

# ML

# Fundamentals



Instituto  
Balseiro

Instituto Balseiro  
26/08/2022



# Announcements

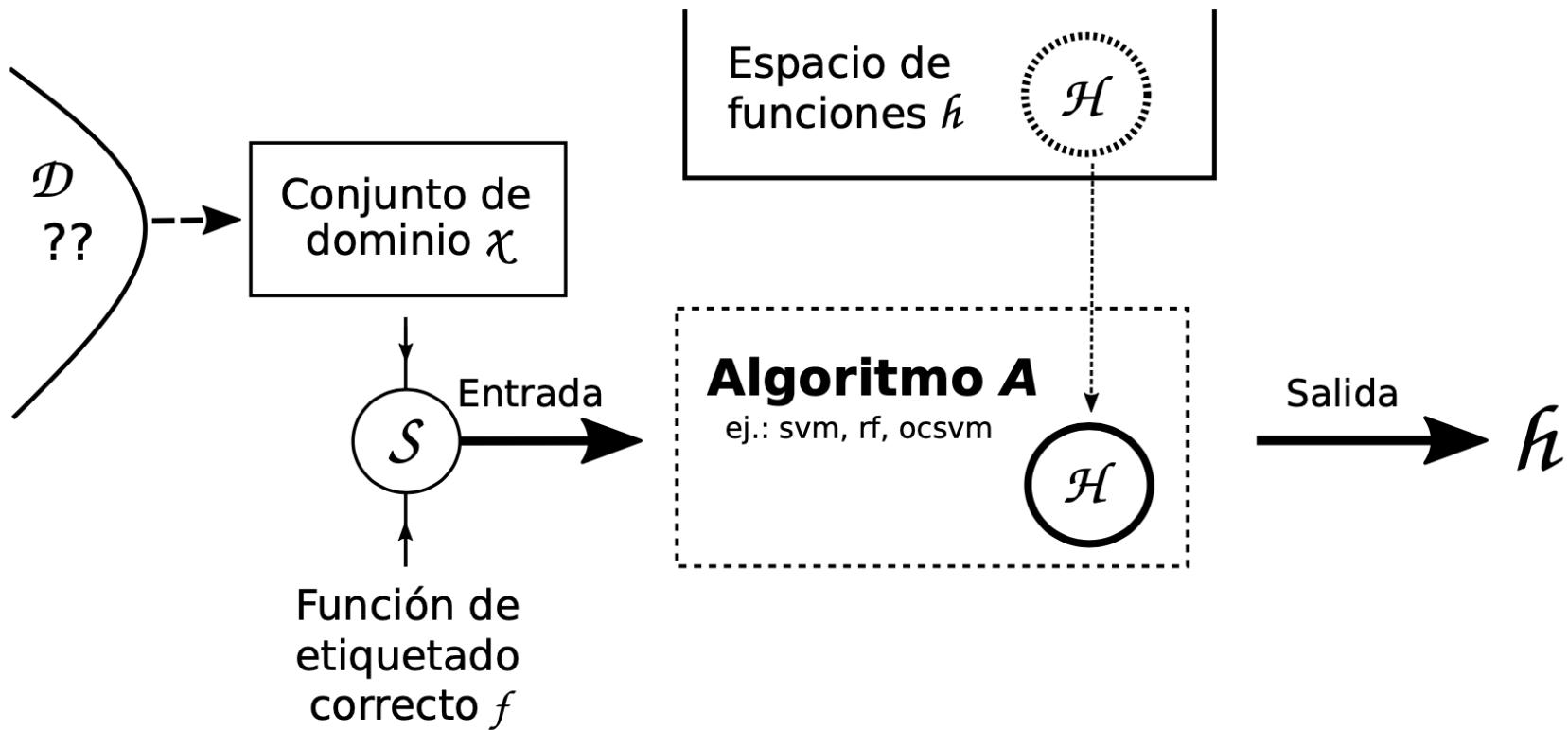
- Inscripciones: hasta el 02/09 !
- Vocacionales: ya cerró. Acerbo y Garrido hablar después conmigo
- Oyentes: confirmar participación
- en Discussions: Cuento corto de Miguel y GPT-3
  - notificaciones de github ➡️ 🚨
- Práctica: Hoy entra la G2, y el lunes la P1
  - P1 se entrega el 09/09.

# Last lecture review



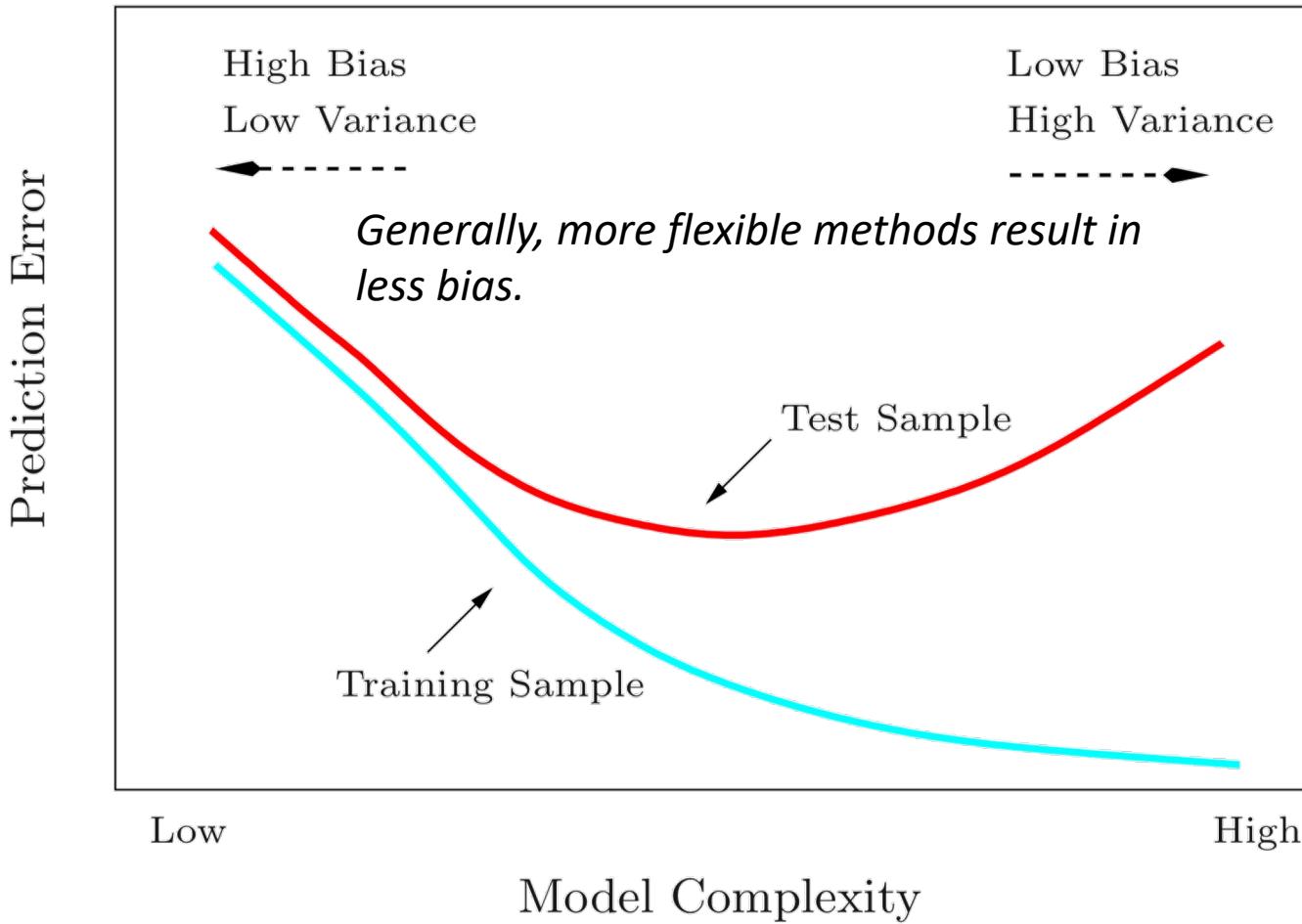
 Permission

# Supervised Learning conceptual model



Técnicas de Aprendizaje Automático Aplicadas a Simulaciones Numéricas de Colisiones de Material Granular Poroso, Daniela N. Rim, Seminario de investigación y/o desarrollo, FCEN, UNCUyo, 2019. Basado en 2014 - Shalev-Shwartz - Understanding machine learning= From theory to algorithms.

# The Bias–Variance Tradeoff



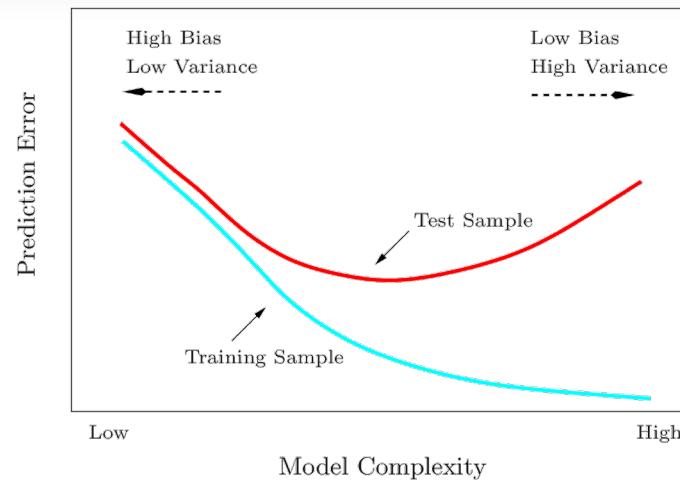
$$E \left[ (y - \hat{f}(x))^2 \right] = \left( \text{Bias} [\hat{f}(x)] \right)^2 + \text{Var} [\hat{f}(x)] + \sigma^2$$

