

The Abdus Salam International Centre for Theoretical Physics





# Intro to ML

### Serafina Di Gioia (postdoc @ ICTP) 05/05/2025

### **Books on ML**

**Bishop, Pattern recognition and ML** 

Elements of Statistical Learning, Springer, 2009, by Hastie et al.

Hands-On Machine Learning with Scikit-Learn and TensorFlow, by Aurelien Geron

**Deep Learning,** by Ian Goodfellow and Yoshua Bengio and Aaron Courville

(free online <a href="http://www.deeplearningbook.org/">http://www.deeplearningbook.org/</a>)

#### What is ML?



What kind of task are solved by AI systems today?



### and what do we use AI for in weather/climate..

- Forecasting
- Downscaling
- Anomaly detection
- Data assimilation
- Hybrid climate modeling
- Optimizing the parametrization of climate simulations

#### From LISP to the DL revolution ...

#### ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.

#### MACHINE LEARNING

Machine learning begins to flourish.



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

# History of ML





How to install it:

pip install seaborn[stats]

conda install seaborn -c conda-forge

Pandas

How to install it:

pip install pandas

conda install -c conda-forge pandas

Documentation at:

https://seaborn.pydata.org/api.html

Documentation at:

https://pandas.pydata.org/docs/getting\_st arted/index.html#getting-started

Documentation at:

https://scikit-learn.org/stable/user g uide.html

### scikit-learn

How to install it:

pip install scikit-learn

conda install -c scikit-learn-intelex

### ML vs AI vs DL (find the differences..)



#### from Google Al

# ML/DL vs classical explicit programming



# DL in a Venn diagram

and why this is not loved by all ML scientists...



# Building blocks of ML algorithms

ML algorithms have three main components

- 1. **decision process**: based on some input data, which can be labeled or unlabeled, your algorithm will produce an estimate about a pattern in the data. This estimate can be used to solve a regression or classification task
- 2. **error function**: it evaluates the prediction of the model. If there are known examples, an error function can make a comparison to assess the accuracy of the model.
- 3. **Model Optimization Process (training)**: If the model can fit better to the data points in the training set, then weights are adjusted to reduce the discrepancy between the known example and the model estimate. The algorithm will repeat this "evaluate and optimize" process, updating weights autonomously until a threshold of accuracy has been met.



### ML branches

	Supervised le	earning	Unsupervised learning	Reinforcen	nent learning
		$  y_{pred} $ $  y_{true} $ $  y_{true} $		Agent	Environment
	ML algorithm learns predicted and actual	by comparing values	ML algorithm learns without labeled data (e.g. clustering, embedding)	Agent (ML alg interacting wi	orithm) learns by th an environment
Learnii	ng type	Model building			Examples
Superv	ised	Algorithms or mod	dels learn from labeled data (task-driven approa	ach)	Classification, regression
Unsupe	ervised	Algorithms or mod	dels learn from unlabeled data (Data-Driven A	pproach)	Clustering, associa- tions, dimensionality reduction
Semi-s	upervised	Models are built u	sing combined data (labeled + unlabeled)		Classification, clustering
Reinfor	cement	Models are based of	on reward or penalty (environment-driven appr	roach)	Classification, control

# Applications of different methods

# Supervised learning



Unsupervised learning

#### Reinforcement Learning





What non-data scientists think of ML:



# What ML is in reality?



- 1. Problem definition
- 2. Data collection and exploration
- 3. Data preparation
- 4. Modeling (model selection & training)
- 5. Model evaluation
- 6. Model deployment and maintenance

### The importance of Exploratory data analysis...

results of surveys compiled by ML scientists some years ago...



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

but still valid!

Source: https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/



#### Preprocessing and Feature Engineering

Data Management for Digital Health, Winter 2019 6

### What are data?

for a scienceFactual information (such as measurements or statistics) usedcommunicatoras a basis for reasoning, discussion, or calculation

for a ML Information in digital form that can be transmitted or processed engineer

for a scientist

Information output by a sensing device or organ that includes both useful and irrelevant or redundant information

# Data type

#### Structured Data

#### Quantitative Data

Also called Numerical Data

Discrete Data

e.g. Population, number of attempts, etc.

#### Continuous Data

e.g. Height, Weight, Age etc.

Qualitative Data

Also called Categorical Data

Nominal Data

e.g. Male/Female, Black/Brown, etc.

#### Ordinal Data

e.g. First/Second, Good/Better/Best, etc.

#### Unstructured Data

Data that doesn't follow a pattern or sequence is called an unstructured data. e.g. Audio, Video, Image.

### Properties of structured vs Unstructured data



### Data containers in Pandas

#### **Series**

A Series is one-dimensional array-like object containing an array of data (of any NumPy data type) and an associated array of data-labels, called its index.

#### Dataframe

Tabular, spreadsheet-like data structure containing an ordered collection of columns of potentially different value types (numeric, string, etc.)

it can be regarded as a dict of Series

What a Dataframe looks like...



