








# INTRODUCTION TO MACHINE LEARNING

ML-INFN HACKATHON 2023

---

Prof. Luca Scrucca

-  Università degli Studi di Perugia
-  [luca.scrucca@unipg.it](mailto:luca.scrucca@unipg.it)
-  <http://www.stat.unipg.it/luca>
-  [luca-scr](#)
-  [luca\\_scr](#)

21 June 2023

All slides © Luca Scrucca

All slides can be used for educational purposes if this copyright notice is included  
Permissions must be obtained from the copyright holder(s) for any other use



- **Leo Breiman** (“Statistical Modeling: The Two Cultures”, Statistical Science, 2001) described “two cultures”:
  1. “generative” modeling culture which seeks to develop stochastic models that fit the data, and then make inferences about the data-generating mechanism based on the structure of those models. Implicit is the notion that there is a true model generating the data, and often a “best” way to analyze the data.
  2. “predictive” modeling culture which focuses on predictions, ignoring the underlying data generating mechanism, and discuss only accuracy of predictions made by different algorithms.
- According to Breiman “Statistics starts with data. Think of the data as being generated by a black box [...]”

Two main goals can be pursued when analyzing data:

- Prediction, i.e to be able to predict what the responses are going to be to future input variables;
- Inference, i.e to infer how nature is associating the response variables to the input variables.

## MACHINE LEARNING

---

- **Machine learning** is the study of how computer algorithms can improve automatically through experience and by the use of data.



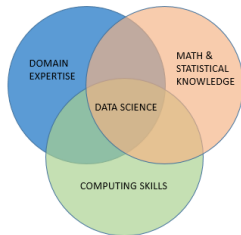
*A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$  if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ . – Tom Mitchell (1997)*

- In the context of predictive modelling, the difference between machine learning and statistical learning is blurred.
  - **Machine learning** (ML) tends to be focused more on developing efficient algorithms that scale to large data in order to optimize a predictive model.
  - **Statistical learning** (SL) generally pays more attention to the probabilistic structure of the model in order to provide an assessment of the uncertainty.

- **Data Science** is a vaguely defined, constantly changing, cross-disciplinary field.  
From a statistician point of view, data science can be seen as a broader view of statistics.

*When physicists do mathematics, they don't say they're doing "number science". They're doing math. If you're analyzing data, you're doing statistics. You can call it data science or informatics or analytics or whatever, but it's still statistics.*  
— Karl Broman (U of Wisconsin)

- **Big Data** refers to data sets that are too large or complex to be dealt with by traditional data analysis software.  
Big data are usually described in terms of three key concepts: [volume](#), [variety](#), and [velocity](#).





- Suppose we collected data for a sample of  $n$  observations. The **training set** is made of pairs of input and output variables:

$$\mathcal{D}_{\text{train}} = \{\mathbf{x}_i, y_i\}_{i=1}^n$$

- Assume there exists a dependency between them, so the output  $y_i$  can be expressed as a function of the input variables  $\mathbf{x}_i$  and some other unobservable (latent) variables  $\mathbf{z}_i$ :

$$y_i = f(\mathbf{x}_i, \mathbf{z}_i)$$

- The aim of **supervised learning** is to fit a model to learn the mapping from the observable input to the output

$$\hat{y}_i = g(\mathbf{x}_i | \boldsymbol{\theta})$$

where  $g(\cdot)$  is a statistical model and  $\boldsymbol{\theta}$  the unknown parameters.

- The learning task corresponds to finding the parameters that minimize a **loss function** measuring the deviation of our prediction  $\hat{y}_i$  from the observed output  $y_i$ :

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n L(y_i, \hat{y}_i) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n L(y_i, g(\mathbf{x}_i | \boldsymbol{\theta}))$$

- Different supervised learning algorithms differ in the models or the loss functions they assume, or the procedures they use in optimization.
- In **regression** problems
  - the output  $y_i$  is a numerical value (quantitative response)
  - $g(\cdot)$  is a regressor function
  - loss is often the squared error, so the aim is to find the best  $\theta$  that minimize the fitting error.
- In **classification** problems
  - the output  $y_i$  is a discrete label (qualitative response)
  - $g(\cdot)$  is a discriminant/classification function
  - if the loss function is the Zero-One loss the aim is to minimize the total number of misclassifications.