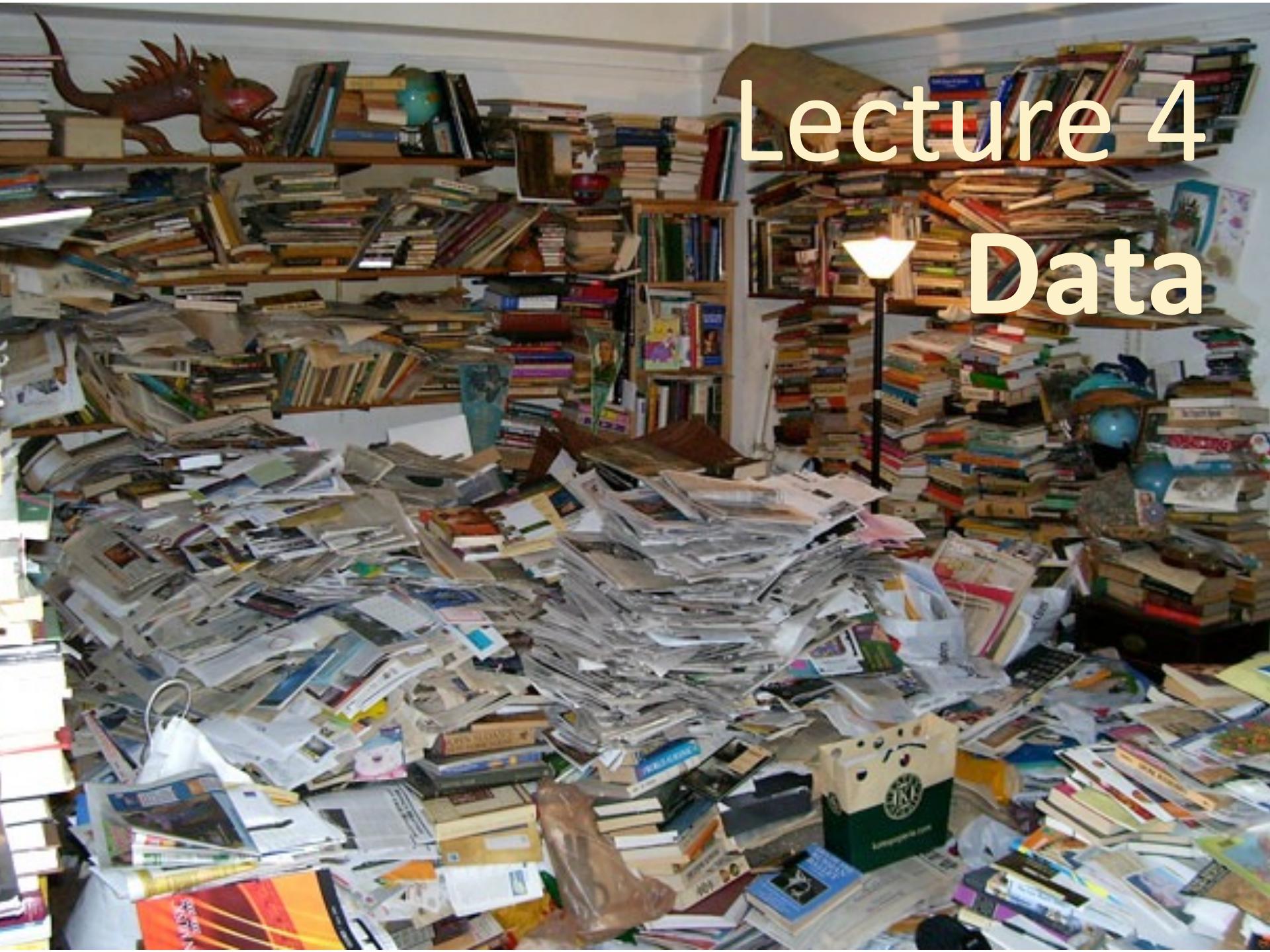
where

$$\begin{aligned} \text{Bias} [\hat{f}(x)] &= E [\hat{f}(x)] - E [f(x)] \\ \text{and} \\ \text{Var} [\hat{f}(x)] &= E[\hat{f}(x)^2] - E[\hat{f}(x)]^2. \end{aligned}$$

# Towards a good model

- Model training: estimating the parameters of the model that lower a given loss function for that particular training data. *Optimization error/Loss/Cost.*
- Model selection: estimating the performance of different models in order to choose the best one: *validation*.
- Model assessment: having chosen a final model, estimating its prediction error (generalization error) on new data. *Test*.



A photograph of a room completely overwhelmed by books. Shelves are stacked high, and the floor is covered in a thick layer of books, papers, and other clutter. A large wooden dragon statue sits on one of the shelves. The lighting is warm, coming from a floor lamp in the center-right.

# Lecture 4

# Data

# ML Fundamentals – Lecture 4

- Data limitations
- Features
- Modeling process and Feature Engineering
- Data cleaning
- Exploratory Data Analysis (EDA)
- Data preprocessing
- Categorical data
- Feature selection

# Data limitations

- Insufficient Quantity of Training Data
- Nonrepresentative Training Data
- Poor-Quality Data
- Irrelevant Features



# Insufficient Quantity of Training Data

- Depending on the algorithm, one might need a huge amount of observations to train properly.
- In general, there is a tradeoff between the algorithm and the amount of data available.
- There are techniques to 'recycle' models (pre-trained models/transfer learning).

# Nonrepresentative Training Data

- A model will be inaccurate if it was trained with data different from the data it's trying to predict.
- Irreducible noise (important with scarce data)
- Sampling bias
- Etc.

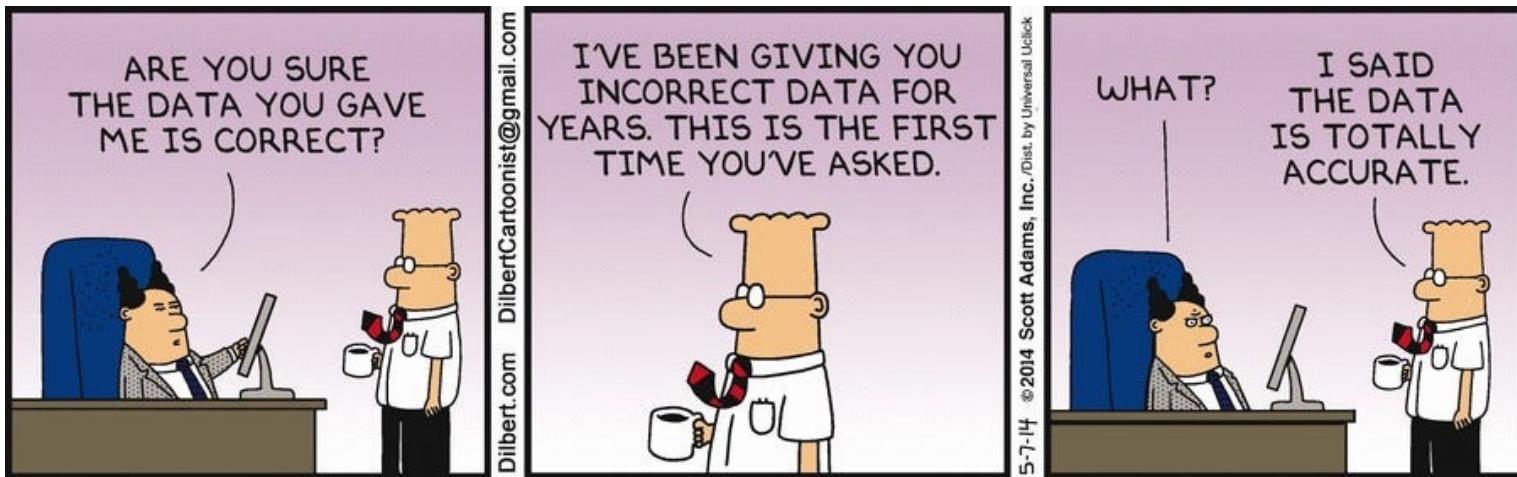


# Poor-Quality Data

- Usually dealt with when data cleaning, but sometimes irremediable.
  - Errors
  - Outliers
  - Other noise (upper/lower, encoding, etc.)
  - Missing data or NA with no clue on its meaning

# Irrelevant Features

- *Garbage in, garbage out*
- Sometimes not trivial to identify which



# What's ML?

A little deeper...

- ML = *representation* + evaluation + optimization

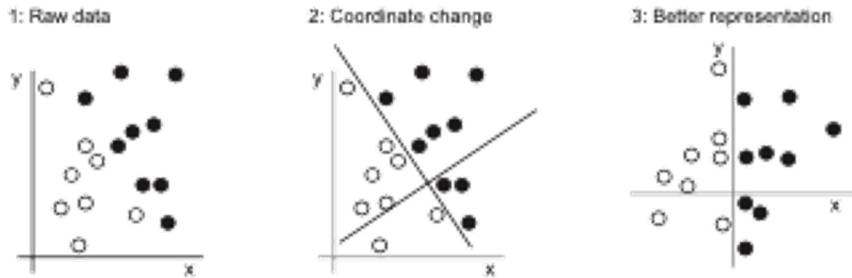


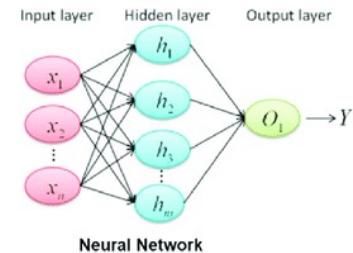
Figure 1.4 Coordinate change

feature representation

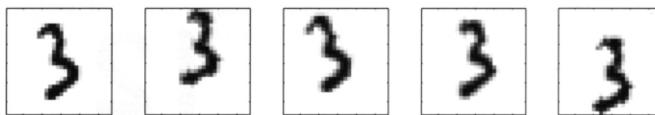
(a)

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \longrightarrow Y$$
$$Y = w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + \cdots + w_n \cdot x_n$$

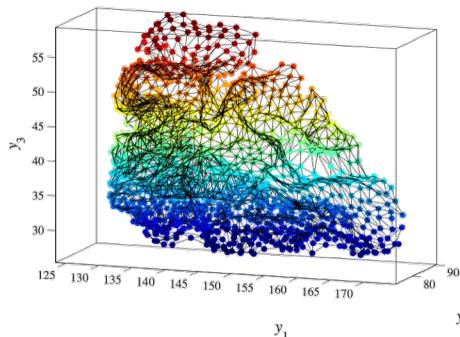
Linear Regression



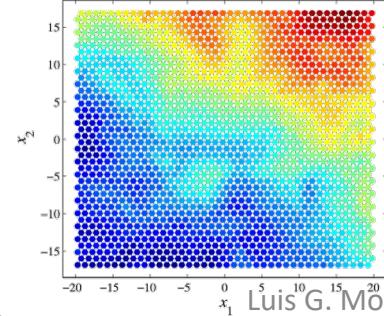
learner space



manifold learning



2006 - Bishop - Pattern Recognition And Machine Learning



Luis G. Moyano - Fundamentos de ML  
Instituto Balseiro

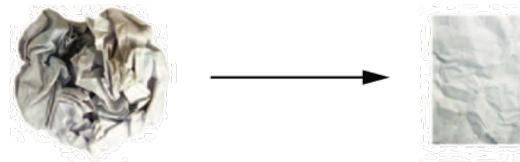


Figure 2.9 Uncrumpling a complicated manifold of data

# Features

- Variables that are discriminant, available and (mostly) informative.

Features									Target
Cust. ID	State	Acct length	Area code	Int'l plan	Voicemail plan	Total messages	Total mins.	Total calls	Churned?
502	FL	124	561	No	Yes	28	251.4	104	False
1007	OR	48	503	No	No	0	190.4	92	False
1789	WI	63	608	No	Yes	34	152.2	119	False
2568	KY	58	606	No	No	0	247.2	116	True

[Feature engineering] is an art like engineering is an art, like programming is an art, like medicine is an art. There are well defined procedures that are methodical, provable and understood.

# Feature Engineering

*aka data wrangling or data munging*

In the words of Pedro Domingos:<sup>84</sup>

... some machine learning projects succeed and some fail. What makes the difference? Easily the most important factor is the features used.

In the words of Jason Brownlee:<sup>51</sup>

[Feature engineering] is an art like engineering is an art, like programming is an art, like medicine is an art. There are well defined procedures that are methodical, provable and understood.

In the words of Andrew Ng:<sup>337</sup>

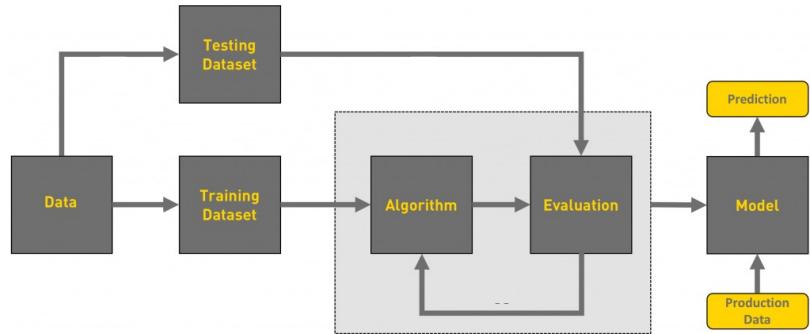
Coming up with features is difficult, time-consuming, requires expert knowledge. “Applied machine learning” is basically feature engineering.