Visit

https://pandas.pydata.org/pandas-docs/stable/getting\_started/intro\_tutorials/01\_ta ble\_oriented.html for a more in-depth walkthrough

### High-level view of dataframes in Pandas

- head() first N observations
- tail() last N observations
- describe() statistics of the quantitative data
- dtypes the data types of the columns
- columns names of the columns
- shape the # of (rows, columns)

# **Exploratory Data Analysis (EDA)**



### Step 0-1: Distinguish attributes and visualizing the data



# Step 2-3 : Statistical analysis

to answer:

Univariate analysis:

- □ histograms/kDE
- summary statistics

Multivariate analysis:

- pair correlation
- multivariate joint distribution properties

#### Questions

- What is the central tendency of the data (mean, median, mode)?
- What is the spread of the data (range, interquartile range, standard deviation)?
- What is the shape of the distribution (symmetrical, skewed)?

# Step 4-5 : Data wrangling

It includes the following tasks:

- □ reformatting/cleaning data (outliers removal)
- data quality assessment
- □ data integration (imputation of missing values)
- □ data preprocessing

Libraries to implement data wrangling:

• Custom code (e.g. Pandas in Python, cudf, dplyr in R)

# EDA is an iterative process

First you explore the data graphically

- 1) Construct graphics and summary statistics to address questions
- 2) Inspect "answer" and assess new questions
- 3) Repeat...

keeping an eye on how to transform data appropriately (e.g., invert, log)

Finally, you perform DATA PREPROCESSING (also called feature engineering)

# Feature scaling for continuous variables

#### Why?

Feature scaling is essential for machine learning algorithms that calculate **distances between data**. If not scaled, the feature with a higher value range starts dominating when calculating distances

#### **Functions in Scikit-learn**

- 1. Min Max Scaler
- 2. Standard Scaler
- 3. Max Abs Scaler
- 4. Robust Scaler
- 5. Quantile Transformer Scaler
- 6. Power Transformer Scaler
- 7. Unit Vector Scaler

#### SMR 3935 - School on parallel programming and parallel architecture for HPC (Kathmandu, 2024)

### Visualizing the rescaling



# Visualizing different non-linear transformation in sklearn



# Feature preprocessing for categorical variables

The two most common techniques

#### Label encoding

- Label encoding returns different values for different classes, bringing in a natural ordering
- Recommended when it's reasonable to assume some sort of (equally spaced) ordering in your categorical feature

country	country	teaching_level	teaching_leve	
France	1	Beginner	0	
Italy	0	Medium Advanced	1	
Spain				

#### **One Hot encoding**

- One hot encoding returns as many dummy variables (0/1 columns) as the classes of the categorical feature
- *Recommended when* no ordering can be assumed in your categorical feature

country	country_France	country_Italy	country_Spain
France	1	0	0
Italy	0	1	0
Spain	0	0	1

special case: dummy encoding

### EDA with Pandas + seaborn

#### Data I/O

pd.read\_csv()

#### **Summary statistics**

df.info()

df.describe()

#### **Missing values**

.isnull(), .isna(), .notnull()

#### **Outlier detection**

#### **Plots**

- **Bar charts** compare different categories.
- Line charts show trends over time or across different categories.
- **Pie charts -** show proportions or percentages of different categories.
- **Histograms** show the distribution of a single variable.
- **Heatmaps** show the correlation between different variables.
- **Scatter plots -** show the relationship between two continuous variables.
- **Box plots -** show the distribution of a variable and identify outliers.
- Violin plots check normality assumption

### **Outlier detection**



IQR= Q3 - Q1

# Iris dataset

small classical ML dataset

first work: Fisher, 1936

UCI repo: https://archive.ics.uc i.edu/dataset/53/iris





### Univariate analysis on Iris data - Distplot

- What is the central tendency of the data (mean, median, mode)?
- What is the spread of the data (range, interquartile range, standard deviation)?
- What is the shape of the distribution (symmetrical, skewed)?


### **Boxplots**



### Violin plot

# Plotting the violinplot
sns.violinplot(x='Species', y='SepalLengthCm', data=data)
plt.show()



### How can we visualize uncertainty on summary statistics?

Using bootstrap...

In [115]: from pandas.plotting import bootstrap\_plot

In [116]: data = pd.Series(np.random.rand(1000))

In [117]: bootstrap\_plot(data, size=50, samples=500, color="grey");



### What is a Bootstrap sampling?



of size B = 5

Bootstrap samples are expected to **approximately** be **representative and independent samples** of the true data distribution (almost *i*. *i*. *d*.)

For n big enough, sampling from Z is like sampling from the true data distribution

If n >> B the samples should not be much correlated



We are thus assuming to fit and then average *m* independent models...