- Popular supervised learning models are:
  - Linear regression
  - Logistic regression
  - Generalized Linear Models (GLM)
  - Generalized Additive Models (GAM)
  - Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA)
  - Naive Bayes methods
  - Mixture models (e.g. Gaussian mixtures)
  - Decision Trees (Regression and Classification Trees)
  - Ensemble methods (Bagging, Random Forests, Boosting)
  - Support Vector Machines (SVM)
  - Neural Networks (NN) and Deep Learning (DL)

- Suppose we collected a dataset  $\mathcal{D}_{\text{train}} = \{\mathbf{x}_i\}_{i=1}^n$  composed of only a set of variables drawn from some unknown probability/density function

$$\mathbf{x}_i \sim p(\mathbf{x})$$

- In **unsupervised learning** for each case only the predictors vector  $\mathbf{x}_i$  is observed, but there is no response  $y_i$  ( $i = 1, \dots, n$ ). Thus, we lack a response variable that can supervise our analysis.

- The aim is to estimate a model with parameters  $\theta$

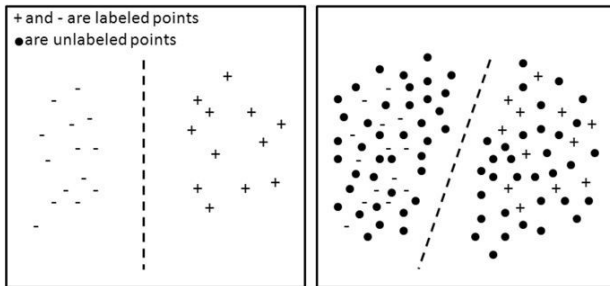
$$\mathbf{x}_i \approx q(\mathbf{x} \mid \theta)$$

where  $q(\cdot)$  is some working distribution depending on parameters  $\theta$ .

- The learning task corresponds to finding the parameters that makes  $q(\cdot)$  as close as possible to the unknown  $p(\cdot)$  and from that understand the relationships between the variables or between the observations.
- Cluster analysis is a typical unsupervised learning task: look for the presence of one or more distinct groups of observations with no explicit assessment criterion because truth is not known (e.g. market segmentation to detect groups of customers).

## SEMI-SUPERVISED LEARNING

- Many problems fall naturally into the supervised or unsupervised learning paradigms.  
However... sometimes the question is less clear-cut.
- There can be situations where for a subset of  $m < n$  observations we have information on both the predictors and the response variable, and for the remaining  $n - m$  observations we have only predictor measurements but no response measurement.



- Such a scenario is referred to as **semi-supervised learning**.

- Reinforcement learning focuses on training agents to make sequential decisions in an environment to maximize a cumulative reward.
- In reinforcement learning, an agent interacts with an environment, learns from its actions, and adjusts its behavior to achieve a specific goal or objective.



- Application domains:
  - recommendation systems (e-commerce websites, streaming services, ...)
  - robotics
  - game playing
  - autonomous vehicles

- In general, suppose that we have observed
  - a quantitative [response](#) (aka dependent variable, output, target, ...) denoted as  $Y$ , and
  - a set of  $p$  [predictors](#) (aka independent variables, covariates, features, ...) collected in the input vector  $X = (X_1, X_2, \dots, X_p)^T$ .

- Further, assume there exists some relationship between  $Y$  and  $X$ , i.e.

$$Y = f(X) + \epsilon$$

where

- $f()$  is unknown and represents the systematic information that  $X$  provides about  $Y$ ;
  - $\epsilon$  is a random error term, which captures measurement errors and other discrepancies, independent of  $X$  and with zero mean.
- There are two main reasons to estimate  $f()$ :
    1. [inference](#)
    2. [prediction](#)

## INFERENCE

---

- In descriptive or explanatory modelling we want to understand how  $Y$  changes as a function of  $(X_1, \dots, X_p)$ .
- Interesting questions:
  - Which predictors are associated with the response?
  - What is the relationship between the response and each predictor?
  - What is the functional form of the relationship between  $Y$  and each predictor?
  - There exists any cause-and-effect relationship?  
(causal inference)

