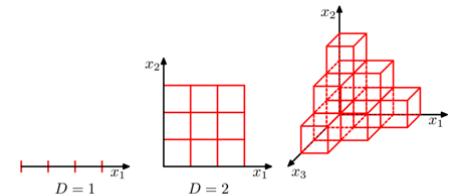
# Curse of dimensionality



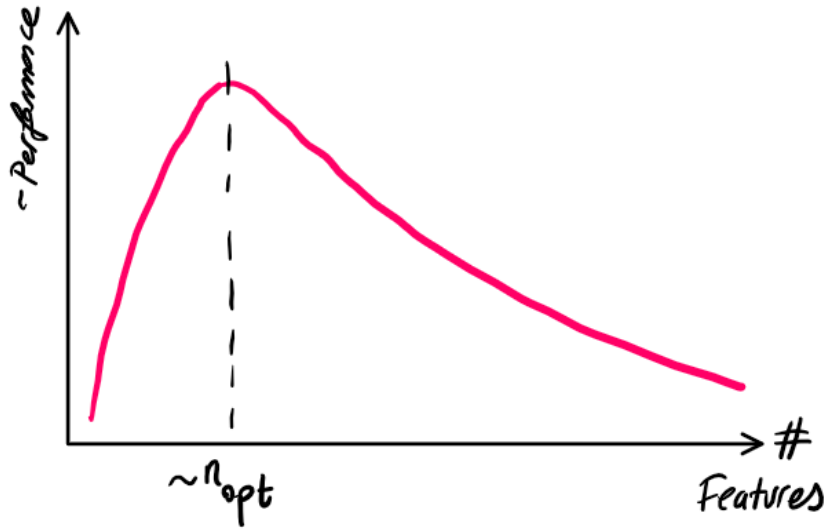
- \* Space gets large  $\Rightarrow$  more data
- \* 'Distance' gets meaningless



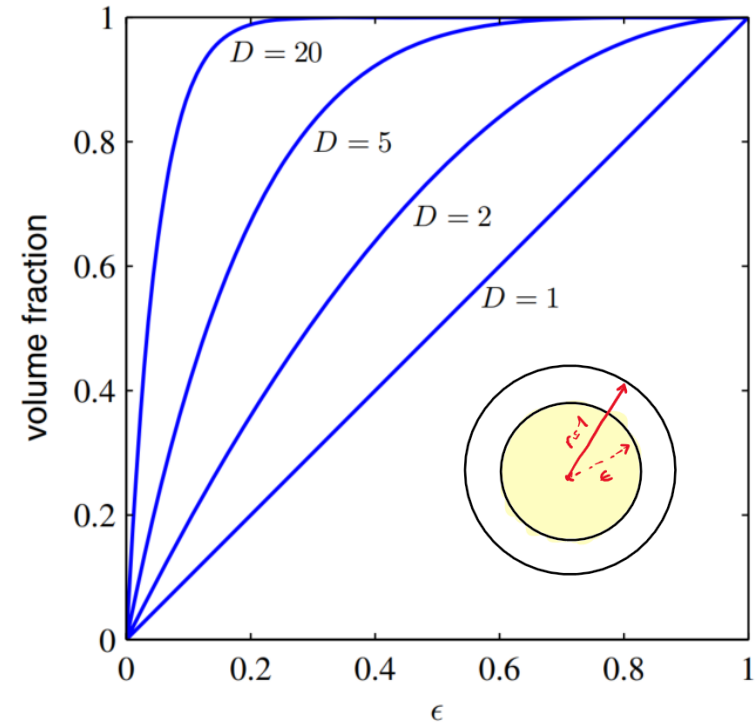
- (1D)  $\Rightarrow$  10 points
- (2D)  $\Rightarrow$  100 points
- (3D)  $\Rightarrow$   $10^3$  points



# Curse of dimensionality



- \* Space gets large  $\Rightarrow$  more data
- \* 'Distance' gets meaningless



## Feature Selection

□ Goal := parsimonious model  $\Rightarrow$  Reduce model complexity

! Some models (SVM, ANN) sensitive to irrelevant features;