### **Pearson Correlation**



Pearson correlation coefficient

$$r = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$
Z-score

$$-1 < r(x, y) < 1$$

### Plots for Feature association analysis: heatmap

plt.figure(figsize=(7,5))
# Plotting the heatmap
sns.heatmap(data.corr(), annot=True)
plt.show()



### Multivariate analysis -> Pair plot (with Seaborn)

it shows joint and marginal distributions for all pairwise relationships and for each variable, respectively

sns.pairplot(data,hue="Species")
plt.show()



### **Traditional ML scenarios**



### Supervised learning framework

- Training :  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- $\boldsymbol{x}_i$  : input vector

$$oldsymbol{x}_i = \left[egin{array}{c} x_{i,1} \ x_{i,2} \ dots \ x_{i,n} \end{array}
ight], \quad x_{i,j} \in \mathbb{R}$$

- *y* : response variable
  - $y \in \{-1, 1\}$ : binary classification
  - $y \in \mathbb{R}$  : regression
  - what we want to be able to predict, having observed some new  $oldsymbol{x}$ .

### The supervised learning problem

**Problem:** Given a set of examples (training set)

$$\mathcal{T} = \left\{ (x^{(i)}, y^{(i)}), i = 1, \dots, n \right\} \qquad \qquad x^{(i)} \in \mathbb{R}^p$$

we wish to estimate a function

$$\hat{y} = f(x)$$

such that  $\hat{y}$  (predictor) is, in some sense, close to y.

### Empirical risk

• Theoretical **risk** of a candidate model *f* is the **expected loss** 

$$\mathcal{R}(f) := \mathbb{E}_{xy}[L(y, f(\mathbf{x}))] = \int L(y, f(\mathbf{x})) \, \mathrm{d}\mathbb{P}_{xy}$$

- Average error we incur when we use f on data from  $\mathbb{P}_{xy}$
- Goal in ML: Find a hypothesis  $f \in \mathcal{H}$  that **minimizes** this

To estimate the empirical risk we Just sum up all losses over training data

$$\mathcal{R}_{emp}(f) = \sum_{i=1}^{n} L\left(y^{(i)}, f\left(\mathbf{x}^{(i)}\right)\right)$$

### Learning as optimization problem

ERM optimization problem:

$$\hat{oldsymbol{ heta}} = rgmin_{oldsymbol{ heta}\in\Theta} \mathcal{R}_{\mathsf{emp}}(oldsymbol{ heta})$$

For (global) minimum  $\hat{\theta}$ :

$$orall eta \in \Theta: \quad \mathcal{R}_{\mathsf{emp}}(oldsymbol{\hat{ heta}}) \leq \mathcal{R}_{\mathsf{emp}}(oldsymbol{ heta})$$

Does not imply that  $\hat{\theta}$  is unique.

- Best numerical optimizer depends on problem structure
- Continuous params? Uni-modal  $\mathcal{R}_{emp}(\theta)$ ?
- Numerical optimization not our focus here, now

### **Stochastic Gradient Descent**

The simple idea of GD is to iteratively go from the current candidate  $\theta^{[t]}$  in the direction of the negative gradient, i.e., the direction of the steepest descent, with learning rate  $\alpha$  to the next  $\theta^{[t+1]}$ :

$$\boldsymbol{\theta}^{[t+1]} = \boldsymbol{\theta}^{[t]} - \alpha \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{R}_{emp}(\boldsymbol{\theta}^{[t]}).$$



We choose a random start  $\theta^{[0]}$  with risk  $\mathcal{R}_{emp}(\theta^{[0]}) = 76.25.$ 

### What is the right model complexity?



### **Bias and variance**

**Bias:** a systematic shift of the entire distribution of points

**Variance:** the spread in data from with respect to the mean position.

Low Variance **High Variance** Low Bias High Bias

# How do we prevent overfitting without reducing the training sample size significantly?

#### TRAINING/DEVELOPMENT/TEST SET

- Highly efficient and easy to implement from a computational point of view
- The dataset has to be inherently homogeneous, otherwise with a single 'unlucky' split you may end up with a not representative training set

#### **CROSS-VALIDATION**

- Efficient sample re-use (all data are used for training, hence all possible sub-populations are covered)
- Adaptable to different sample sizes
- Can become computationally expensive when too many folds need to be evaluated

### Training/Validation/Test split (to build robust models)

To deploy a ML algorithm parametrized by a set of hyper-parameters we do:

- 1. Select HP values
- 2. Split input data set into training/validation/test set
- 3. Train on the training set
- 4. Evaluate performances on the validation set
- 5. Go back to 1 (select new HP values)
- 6. Select HP leading to the smallest validation error
- 7. Compute error estimation on the test set



### **Cross-validation**

• Exhaustive cross-validation methods learn and test on all possible ways to divide the original sample into a training and a validation set

#### LEAVE-p-OUT CROSS-VALIDATION

- Given a dataset containing n events
- Use p events out of the n for the validation and (n p) for the training
- Repeat the training and test for all possible combinations of the p events
- Average the results of each train-test to obtain the overall results



The number of possible combination is determined by the Binomial coefficient  $C_p^n = {n \choose p}$ . Therefore:  $C_{30}^{100} = {100 \choose 30} \cong 3 * 10^{25}$ 



For *p* > 1 and for even moderately large *n*, it may quickly become computationally infeasible

Non-exhaustive cross-validation do not compute all ways of splitting the original sample. Those methods are approximations of leave-p-out cross-validation.