In the words of Yoshua Bengio:<sup>337</sup>

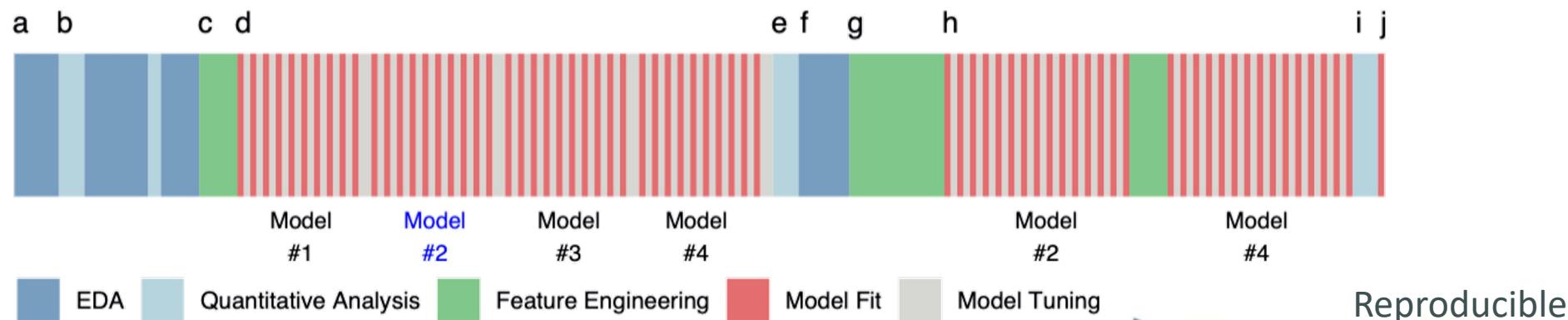
Good input features are essential for successful ML. Feature engineering is close to 90% of effort in industrial ML.

# Model != Modeling Process

1. iterative
2. heuristic/creative
3. documented



A typical modeling process, comprising several ML cycles:



2019 - Book - Kuhn - Feature Engineering and Selection=A Practical Approach for Predictive Models, Fig 1.4, p9

Reproducible  
Research



# Feature Engineering

aka data wrangling or data munging

- **EDA:** exploratory analysis of the raw data.
  - identifying good features and expanding them
- **Error analysis:** which features are hurting performance more than helping?
  - identifying redundant/uninformative features and dropping them.

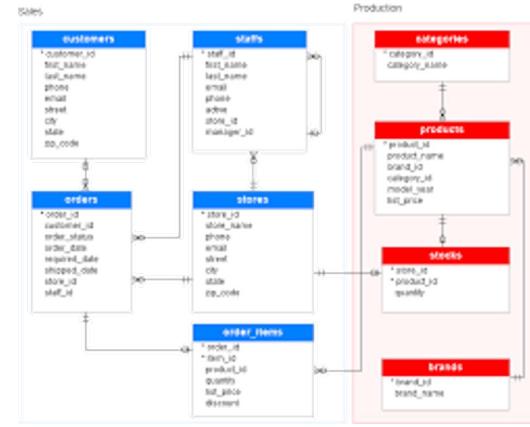
# Domain Knowledge



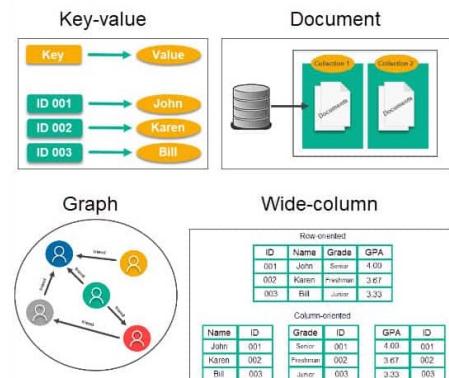
Luis G. Moyano - Fundamentos de ML -  
Instituto Balseiro

# Feature types

Type	Description / Note	Example
<b>Binary</b>	One of two values <i>simplest</i>	Did the customer leave a tip? Did the turbine make noise? Did the student turn in the exam before the end of the allocated time? Does the person like cats? Did the licence plate have a match in the DMV database?
<b>Categorical</b>	One among many values <i>values known beforehand</i>	Colour of the car, Brand of the car, Credit card type, Number of bedrooms ( <i>also a number</i> ) Type of mattress ( <i>also mattress surface area</i> )
<b>Discrete</b>	Can be mapped to the integers <i>order is important</i>	Number of pizzas eaten, Times visited the gas pump, Times the turbine got serviced last year, Steps walked last week ( <i>compare: distance walked, which is continuous</i> ), Number of people living in the house.
<b>Continuous</b>	Can be mapped to a real number <i>representation issues</i>	Engine temperature, Latitude, Longitude, Speech length, Colour intensity at centre of image, Distance walked, Temperature at centre of dish ( <i>as measured by an IR camera</i> )
<b>Complex</b>	Records, lists, sets <i>challenging</i>	Date (year, month, day; note that it can be expressed as number of days from a fixed date), Product name ( <i>it might include brand, product type, colour, size and variant</i> ), Location (latitude and longitude), Complaint (a sequence of characters), Countries visited (it is a set of categories)



**Structured:** Spreadsheets, csv, arrays, dataframes, lists, SQL



**Unstructured:** twitter, JSON, MongoDB, NoSQL

# Dataset types

- Structured vs unstructured data
- File types (txt, CSV, SQLite, JSON, BigQuery)
- Sizes (Small, Medium, Large, Big data)
- Licenses (Creative Commons, GPL, Other Database, Other)

# Missing values

blanks, NaNs, NAs, <special holder>, etc.

PassengerId	Survived	Pclass	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	Male	22	1	0	A/5 21171	7.25		S
2	1	1	Female	38	1	0	PC 17599	71.2833	C85	C
3	1	3	Female	26	0	0	STON/O2. 3101282	7.925		S
4	1	1	Female	35	1	0	113803	53.1	C123	S
5	0	3	Male	35	0	0	373450	8.05		S
6	0	3	Male		0	0	330877	8.4583		Q

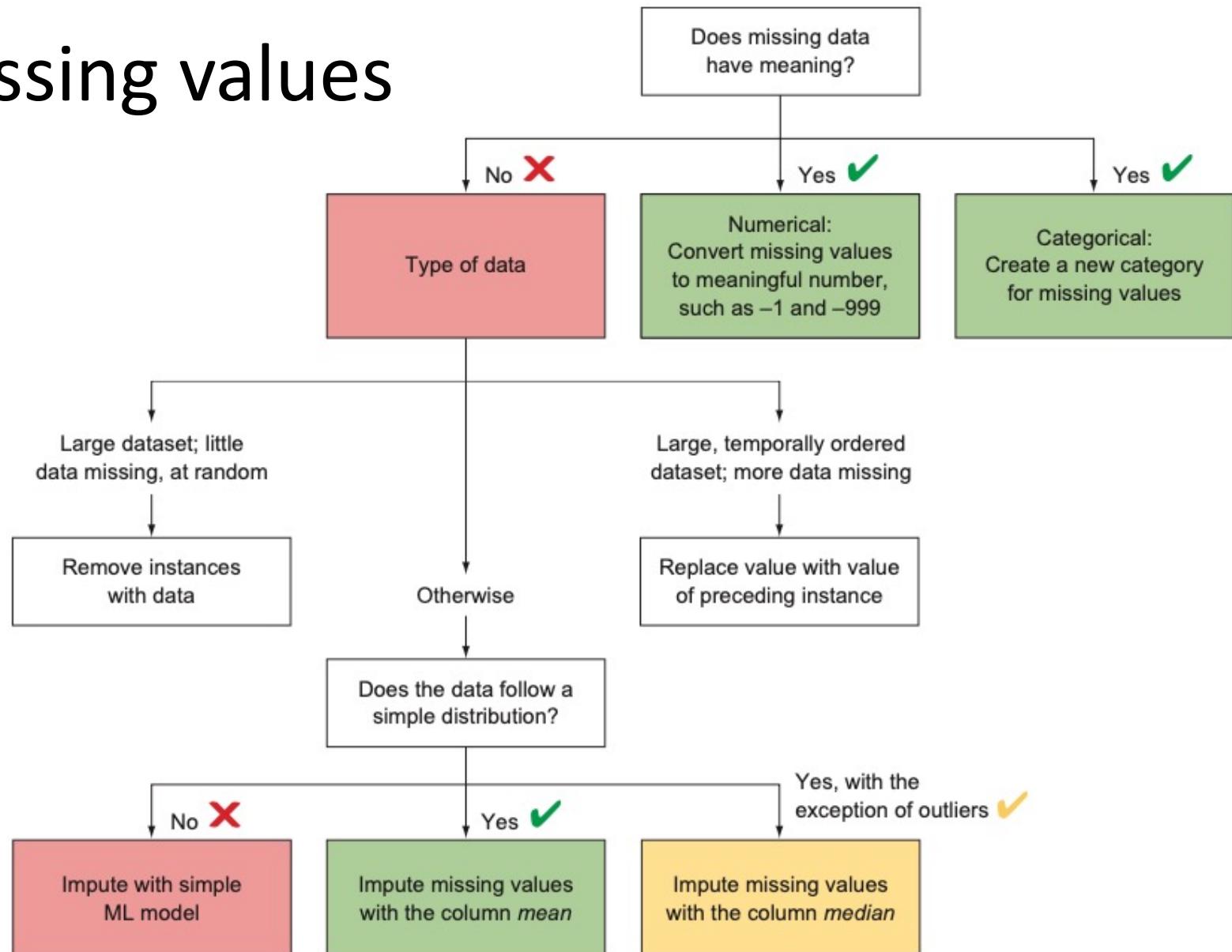
Missing values

- Causes: structural, random (MCAR, MAR), specific
- Is the NA informative?
- Should NAs stay? Or should they go?
- Imputation



<https://scikit-learn.org/stable/modules/impute.html#impute>

# Missing values



Henrik Brink, Joseph Richards, Mark Fetherolf - Real-World Machine Learning-Manning Publications (2016)  
Josse et al. *On the consistency of supervised learning with missing values*. 2019. ([hal-02024202v2](https://hal-02024202v2))

# EDA: Exploratory data analysis

- Load and inspect data
- Serialization formats
  - JSON
  - HDF5
  - Pickle
    - most objects!
    - easy to use!
    - ... but be careful

```
import numpy as np
import pandas as pd
import pickle

>>> pdDf = pd.read_csv('100 Sales Record.csv')
>>> pdDf.head()

with open('test.pkl','wb') as f:
    pickle.dump(pdDf, f)

with open("test.pkl", "rb") as f:
    d4 = pickle.load(f)

>>> d4.head()
```

**Warning:** The `pickle` module is **not secure**. Only unpickle data you trust.

It is possible to construct malicious pickle data which will **execute arbitrary code during unpickling**. Never unpickle data that could have come from an untrusted source, or that could have been tampered with.

Consider signing data with `hmac` if you need to ensure that it has not been tampered with.

Safer serialization formats such as `json` may be more appropriate if you are processing untrusted data. See [Comparison with json](#). <https://docs.python.org/3/library/pickle.html>