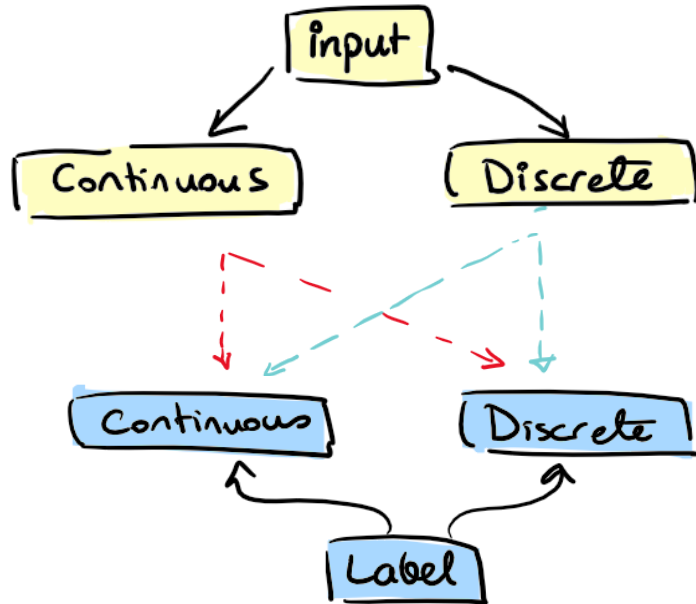
! Some models (LR, Logit) vulnerable to correlated features.

□ Curse of Dimensionality  $\Leftrightarrow$  Data Density to learn

! Which metric & method to use based on data type?

# Feature Selection

- Goal := parsimonious model  $\Rightarrow$  Reduce model complexity

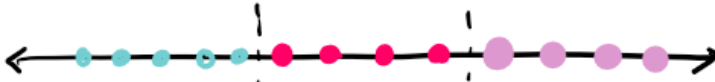


- ① Feature selection methods can depend on data type.
- ② Special care is needed for heterogeneous feature space.

## Case I: Numerical $\Rightarrow$ Categorical

### Anova - F Score

- Compares means & variances of categories.

eg.  $X_1 \Rightarrow$    
values are different for each cat. value.

### Mutual Information

- $MI_{ij} = Entropy(i) - Entropy(i|j)$
- Used for "feature & label" couple
- Generalize well to multiclass

Notebook



## Case I: Numerical $\Rightarrow$ Numerical

### Pearson Correlation

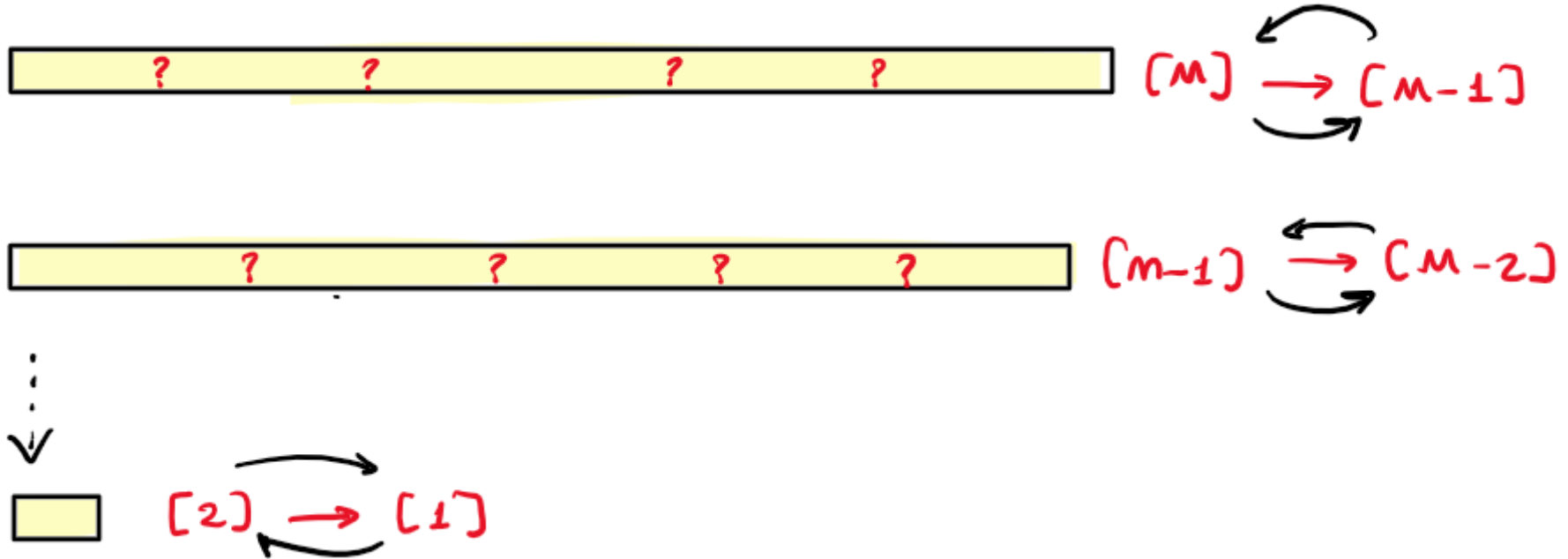
- \* What we have in base Corr. matrix.
- \* Lin. Corr. between feature & label.