### Other causes of overfitting: data leakage

**Preprocessing Leakage**: If data preprocessing (like normalisation, feature selection) uses the entire dataset rather than just the training set, information from the test set can leak into the model.

**Model Selection Leakage**: When model selection or hyperparameter tuning is done using the test set, it leads to an overfitting on the test data, providing a misleadingly high performance estimate (more on that later).

**Temporal Data Leakage**: In time-series data, using future data in the training process leads to leakage as the model gets access to information it would not have in a real-world scenario.

### **Regression models**

- KNN regression
- Linear regression
- Regularized/penalized linear regression (Lasso, Ridge)
- Bayesian linear regression
- Gaussian process regression (non-parametric)

The first two regression can be implemented with OLS method or with SGD.

### **KNN** regression

Given a dataset  $D = \{(x^{(1)}, y^{(1)}), ..., (x^{(N)}, y^{(N)})\}$  for every new x do:

- find the k-number obs in D most similar to X (these k obs. points are called k-nearest neighbours of x)
- 2. average the output of the k-nearest neighbors of x

$$\hat{y} = \frac{1}{K} \sum_{k=1}^{K} y^{(n_k)}$$

### Pros -Cons of KNN

As with K-NN classification (or any prediction algorithm for that matter), K-NN regression has both strengths and weaknesses. Some are listed here:

**Strengths:** K-nearest neighbors regression

- 1. is a simple, intuitive algorithm,
- 2. requires few assumptions about what the data must look like, and
- 3. works well with non-linear relationships (i.e., if the relationship is not a straight line).

Weaknesses: K-nearest neighbors regression

- 1. becomes very slow as the training data gets larger,
- 2. may not perform well with a large number of predictors, and
- 3. may not predict well beyond the range of values input in your training data.

### Linear regression



Linear least squares fitting with  $X \in \mathbb{R}^2$ . We seek the linear function of X that minimizes the sum of squared residuals from Y.

From: Hastie, Tibshirani, Friedman, "The Elements of Statistical Learning", Springer,

### **Linear regression**

 $\bigcirc$ 

Model definition: 
$$y = \beta_0 + \beta_1 x + \varepsilon$$
, where:

- $\beta_0$  is the intercept of the model,
- $\beta_1$  represents the slope

 The best model is the one that minimises the sum of squared (SS) residuals over the n observations:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^{n} [y_i - f(\beta, x_i)]^2 \quad \boxed{\begin{array}{c} \text{COST} \\ \text{FUNCTION} \end{array}}$$

• The algorithm used to find the optimal values  $\hat{\beta}_0$  and  $\hat{\beta}_1$  is called **OLS** (Ordinary Least Squares)



### **Ordinary Least Squares**

$$\hat{y} = w^T \mathbf{x} + b = \sum_{i=1}^p w_i x_i + b$$

$$\min_{w \in \mathbb{R}^p, b \in \mathbb{R}} \sum_{i=1}^n ||w^T \mathbf{x}_i + b - y_i||^2$$

Unique solution if  $\mathbf{X} = (\mathbf{x}_1, \dots \mathbf{x}_n)^T$  has full column rank.

### Linear regression from a probabilistic perspective

If we assume

 $Y \in \mathbb{R}$  and  $\mathbf{X} \in \mathbb{R}^p$ 

$$f_{Y|\mathbf{X}}(y|\mathbf{x}) = \mathcal{N}(y|\mathbf{w}^T\mathbf{x} + w_0, \sigma^2),$$

Parameters of  $f_{Y|\mathbf{X}}$  are unknown; instead, i.i.d. training data:

 $\mathcal{D} = ((\mathbf{x}_1, y_1), ..., (\mathbf{x}_n, y_n))$ 

the likelihood is :

$$\mathbb{L} = \prod_{i=1}^{T} \mathcal{N}(y_i; \mathbf{w}^T \mathbf{x}_i + w_0, \sigma^2)$$

$$\log(\mathbf{L}) = K - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \mathbf{w}^T \mathbf{x}_i - w_0)^2$$

Maximum likelihood estimate of w:

$$(\hat{\mathbf{w}}, \hat{w_0})_{\mathsf{ML}} = \arg\min_{\mathbf{w}, w_0} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - w_0)^2$$