# EDA: Exploratory data analysis

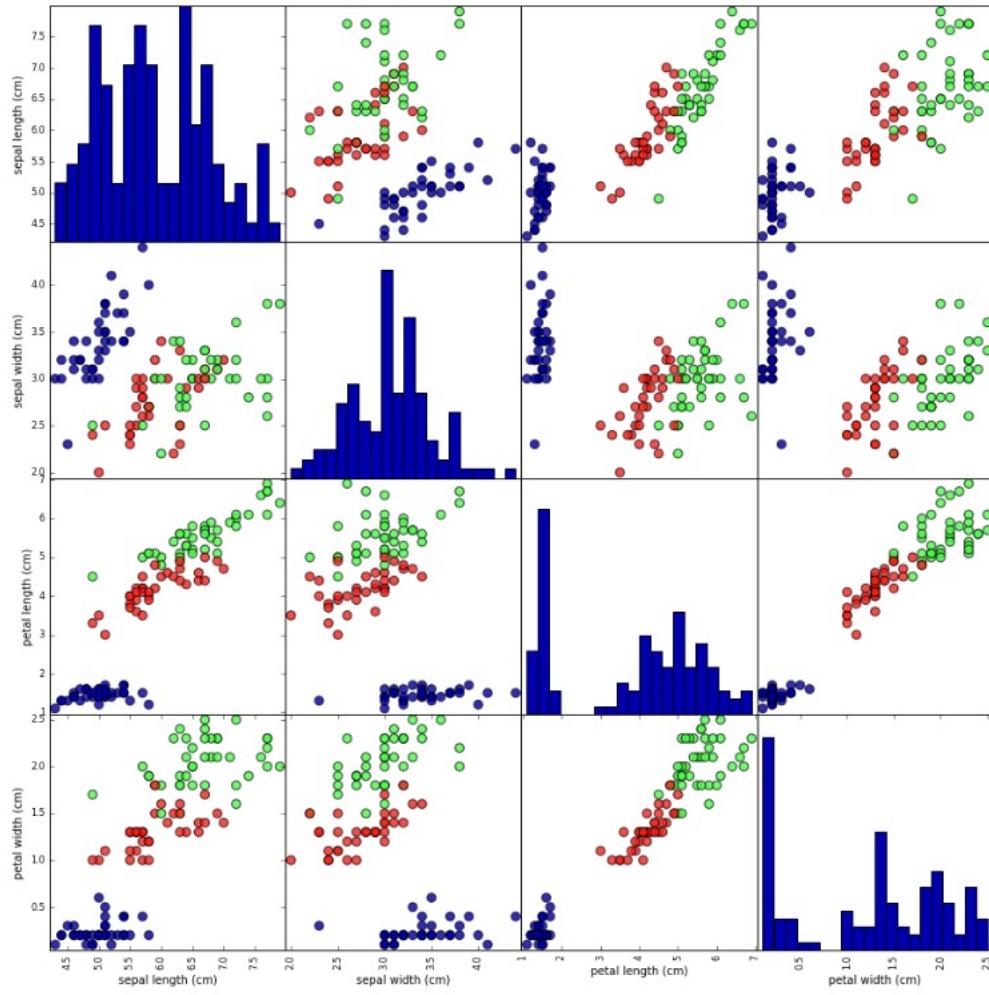
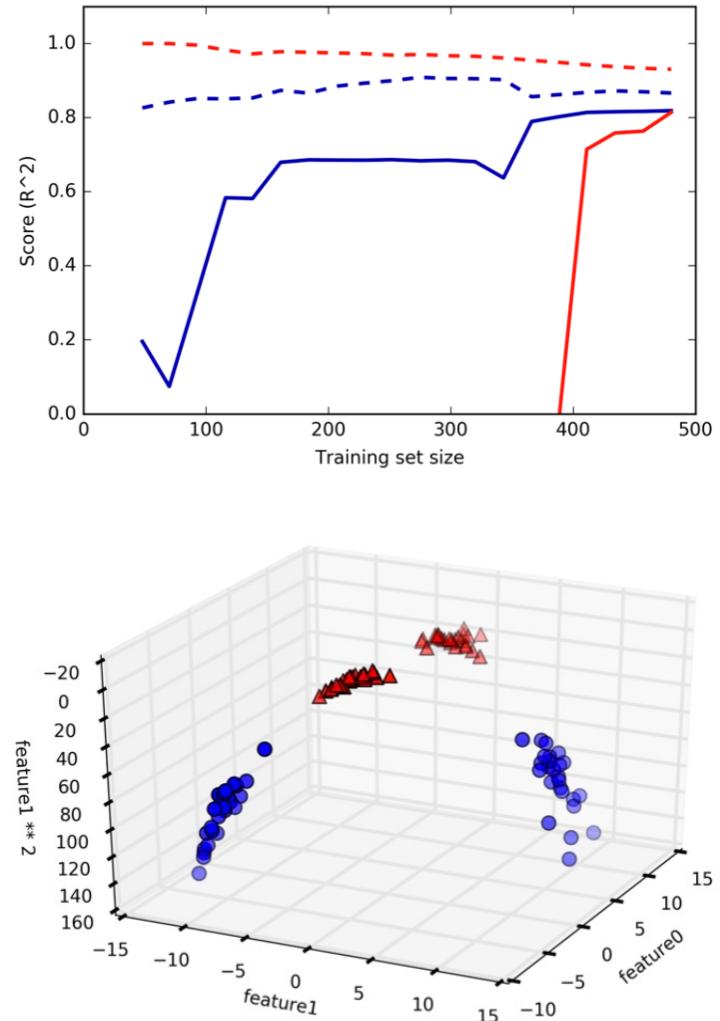
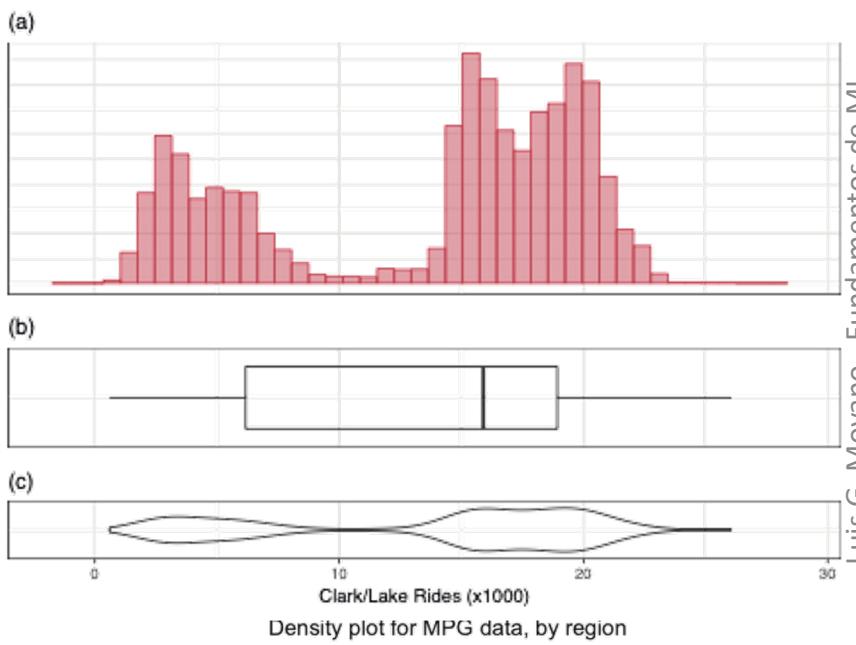


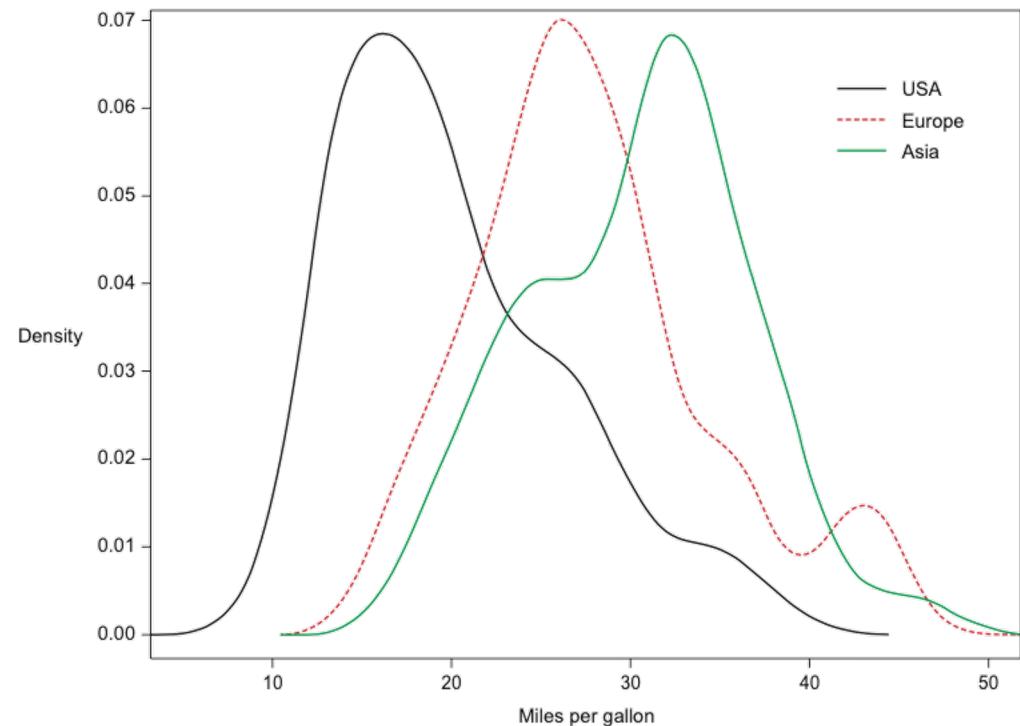
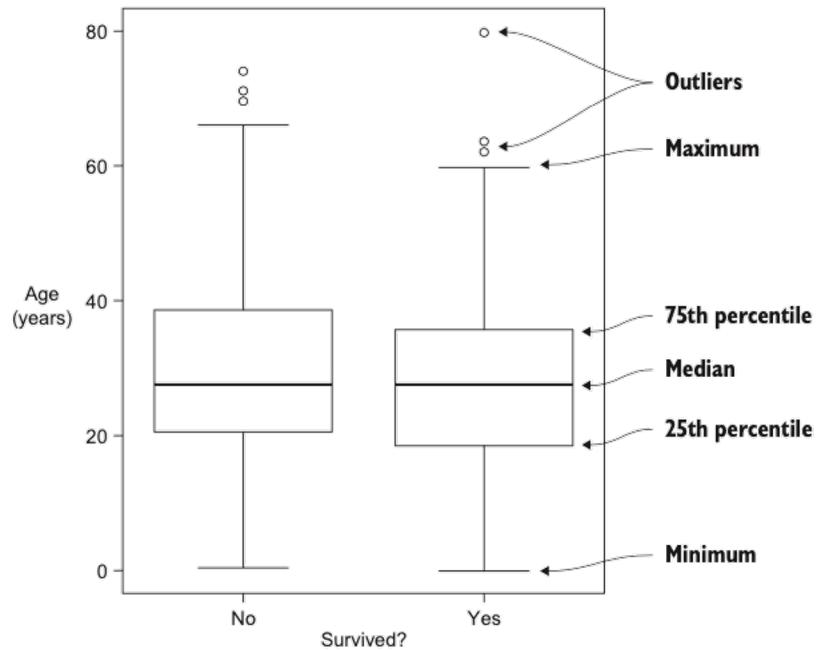
Figure 1-3. Pair plot of the Iris dataset, colored by class label