## PREDICTION

---

- In predictive modelling the goal is to predict the response variable based on the observed values of the predictors:

$$\widehat{Y} = \widehat{f}(X)$$

- $\widehat{f}()$  is often treated as a **black box**: we are not interested in knowing the exact form of  $f()$ , provided that it yields accurate predictions for  $Y$ .
- The prediction error of estimating  $Y$  using  $\widehat{Y}$  can be decomposed as

$$Y - \widehat{Y} = f(X) + \epsilon - \widehat{f}(X) = (f(X) - \widehat{f}(X)) + (Y - f(X))$$

- Suppose that both  $f()$  and  $X$  are fixed, then recalling that  $E[\epsilon] = 0$ , the expected prediction error (under squared error loss) is given by

$$\begin{aligned} E[(Y - \widehat{Y})^2] &= E[(f(X) + \epsilon - \widehat{f}(X))^2] \\ &= E[(f(X) - \widehat{f}(X))^2] + E[\epsilon^2] + 2E[\epsilon(f(X) - \widehat{f}(X))] \\ &= \underbrace{(f(X) - \widehat{f}(X))^2}_{\text{reducible error}} + \underbrace{V[\epsilon]}_{\text{irreducible error}} \end{aligned}$$

## ESTIMATING $f()$

---

- Estimation (or learning in ML) is the process of applying a [statistical/machine learning method to the training data to estimate the unknown function  \$f\(\)\$](#) .
- Main goal: estimate  $f()$  with the aim of [minimizing the reducible error](#).
- The [irreducible error](#) provides a lower bound on the accuracy of our prediction for  $Y$ , and it is almost always unknown in practice.
- Several approaches are available, both parametric and non-parametric.



## PARAMETRIC METHODS

---

A two-step model-based approach:

1. Select the functional form, or shape, of  $f()$ .

For example, the **linear model** assumes that  $f()$  is linear in  $X$ :

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

2. Select a procedure that uses the training data to fit or train the model.

For example, in the linear model case we only need to **estimate the parameters** ( $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ ). A popular approach is (ordinary) least squares (OLS), but many other exists (maximum likelihood, regularized ML, Bayesian estimation, ...).

- This model-based approach is called parametric because it reduces the problem of estimating  $f()$  down to one of estimating a set of parameters (the coefficients of the model).

Pros: generally is much easier to estimate a set of parameters than it is to fit an entirely arbitrary function  $f()$ .

Cons: the selected model can be a poor approximation of true unknown form of  $f()$ .

## NON-PARAMETRIC METHODS

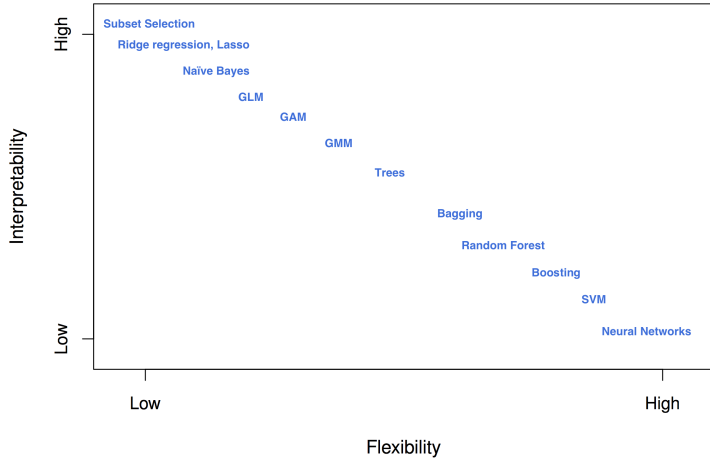
---

- Non-parametric methods do not make explicit assumptions about the functional form of  $f()$ .
- They try to estimate  $f()$  getting as close to the data points as possible without being too rough or wiggly.

Pros: avoid the assumption of a particular functional form for  $f()$ , so they have the potential to accurately fit a wider range of possible shapes for  $f()$ .