### Mutual Information

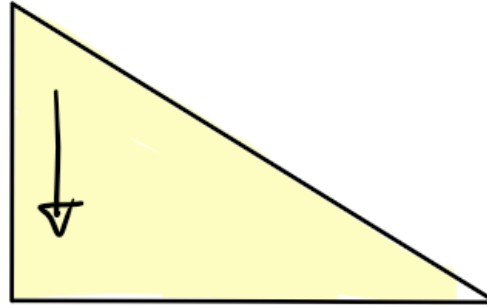
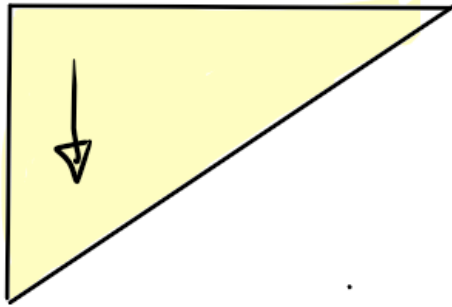
- $MI_{ij} = \text{Entropy}(i) - \text{Entropy}(i|j)$
- Used for "feature & label" couple
- Generalize well to multiclass

Notebook

# Wrappers - I : Greedy Approach



## Wrappers - I : Greedy Approach

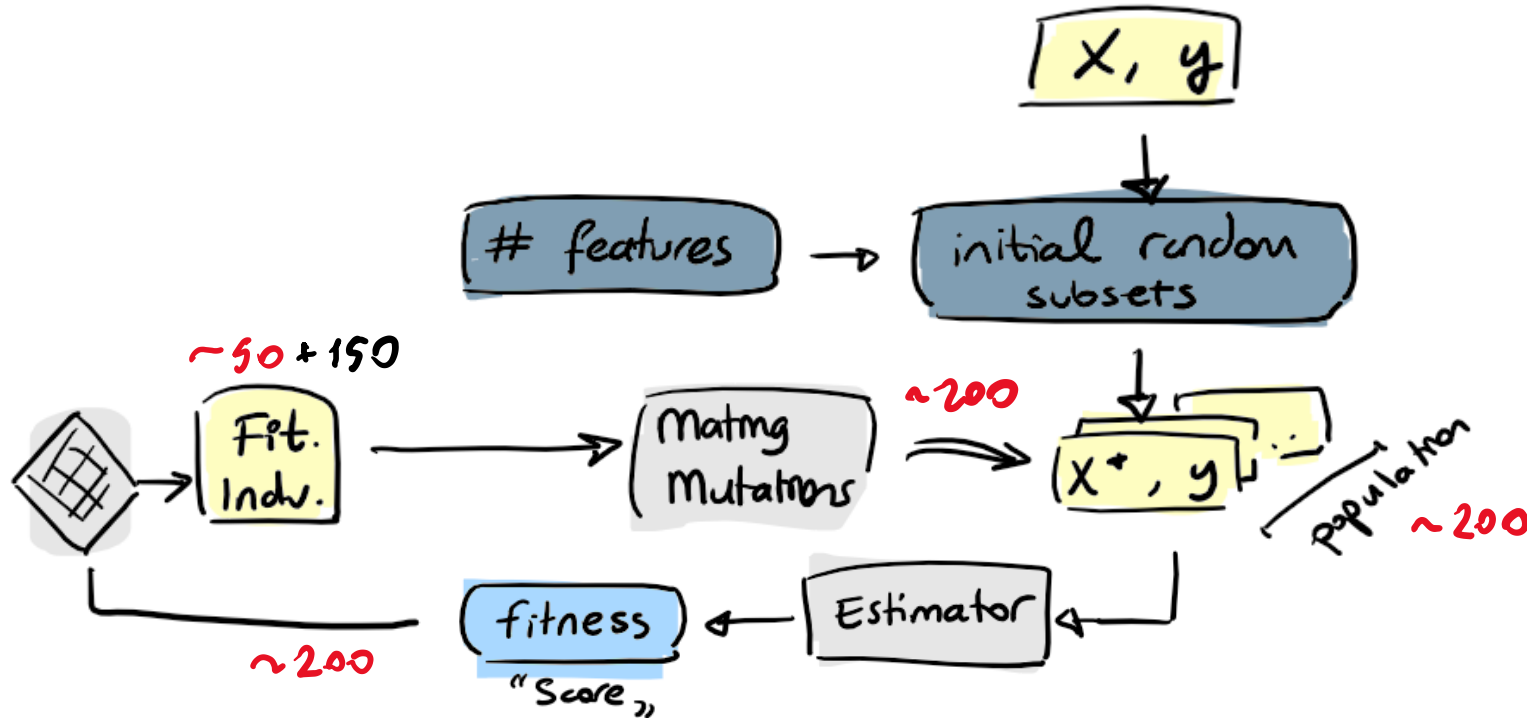


{ NP-hard }

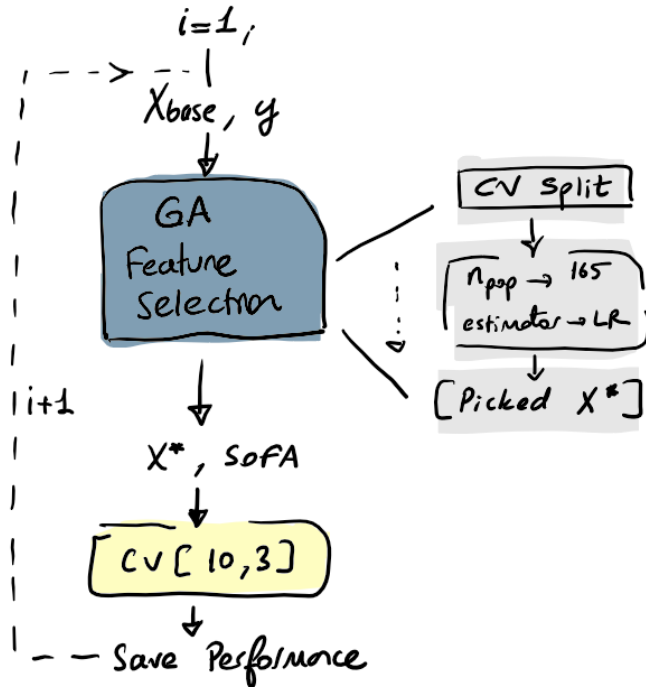
- \* Execute many times on random subsets  $\Rightarrow$  Check freq. of selection
- \* Stochastic search; simulated annealing; gradient decent; ...

Notebook

## Genetic Algorithm for Feature Selection



# Genetic Algorithm for Feature Selection





# Simulated annealing algorithm

```
1 Create an initial random subset of features and specify the number of iterations;
2 for each iteration of SA do
3   | Perturb the current feature subset;
4   | Fit model and estimate performance;
5   | if performance is better than the previous subset then
6   |   | Accept new subset;
7   | else
8   |   | Calculate acceptance probability;
9   |   | if random uniform variable > probability then
10  |   |   | Reject new subset;
11  |   | else
12  |   |   | Accept new subset;
13  |   | end
14  | end
15 end
```