Another view: loss function  $L(y, \hat{y}) = (y - \hat{y})^2$ ; Bayes/expected risk for  $\hat{y} : \mathbb{R}^p \to \mathbb{R}, \ \hat{y}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ :

$$R[\mathbf{w}, w_0] = \mathbb{E}[(Y - w^T X - w_0)^2] = \int \int (y - \mathbf{w}^T \mathbf{x} - w_0)^2 \underbrace{f_{Y, \mathbf{X}}(y, \mathbf{x})}_{\text{unknown}} d\mathbf{x} \, dy$$

The empirical risk is, in this case, the residual sum of squares (RSS)

$$R_{\mathsf{emp}}[\mathbf{w}, w_0] = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - w_0)^2 = \mathsf{RSS}(\mathbf{w}, w_0)$$

### **Ridge Regression**

$$\min_{w \in \mathbb{R}^p, b \in \mathbb{R}} \sum_{i=1}^n (w^T \mathbf{x}_i + b - y_i)^2 + lpha ||w||^2$$

Always has a unique solution. Tuning parameter alpha. L2 regularization

### **Lasso Regression**

$$\min_{w \in \mathbb{R}^p, b \in \mathbb{R}} \sum_{i=1}^n ||w^T \mathbf{x}_i + b - y_i||^2 + \alpha ||w||_1$$
  
L1 regularization

- Shrinks w towards zero like Ridge
- Sets some w exactly to zero automatic feature selection!

- Ridge results in smooth solution path with non-sparse params
- Lasso induces sparsity, but only for large enough  $\lambda$



### Bayesian Learning for supervised ML

**Step 1:** Given *n* data,  $D = x_{1:n} = \{x_1, x_2, ..., x_n\}$ , write down the expression for the likelihood:

 $p(\mathbf{D}|\boldsymbol{\theta})$ 

**Step 2:** Specify a **prior**:  $p(\theta)$ 

**Step 3:** Compute the **posterior**:

$$\frac{p(\boldsymbol{\theta} \mid \boldsymbol{D})}{p(\boldsymbol{D})} = \frac{p(\boldsymbol{D} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\boldsymbol{D})}$$

Bayesian Learning for supervised ML



### **Bayesian Linear regression**

The likelihood is a Gaussian,  $\mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_n)$ . The conjugate prior is also a Gaussian, which we will denote by  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_0, \mathbf{V}_0)$ .

Using Bayes rule for Gaussians, the posterior is given by

$$p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y}, \sigma^2) \propto \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_0, \mathbf{V}_0) \mathcal{N}(\mathbf{y}|\mathbf{X}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_n) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\theta}_n, \mathbf{V}_n)$$
$$\boldsymbol{\theta}_n = \mathbf{V}_n \mathbf{V}_0^{-1} \boldsymbol{\theta}_0 + \frac{1}{\sigma^2} \mathbf{V}_n \mathbf{X}^T \mathbf{y}$$
$$\mathbf{V}_n^{-1} = \mathbf{V}_0^{-1} + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X}$$



### Special case:

Consider the special case where  $\boldsymbol{\theta}_0 = \mathbf{0}$  and  $\mathbf{V}_0 = \tau_0^2 \mathbf{I}_d$ , the posterior mean reduces to

$$\boldsymbol{\theta}_n = \frac{1}{\sigma^2} \mathbf{V}_N \mathbf{X}^T \mathbf{y} = \frac{1}{\sigma^2} \left( \frac{1}{\tau_0^2} \mathbf{I}_d + \frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

$$= \left( \lambda \mathbf{I}_d + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

-

where we have defined  $\lambda := \frac{\sigma^2}{\tau_0^2}$ . We have therefore recovered **ridge re-gression** again!

## Extending regression to non-normal /non-linear problems

GLM

$$g(\mathbb{E}[Y])=eta_0+eta_1X_1+eta_2X_2$$

#### GAM

Further relax GLMs by modeling each predictor with a **smooth function**:

$$g(\mathbb{E}[Y]) = eta_0 + f_1(X_1) + f_2(X_2)$$

Common smooth functions: splines, loess...

Advantages: captures non-linear effects, retains interpretability
# **Gaussian Process regression**

There are two possible definitions:

- 1. A Gaussian process is a collection of random variables, any finite number of which have consistent Gaussian distributions.
- 2. A Gaussian process is a probability distribution over possible functions that fit a set of points.

$$f(\mathbf{x}) \sim GP(\underline{m}(\mathbf{x}), \underline{\kappa}(\mathbf{x}, \mathbf{x}'))$$
  

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$$
  

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))^T] \quad \left[k(x, x') = \exp(-\frac{1}{2}(x - x')^2)\right]$$

## Important hyper-pars of GP regressor

Non-parametric ( $\infty$ -parametric)