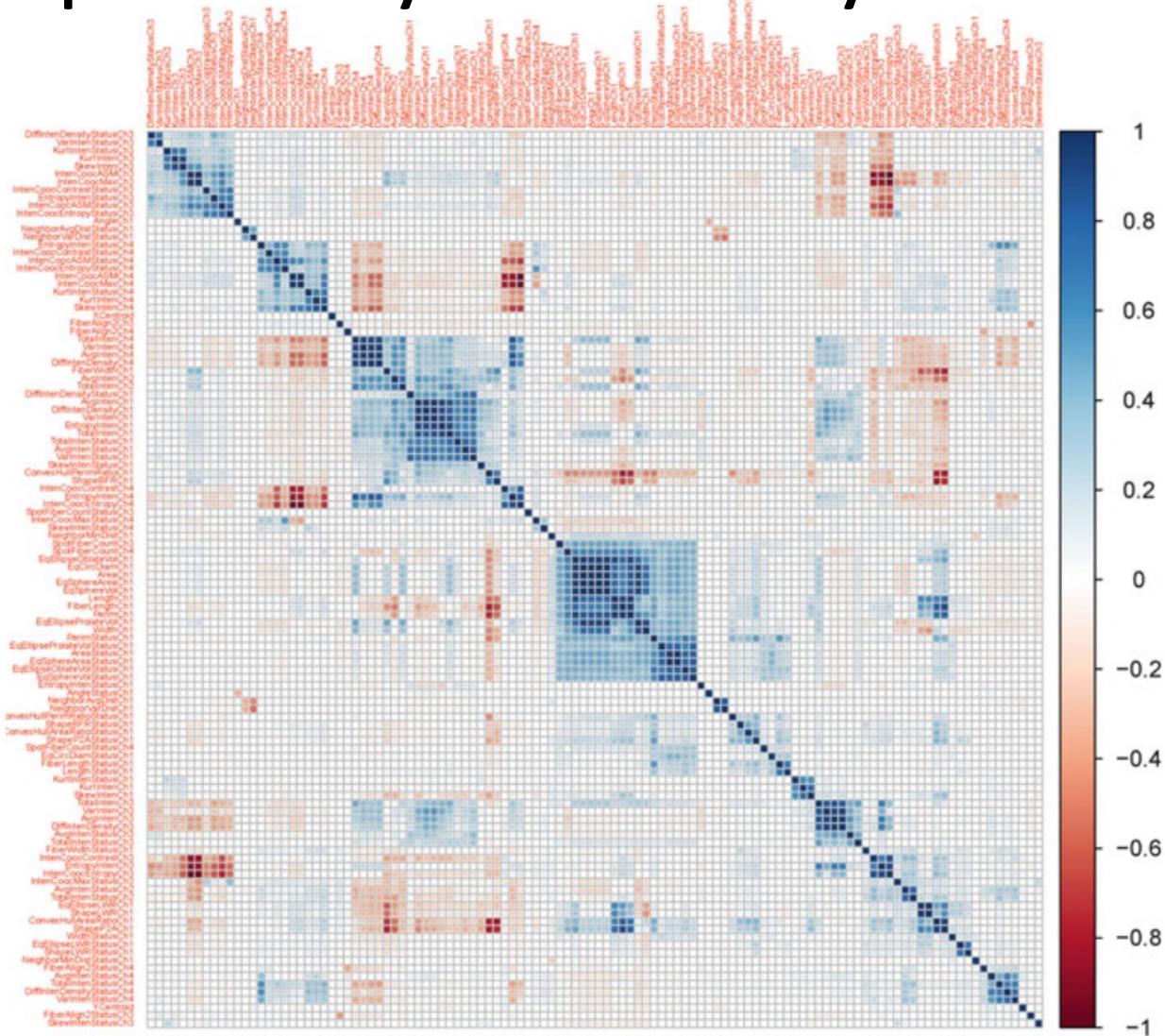
# EDA: Comparing data



Box plot for Titanic data: Passenger age vs. survival



# EDA: Exploratory data analysis



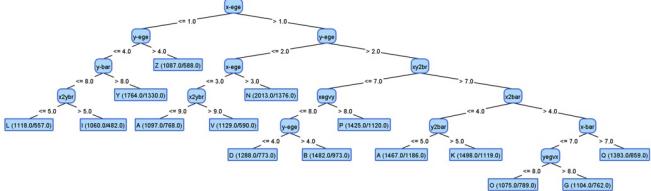
# (Multi)collinearity

- Redundant predictors frequently add more complexity to the model than information they provide to the model.
- In some models, may lead to unstable behavior, numerical errors or poor interpretability.
- Other models (e.g. trees) deal well with this.
- For prediction, it's generally **not** an issue.



# No-free-lunch theorem

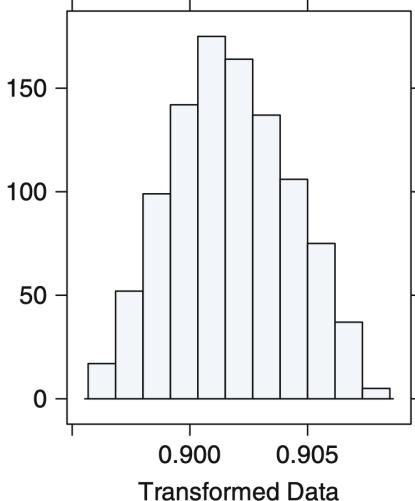
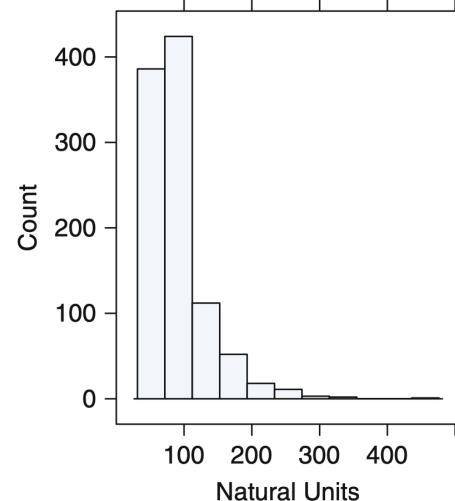
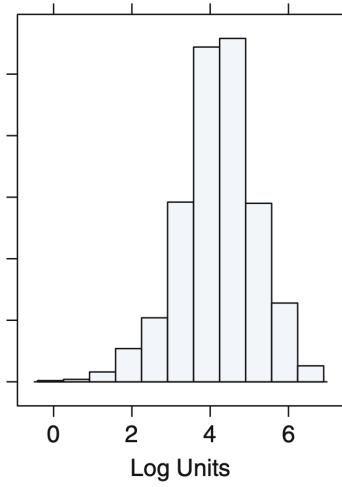
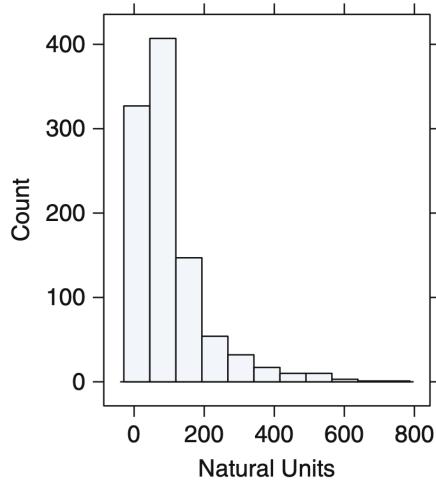
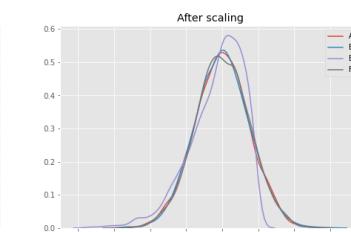
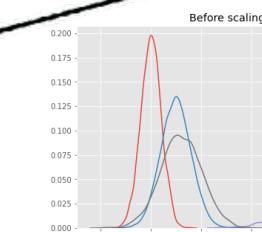
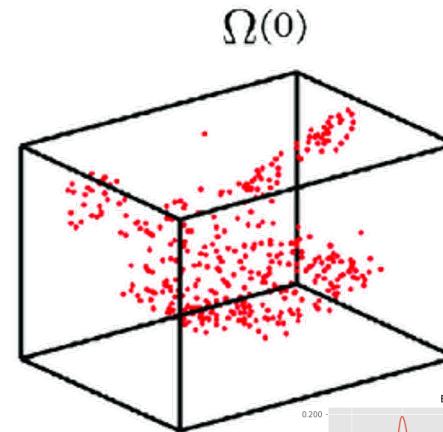
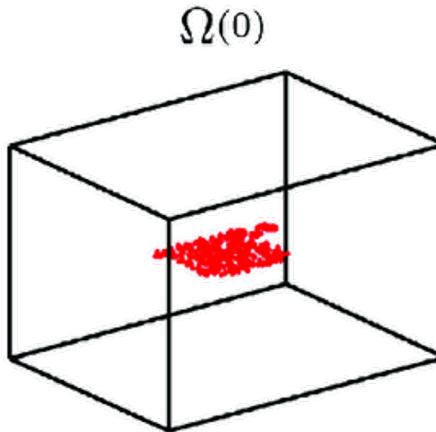
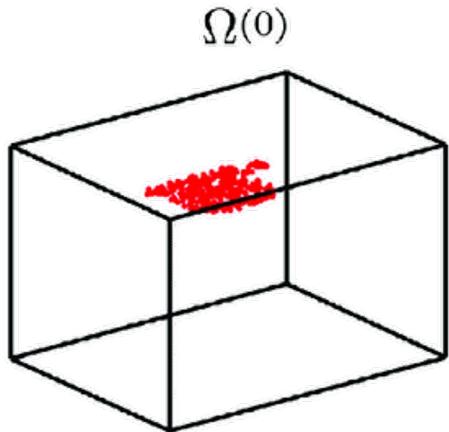
- Always check your assumptions before relying on a model or search algorithm.
  - There is no “super algorithm” that will work perfectly for all datasets.



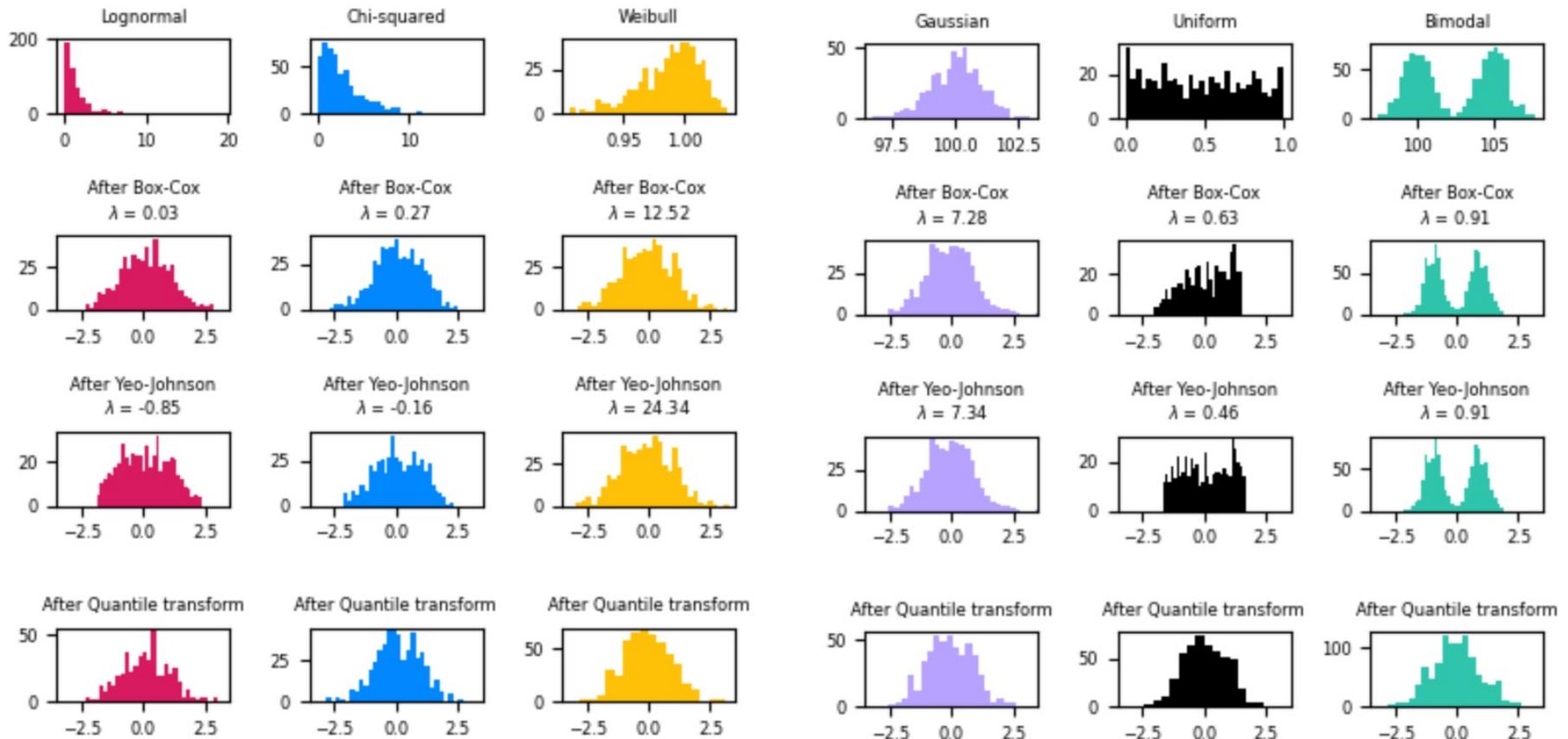
# Data pre-processing

- Normalization: centering and scaling
  - [0, 1]?
    - <http://www.faqs.org/faqs/ai-faq/neural-nets/part2/section-16.html>
    - Whitening (decorrelation, identity covariance)
- dimensionality reduction (PCA, t-SNE, word2vec)
- outlier removal
- near-zero variance (bad, but check target correlation)
- collinear features

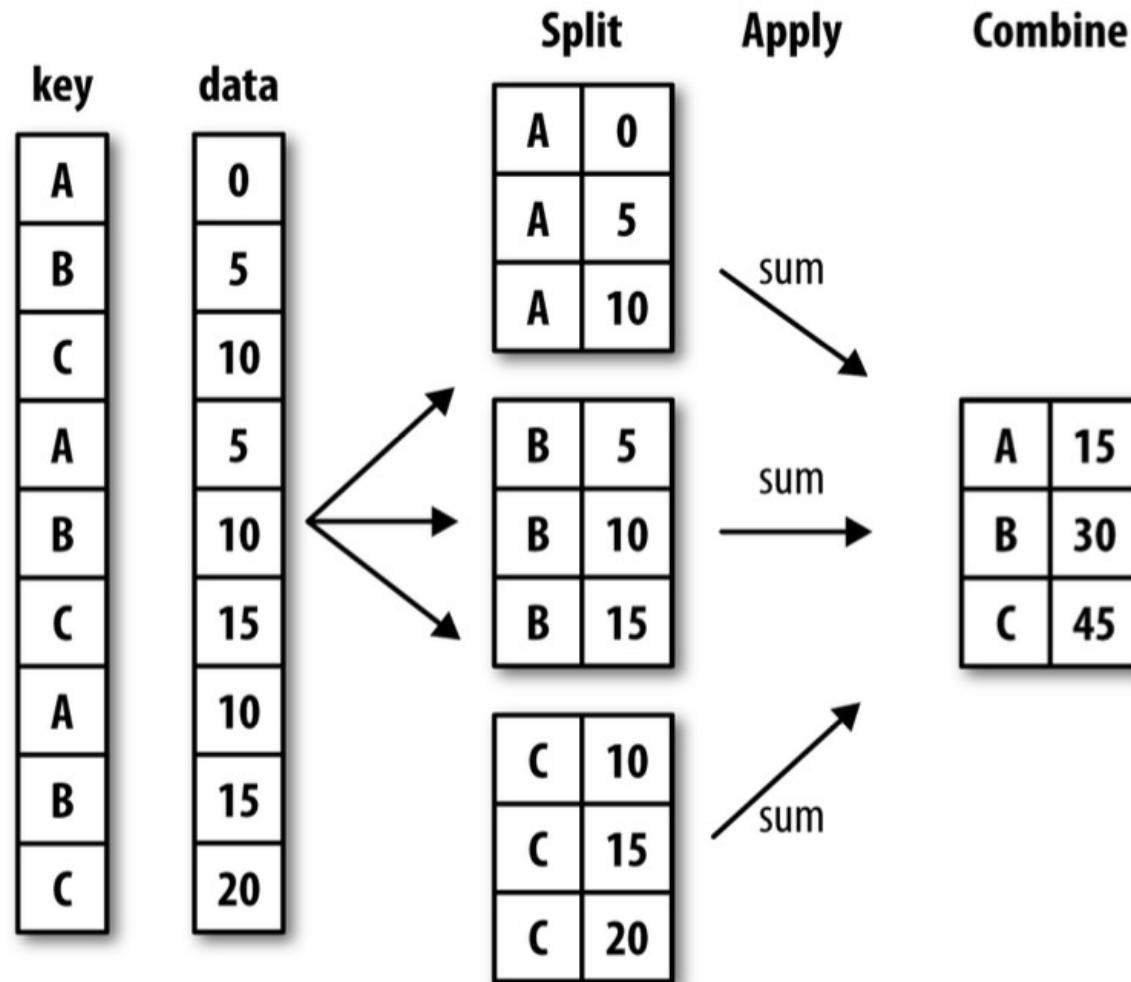
# EDA: Exploratory data analysis



# Box-Cox, Yeo-Johnson, Quantile



# Group data (Split + Apply + Combine)

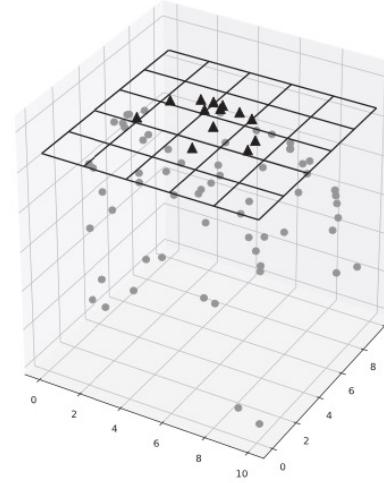
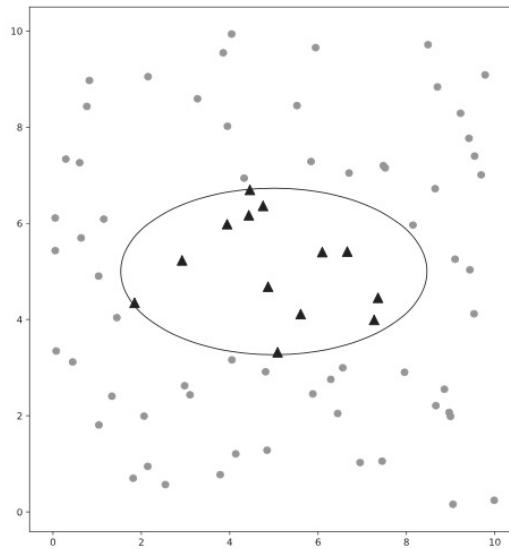


# Other transformations

- pivoting
- decomposing
- discretizing
- kernel transformations

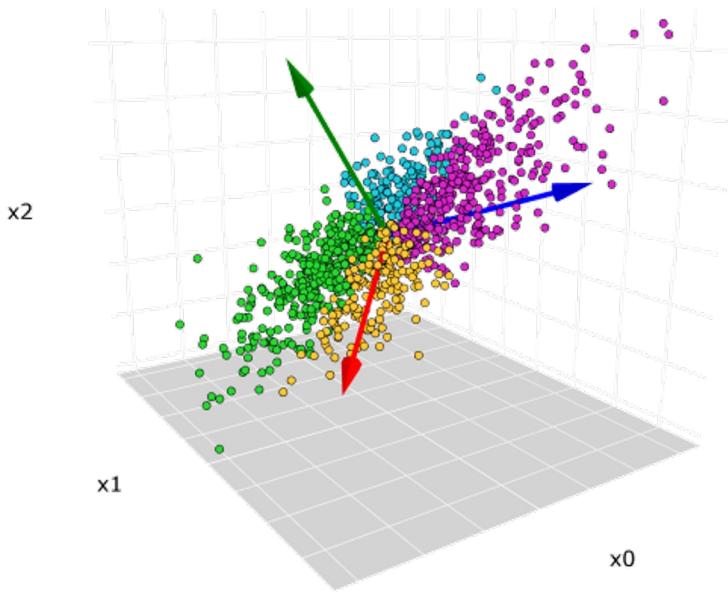
Original Data			Data after Pivot	
ID	Activity	Target	ID	Read?
1	read	n	1	y
1	save	n	2	y
1	search	y	3	n
2	save	n	4	n
2	read	n	5	y
3	save	n		
4	search	n		
5	read	y		
5	search	y		