Parametric model



# GP algorithm

**input**: X (inputs), y (targets), k (covariance function),  $\sigma_n^2$  (noise level),  $\mathbf{x}_*$  (test input) 2:  $L := \text{cholesky}(K + \sigma_n^2 I)$  $\boldsymbol{\alpha} := L^{\top} \backslash (L \backslash \mathbf{y})$  $\left.\right\}$  predictive mean eq. (2.25) 4:  $\bar{f}_* := \mathbf{k}_*^\top \boldsymbol{\alpha}$  $\mathbf{v} := L \setminus \mathbf{k}_*$  $\left. \right\}$  predictive variance eq. (2.26) 6:  $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^\top \mathbf{v}$  $\log p(\mathbf{y}|X) := -\frac{1}{2}\mathbf{y}^{\top}\boldsymbol{\alpha} - \sum_{i} \log L_{ii} - \frac{n}{2} \log 2\pi$ eq. (2.30)8: return:  $\bar{f}_*$  (mean),  $\mathbb{V}[f_*]$  (variance),  $\log p(\mathbf{y}|X)$  (log marginal likelihood)

## Estimating model performance



## Pros/cons of MSE vs RMSE vs MAE

• 
$$MSE = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2$$

• 
$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)^2}$$

• 
$$MAE = \frac{1}{n} \sum_{i=1}^{n} |\hat{y}_i - y_i|$$

It assigns more weight to bigger errors, useful when interested in **detecting outliers** 

It maintains the same properties as the MSE, but it is in the same scale as y

It is simply the average difference observed in the predicted and actual values across the whole test set, useful if not particularly interested in outliers

# **Classification models**

- They all share the same schema, by having:
  - a response variable (target) of categorical nature,
  - a number of explanatory variables (features) which are assumed to be helpful in classifying each observation

target

$$y = \begin{cases} 0, & negative \ class \\ 1, & positive \ class \end{cases}$$



# Approaches to classification tasks

Two fundamental approaches exist to construct a classifier:

- **Discriminant approach** asks "What is the best prediction for the class given these data?" (uses loss functions and empirical risk minimization)
- Generative approach asks "Which class tends to have data like these?" (models the feature distributions in each class separately)



observations are assigned to the same class.

**Decision Boundary:** Points where all classes have the same probability/score.

# Logistic regression for binary classification

The logistic regression aims at modelling a function with this shape (also called **logistic or sigmoid function**):



- It is bounded in [0, 1] as probabilities also are
- It is *S*-shaped, being is a cumulative distribution along the variable  $\int_{x0}^{x1} pdf(x) dx$

$$p(x;\beta) = \frac{e^{\beta_0 + \vec{\beta} \cdot \vec{x}}}{1 + e^{\beta_0 + \vec{\beta} \cdot \vec{x}}}$$

# The logit function

$$z = \left(\sum_{i=1}^n w_i x_i\right) + b$$



# Support vector classifier

**Basic idea:** maximize the distance between the support vector points and the decision boundary between classes



# **Decision trees**

**TASK**: fit a 2-depth tree on the following data 1.5 1.0 0.5 Q 0.0 -0.5 -1.0-1.5 -1.5 0.5 10 15 -1.0-0.5 0.0 x1

 The algorithm aims at partitioning the two-dimensional features space by recursive binary splitting (sub-optimal solution)



- At each node, all the available features are evaluated
- A given objective function is then optimized to find the most informative feature and point for the split

# How do decision trees actually works?

## Splitting algorithm

- 1. Calculate the I(D) of the parent node
- 2. Calculate the weighted average I(D) of child nodes
- 3. Subtract the two quantities (estimate IG)
- 4. Select the split with the highest information gain

$$IG(D_p) = I(D_p) - \sum_{j=1}^{m} \frac{n_j}{n_p} I(D_j)$$
  
Gini index: 
$$I(D_j) = \sum_{k=0,1} p_{jk} (1 - p_{jk})$$
  
Entropy: 
$$I(D_j) = -\sum_{k=0,1} p_{jk} \log_2(p_{jk})$$



# What is the best depth for our task?

#### **Decision boundaries** 2-DEPTH TREE Value represents the 2.0 $x1 \le 0.545$ У number of samples gini = 0.5• 0 15 samples = 1000in the class 0 and 1 • 1 value = [500, 500] respectively class = 010 0.5 $x1 \le 0.757$ $x1 \le -0.552$ Q 0.0 qini = 0.485qini = 0.094samples = 839samples = 161value = [347, 492] value = [153, 8]-0.5class = 1class = 0-1.0-1.5qini = 0.159qini = 0.416qini = 0.256qini = 0.0samples = 161samples = 678samples = 53samples = 108-2.0 value = [147, 14] value = [200, 478] value = [108, 0]-2.0 -1.5 -1.0 -0.5 0.0 0.5 10 15 20 value = [45, 8]x1 class = 1class = 0class = 0class = 0

# The power of collective knowledge

Ensemble models can be used both in regression and classification tasks

- REGRESSION: the final predictions can be the average or weighted average of the single models' predictions
- CLASSIFICATION: as before, the average or weighted average may be taken to compute the scores.
   The final classification is based on the majority vote among the single classifiers