$$K_{\text{rbf}}(\vec{x}, \vec{y}) = \exp(-\gamma \|\vec{x} - \vec{y}\|^2)$$

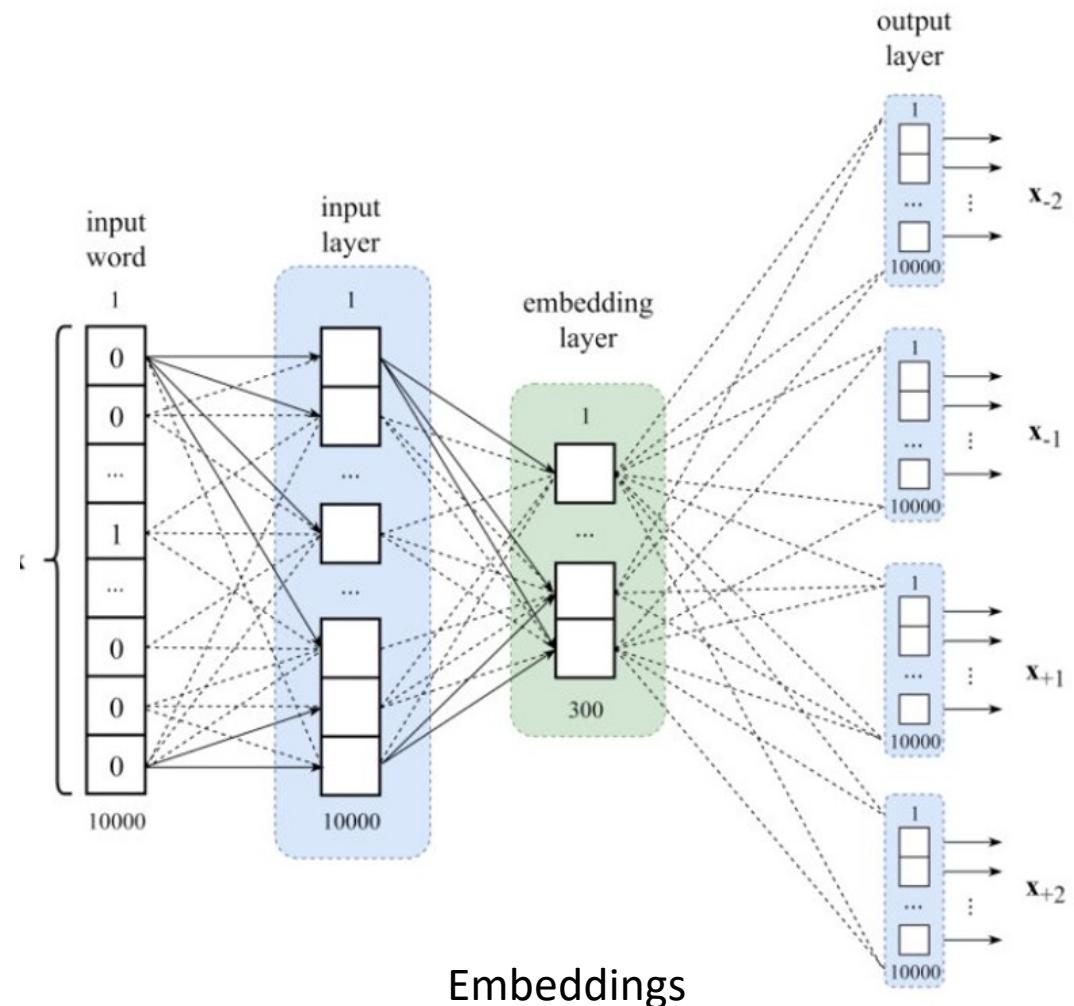


# Learning representations

## Dimensionality reduction and embeddings



PCA

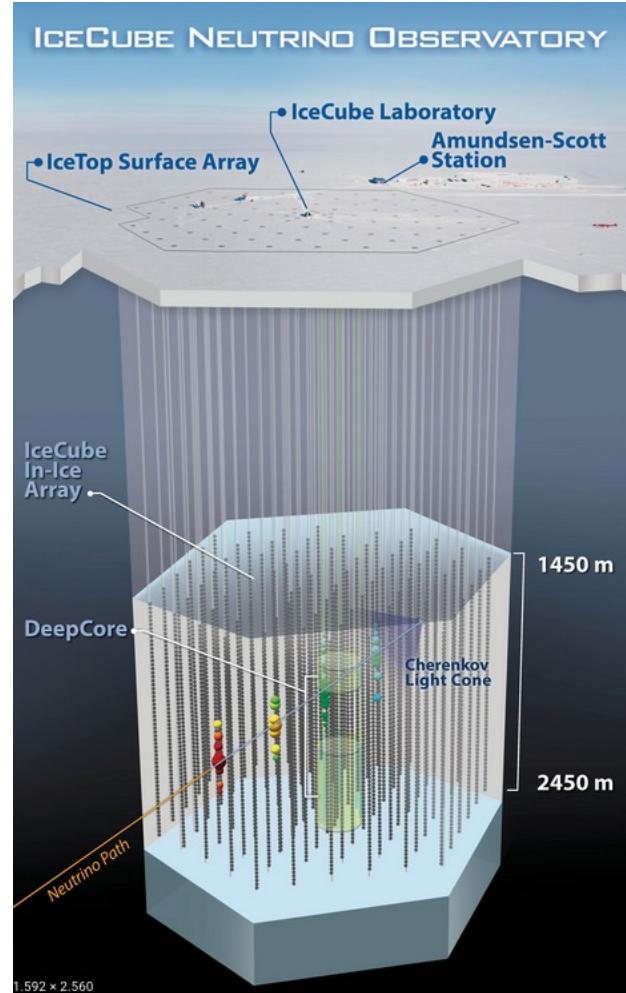
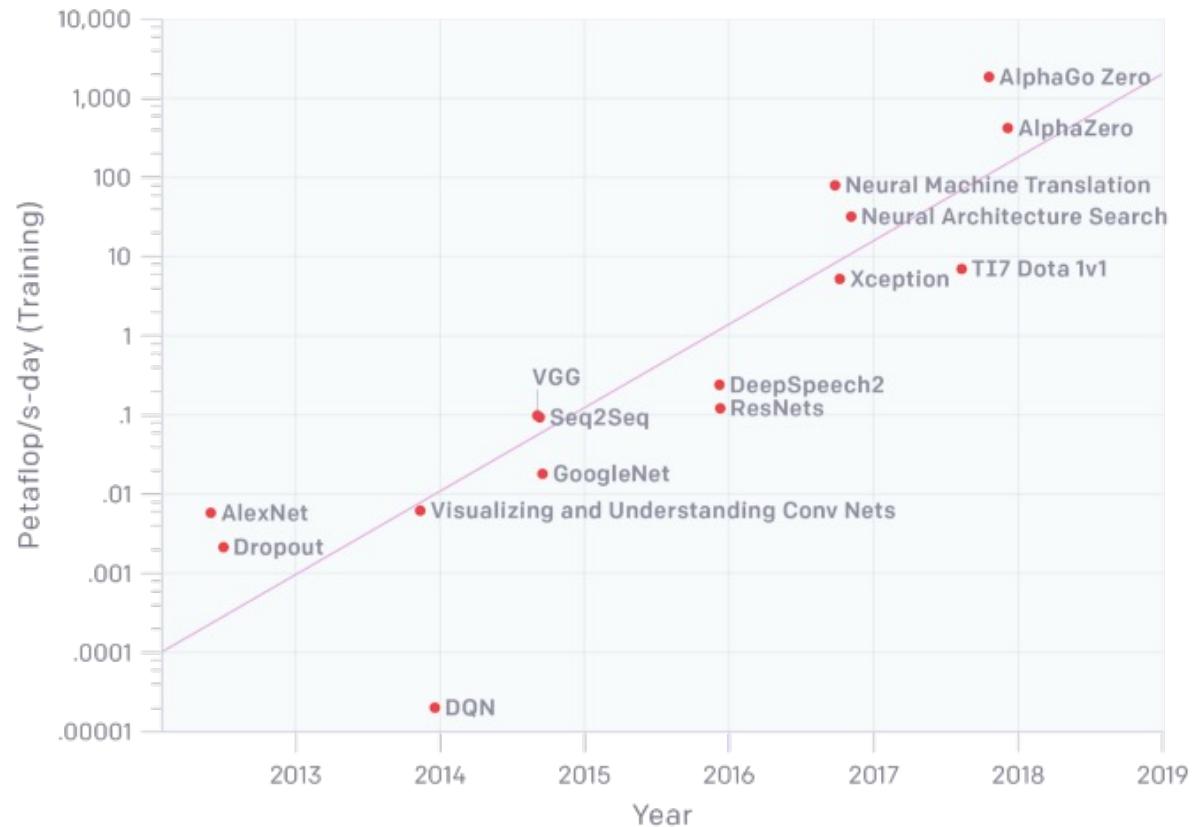


# Deep Learning doesn't need FE

- DL offers the possibility to eliminate FE by trading large amounts of training data for the painstaking process of feature improvement.  
2020 Duboue p127
- Pero, pero, pero... Good features solve problems using fewer resources.

# Big Data

AlexNet to AlphaGo Zero: A 300,000x Increase in Compute



# Categorical features

The diagram illustrates the concept of categorical features. At the top, a small table shows two rows of data with columns: Person, Name, Age, Income, and Marital status. The 'Name' and 'Marital status' columns are highlighted in light blue, representing categorical features. Below this, a large table displays six rows of passenger data with columns: PassengerId, Survived, Pclass, Gender, Age, SibSp, Parch, Ticket, Fare, Cabin, and Embarked. A bracket labeled 'Categorical features' points from the small table to the 'Gender' column in the large table, indicating that the categories in the small table correspond to the categorical feature 'Gender' in the larger dataset.

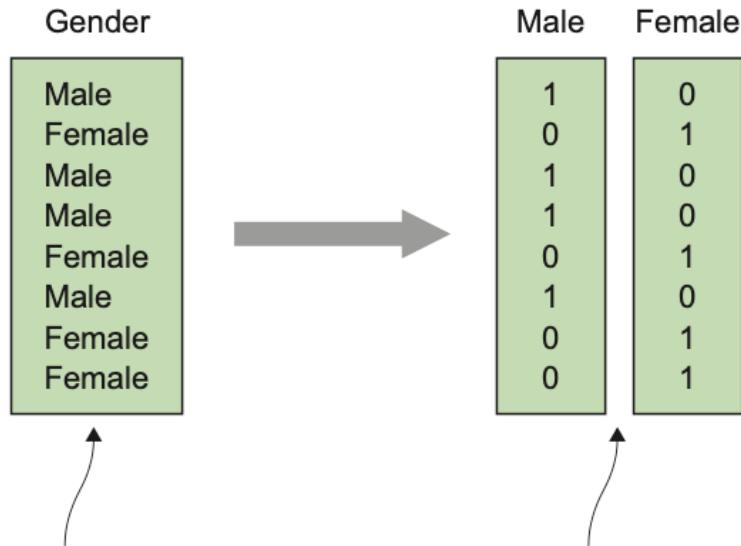
Person	Name	Age	Income	Marital status
1	Jane Doe	24	81,200	Single
2	John Smith	41	121,000	Married

PassengerId	Survived	Pclass	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	Male	22	1	0	A/5 21171	7.25		S
2	1	1	Female	38	1	0	PC 17599	71.2833	C85	C
3	1	3	Female	26	0	0	STON/O2. 3101282	7.925		S
4	1	1	Female	35	1	0	113803	53.1	C123	S
5	0	3	Male	35	0	0	373450	8.05		S
6	0	3	Male		0	0	330877	8.4583		Q

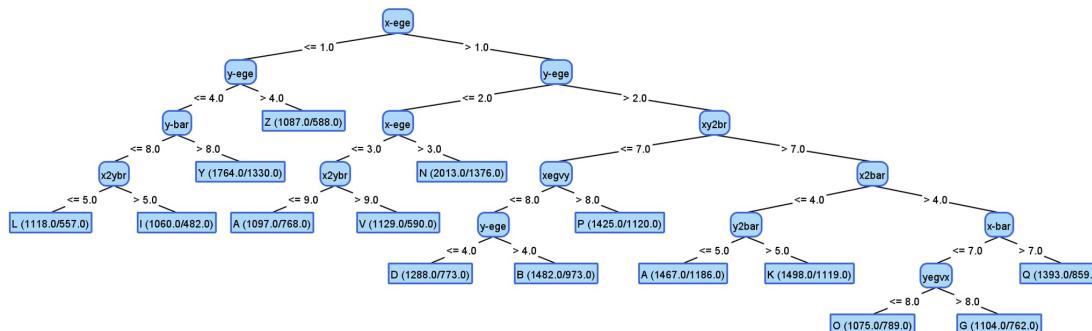
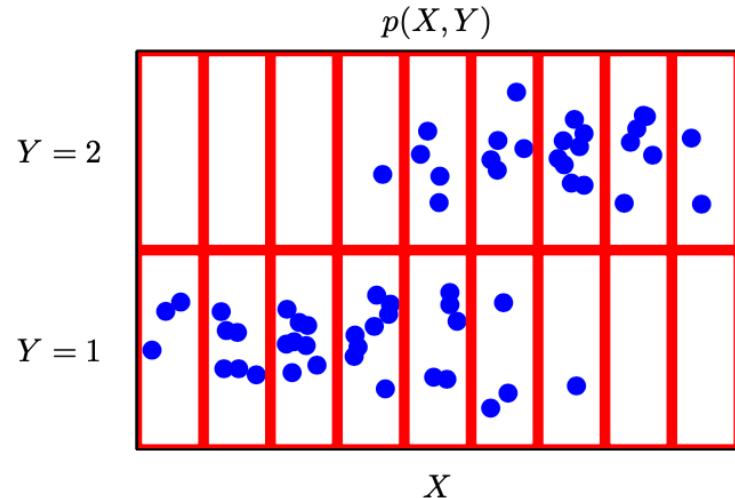
<https://scikit-learn.org/stable/modules/preprocessing.html#encoding-categorical-features>

# One-hot encoding, aka dummy variables



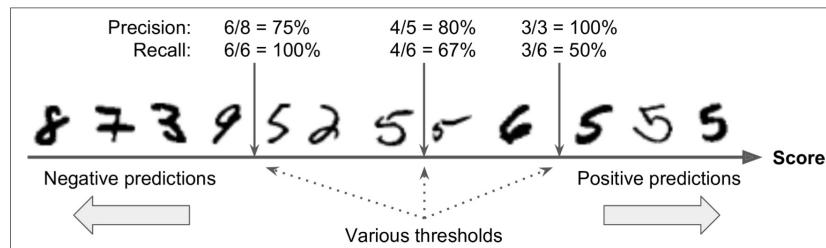
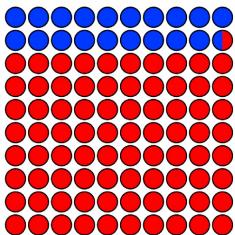
Categorical feature  
with two categories:  
“Male” and “Female”

Categorical feature  
converted to two binary  
features: one per category

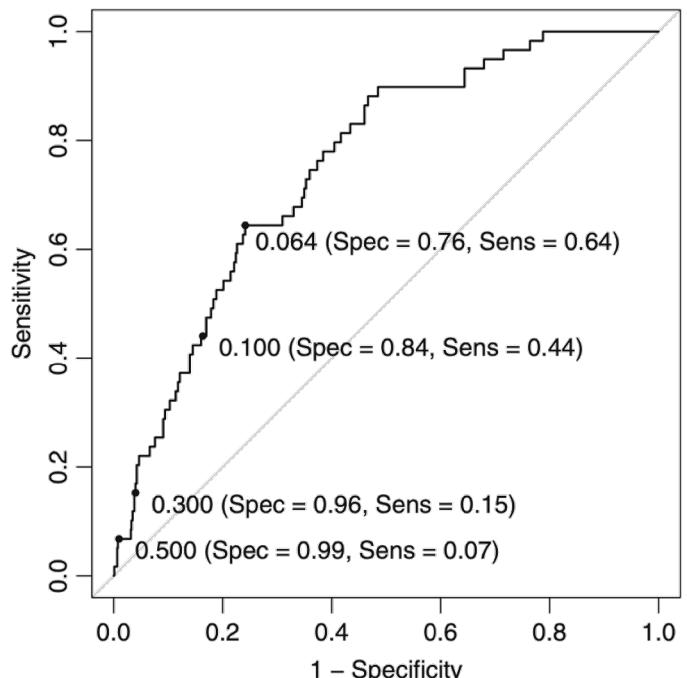
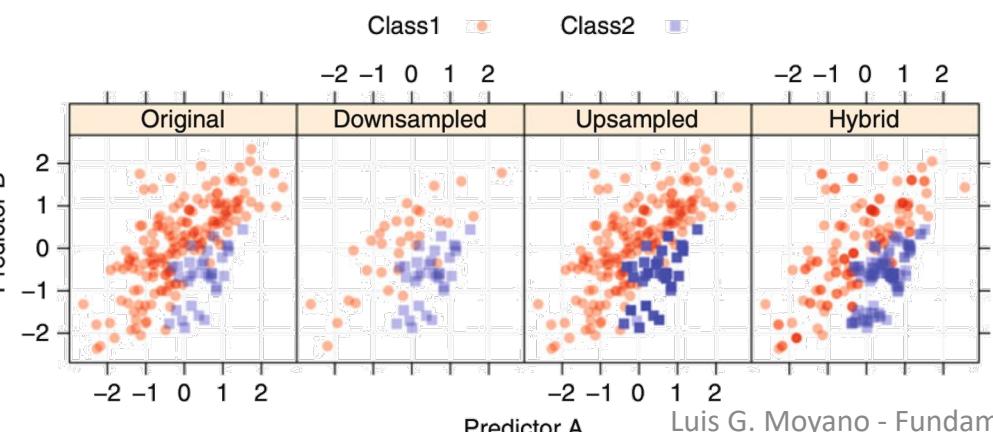


# Severe Class Imbalance

1. Alternate Cutoffs
2. Sampling methods
  - SMOTE
  - ROSE
3. Cost-sensitive training
  - error/class penalty



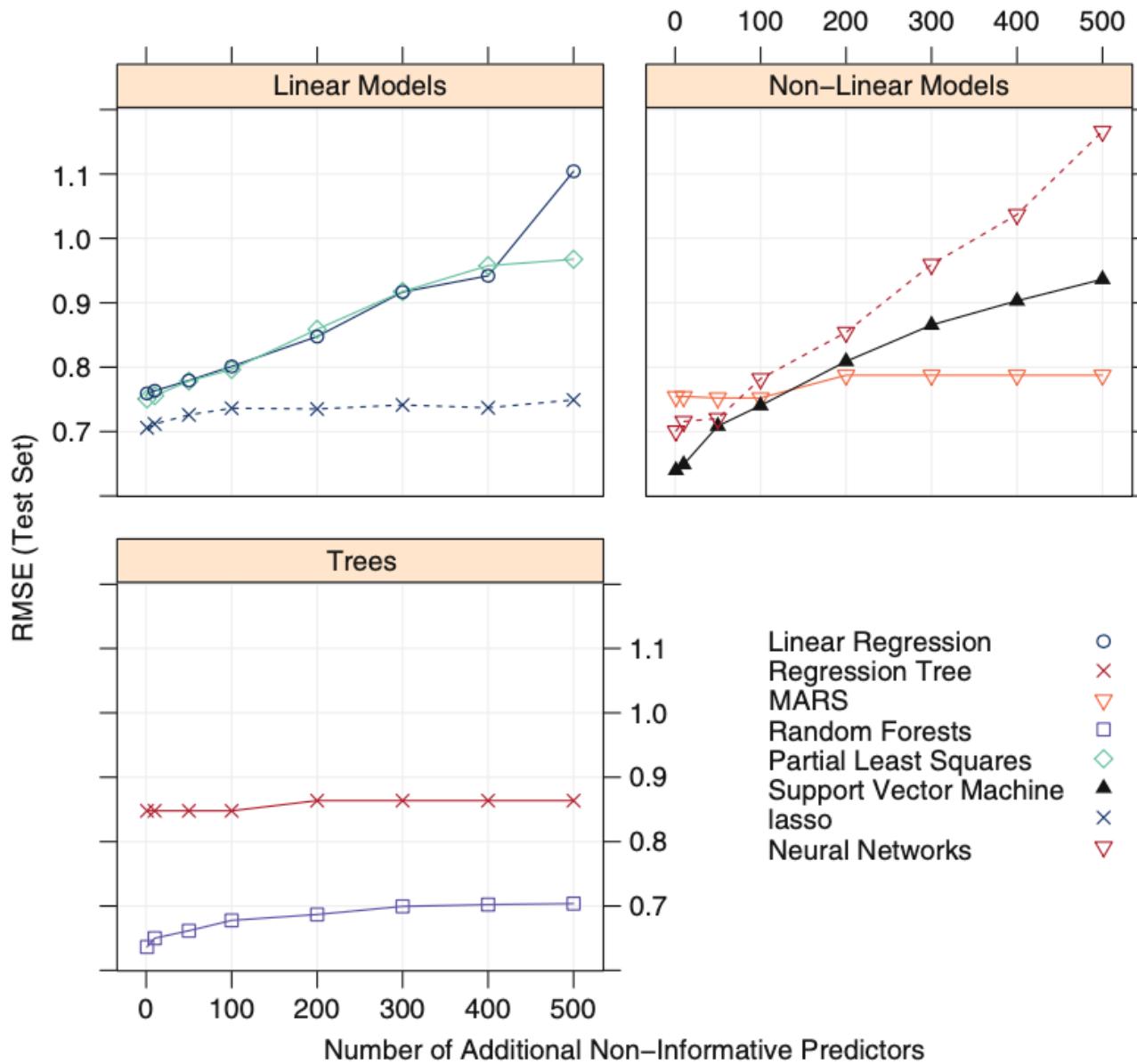
0.50 Cutoff		0.064 Cutoff	
	Insurance no insurance		Insurance no insurance
Insurance	11	19	71
Noinsurance	105	1827	45
			1,405



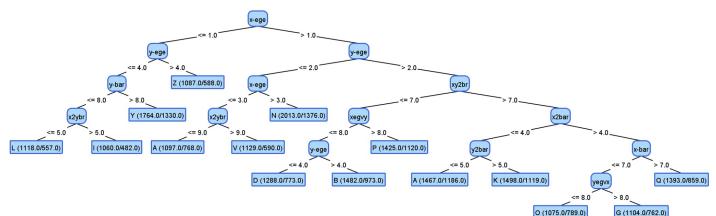
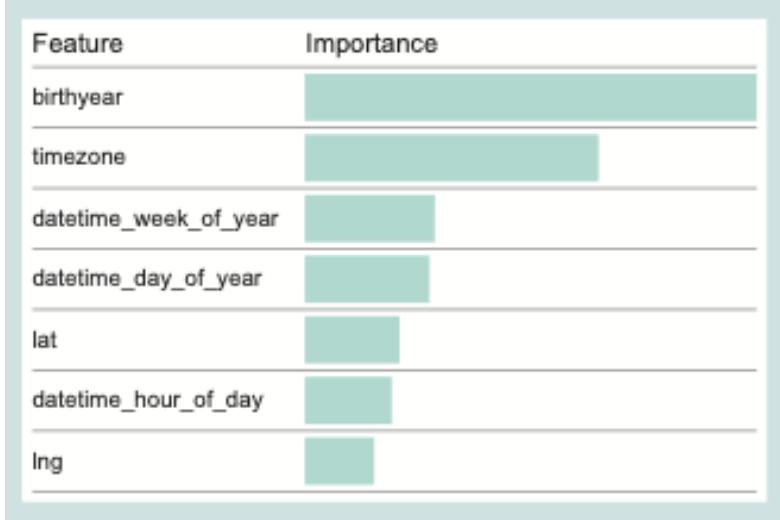
# Feature engineering workflow (#ponele)

- Include all the features that you suspect to be predictive of the target variable.  
Fit an ML model. If the accuracy of the model is sufficient, stop.
- Otherwise, expand the feature set by including other features that are less obviously related to the target.  
Fit an ML model. If the accuracy of the model is sufficient, stop.
- Otherwise, starting from the expanded feature set, run an ML *feature selection algorithm* to choose the best, most predictive subset of your expanded feature set.

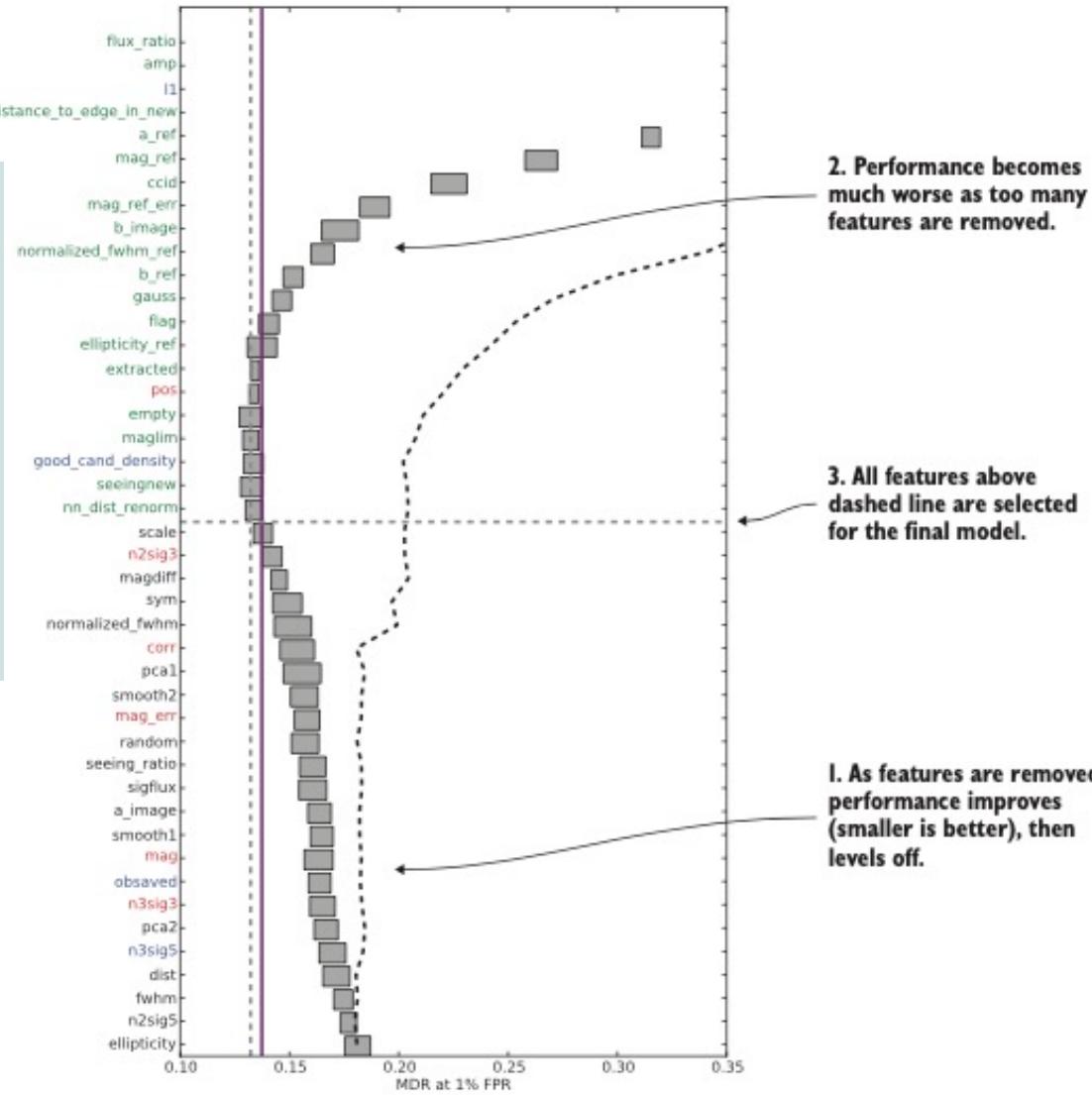
# Non-informative Predictors



# Feature selection



[https://scikit-learn.org/stable/modules/feature\\_selection.html](https://scikit-learn.org/stable/modules/feature_selection.html)



# Lecture 4 – Data – Summary

- Data limitations
- Features
- Modeling process and Feature Engineering
- Data cleaning
- Exploratory Data Analysis (EDA)
- Data preprocessing
- Categorical data
- Feature selection

Next: Linear algorithms