#### Wisdom of Crowds (Surowiecki, 2004)

The collective knowledge of a diverse and independent body of people typically exceeds the knowledge of any single individual, and can be harnessed by voting





- We call base or weak classifier the model used as building block of the ensamble model
- The base classifier is trained on **subsamples** of the training data, so that:
  - it should learn the details of the specific subsample
  - b does not need to be much granular (e.g. a very deep decision tree) to learn details because the subsample has a relatively small size, hence less information to draw from (→ weak)
- As it is trained multiple times, the **computational cost** of training the model is a relevant aspect to consider
- A very widespread choice is the decision tree
  - It is a highly unstable model (high variance: any little perturbation of the training data may result in significantly different tree structures)
  - Notoriously very fast to run

# **Ensemble models**



# Boosting

What differs is the **strategy** used to generate the subsets:

- bagging uses bootstrap sampling,
- boosting, at each step, increases the probability of picking events that were misclassified at the previous step

In the boosting model, each *weak classifier* is mostly focused on **learning how to correct the errors** done by the previous *weak classifier* 



# Bagging



DATASET



SAMPLING

## weak learner





Let  $Z = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ be the complete training set Create m subsamples  $Z_i$  of size B out of Zusing **bootstrapping** techniques (random sampling with replacement) Let  $C_i(x)$  be the classifier trained on the subsample  $Z_i$ (base classifier)

Finally define the ensemble classifier as the average prediction over the *m* base learners:

$$C_{bag}(x) = \frac{1}{m} \sum_{i=1}^{m} C_i(x)$$

# **Random Forest**

 Random forest is a substantial modification of *bagging* that builds a large collection of *de-correlated trees* and then averages them (it is therefore specifically developed for trees as base learners)

> The idea is to **improve the variance reduction** of *bagging* **by reducing the correlation of the trees**

> > To do so, in the tree-growing process features are randomly selected

• Before each split,  $m \le p$  of the input features are selected at random as candidates for splitting (a common choice is  $m = \sqrt{p}$  for classification tasks and  $m = \lfloor p/3 \rfloor$  for regression tasks)

# What kind of problems we solve with ensemble models?

Averaging weak learners outputs do not change the expected value (bias of the model) but reduces its variance

(just like averaging i. i. d. random variables preserve expected value but reduces variance)



A variance reduction **decreases the prediction error** of the final model (bias-variance tradeoff)

Err(x)=Bias<sup>2</sup>+Variance+Irreducible Error

On which models bagging will perform especially well?

On low bias, high variance models (e.g. trees)

# **Confusion matrix**

		PREDICTION		
		Positive	Negative	
TRUTH	Positive	True Positive <b>TP</b>	False Negative <b>FN</b>	Type I error
	Negative	False Positive <b>FP</b>	True Negative TN	

## Type I error

# **Evaluation metrics for classifiers**

## F1 score

$$2\frac{precision \cdot recall}{precision + recall} = \frac{2 \cdot TP}{2 \cdot TP + FP + FN}$$

$$extbf{accuracy}(y, \hat{y}) = rac{1}{n_{ ext{samples}}} \sum_{i=0}^{n_{ ext{samples}}-1} 1(\hat{y}_i = y_i)$$

#### FPR (False Positive Rate)

- What is the proportion of misclassified instances out of all the true negatives?
- $\frac{FP}{FP+TN}$  (should be as low as possible)

#### **RECALL or TPR (True Positive Rate)**

- What proportion of actual positives is identified correctly?
- $\frac{TP}{TP+FN}$  (should be high as possible)

#### PRECISION

- What proportion of positive predictions is actually correct?
- $\frac{TP}{TP+FP}$  (should be high as possible)

# How these metrics are called in scikit-learn?

Scoring string name	Function	Comment
Classification		
'accuracy'	<pre>metrics.accuracy_score</pre>	
'balanced_accuracy'	<pre>metrics.balanced_accuracy_score</pre>	
'top_k_accuracy'	<pre>metrics.top_k_accuracy_score</pre>	
'average_precision'	<pre>metrics.average_precision_score</pre>	
'neg_brier_score'	<pre>metrics.brier_score_loss</pre>	
'f1'	<pre>metrics.f1_score</pre>	for binary targets

## **Unsupervised Learning algorithm**

- Manifold Learning (Dim. Red. & ID est.):
  - PCA
  - K-PCA
  - ISOMAP
  - t-SNE
  - Autoencoders
- Density Estimation:
  - Histograms
  - Kernel Density Estimation
  - k-Nearest Neighbor
  - Generative Adversarial NN
- Clustering:
  - k-means/c-means, kernel k-means, spectral clustering...
  - Hierarchical clustering
  - Density Based clustering
  - Self Organizig Maps



# PCA steps

Given m observations of n-dim column vectors , to perform PCA we follow the 6 steps below:

1. Compute the mean vector

$$\vec{\mu} = \frac{1}{m} \sum_{j=1}^{m} \vec{x}^{(j)}$$

- 2. Compute the covariance matrix ( $C = X^T X$ ,  $S = C C_mean$ )
- 3. Compute the sorted eigenvalue/eigenvector pairs (Diagonalization/SVD)
- 4. Choose the number of dimensions (k) in which project the data
- 5. Project the input data onto the selected k eigenvectors

$$\mathbf{X} pprox \sum_{k=1}^r \mathbf{w}_k \mathbf{c}_k^T \quad ext{or} \quad \mathbf{X} pprox WC^T$$



**Example reconstruction of data with 3 principal components.** A data matrix (*left*) is approximated by the product of a  $n \times r$  matrix and a  $r \times p$  matrix (i.e.  $WC^T$ ). This product is at most a rank-r matrix (in this example, r = 3). Each paired column of W and row of  $C^T$  form an outer product, so the full reconstruction can also be thought of as a sum of r rank-one matrices.



How do I choose the number of PCA components?

The scree plot (Cattel method)



# CLUSTERING

Clustering tries to separate data "naturally", in such a way that *similar* elements lay in the same cluster while *dissimilar* elements belong to a different one

## **Similarity**

pairwise function of the features. Its definition depends on the nature of input data. It can be seen as a distance (but it does not always corresponds to metric distances!)

Clustering ≠ classification it generates groups without labels

# Considerations on clustering

The result of the clustering process depends on:

- your cluster definition
- the metric adopted to measure distances
- the feature you choose to consider to cluster data

# What is a cluster?



# Other cluster examples...



# Clustering stone collections

Cluster by light wavelength



Are they nice or ugly? Color classification?

Most common flat clustering algorithm:

# K-means: A flat clustering algorithm

MacQueen, J. (1967, June). Some methods for classification and analysis of multivariate observations. In Proceedings of the fifth Berkeley symposium on mathematical statistics and probability (Vol. 1, No. 14, pp. 281-297).

# K-means algorithm

**Input:** cluster size k, instances  $\{\vec{x_i}\}_{i=1}^n$ **Output:** cluster membership assignments  $\{z_i\}_{i=1}^n$ 

- 1. Initialize k cluster centroids  $\{\vec{c}_l\}_{l=1}^k$  (randomly from the data set).
- 2. Repeat until no instance changes its cluster membership:
  - Decide the cluster membership of instances by assigning them to the nearest cluster centroid

 $z_i = argmin_l(d(\vec{c_l}, \vec{x_i}))$  Minimize intra distance

Update the k cluster centroids based on the assigned cluster membership

$$\overrightarrow{c_l} = \frac{\sum_{i=1}^n \delta_{z_{il}} \overrightarrow{x_i}}{\sum_{i=1}^n \delta_{z_{il}}}$$

Maximize inter distance

## Loss function/Objective function

• Loss function:

$$O(z) = \sum_{l=1}^{k} \sum_{i=1}^{n} \delta(z_{i,l}) \quad \|\overrightarrow{x_{i}} - \overrightarrow{c_{l}}\|^{2}$$

- z is an array with n components that reflects the assignation in clusters.
- $\overrightarrow{c_l}$  is the vector with the coordinates of the *l*-th cluster centroid.  $\overrightarrow{c_l} = \frac{\sum_{i=1}^n \delta_{z_{il}} \overrightarrow{x_i}}{\sum_{i=1}^n \delta_{z_{il}}}$
## K-means steps

- Randomly pick k centers.
- Assign each point to its nearest center.
- Recompute centers.
- Iterate...



#### SMR 4067- AI and Climate Modeling



# **Clustering proteins**

Hierarchical clustering of the secretome reveals clusters of secreted protein families as high priority effector candidates.

Saunders DGO, Win J, Cano LM, Szabo LJ, Kamoun S, et al. (2012) Using Hierarchical Clustering of Secreted Protein Families to Classify and Rank Candidate Effectors of Rust Fungi. PLoS ONE 7(1): e29847. doi:10.1371/journal.pone.0029847 http://journals.plos.org/plosone/article?id=info:doi/10.137 1/journal.pone.0029847

By clustering and classifying plant proteins detect protein groups (families) that will probably act against mold

# Outliers and novelty detections (2 faces of anomaly detection)

Many applications require being able to decide whether a new observation belongs to the same distribution as existing observations (it is an *inlier*), or should be considered as different (it is an *outlier*). Often, this ability is used to clean real data sets. Two important distinctions must be made:

outlier detection:	The training data contains outliers which are defined as observations that are far from the others. Outlier detection estimators thus try to fit the regions where the training data is the most concentrated, ignoring the deviant observations.
novelty detection:	The training data is not polluted by outliers and we are interested in detecting whether a <b>new</b> observation is an outlier. In this context an outlier is also called a novelty.

### other methods for outliers detection in scikit-learn



SMR 4067- AI and Climate Modeling