Cons: a large number of observations is required to accurately estimate  $f()$ .

- If we are mainly interested in explanatory inference, then simple models (e.g. Linear Models, Logistic Regression) are much more interpretable than black-box models (e.g. Random Forest, SVM, Neural Networks).
- Flexible models allow to fit many different possible functional forms for  $f()$ , but usually require estimating a larger number of parameters.
- In general, as the flexibility of a model/algorithm increases, its interpretability decreases.
- **Overfitting** is the main risk, i.e. to follow the observed data (including the error/noise component) too closely.
- If we are only interested in prediction, then the interpretability of the predictive model may be simply not of interest.
- Flexible models may provide good fit but there is the risk of overfitting.
- Models involving fewer variables are often preferred over more complicated models involving several variables or features.



Interpretability vs flexibility using different statistical/machine learning methods

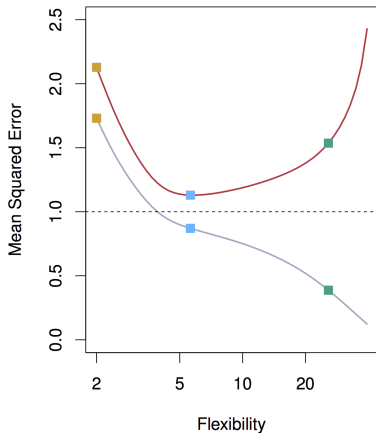
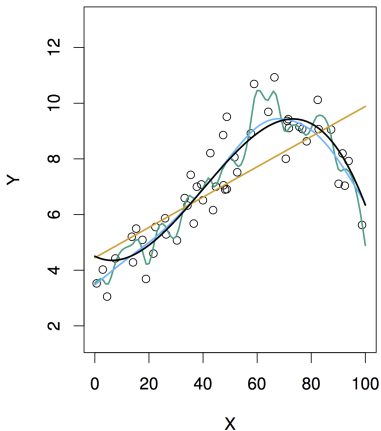
- Suppose we have fit a model  $f(x)$  to some training data  $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , and we wish to assess its performance.
- Compute the **average squared prediction error** over  $\mathcal{D}_{\text{train}}$ :

$$\text{MSE}_{\text{train}} = \frac{1}{n} \sum_{i \in \mathcal{D}_{\text{train}}} [y_i - \hat{f}(\mathbf{x}_i)]^2$$

- Since the same data is used both for “learning” and for “evaluating” the fit of a model, this gives an **optimistic** evaluation of model accuracy.
- If used for selecting the complexity of a statistical model, it is biased toward **overfitting** models.
- Compute the MSE on a **test set**  $\mathcal{D}_{\text{test}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , i.e. a fresh dataset not used for parameters estimation:

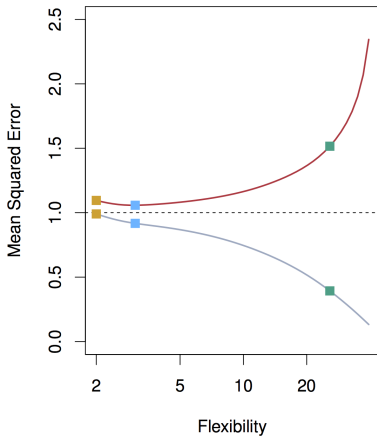
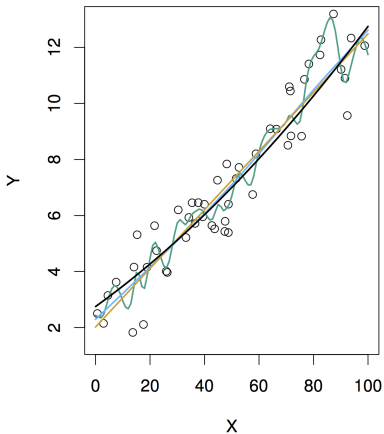
$$\text{MSE}_{\text{test}} = \frac{1}{m} \sum_{i \in \mathcal{D}_{\text{test}}} [y_i - \hat{f}(\mathbf{x}_i)]^2$$

- This is a more realistic measure of how accurately an algorithm is able to predict outcome values for previously unseen data.



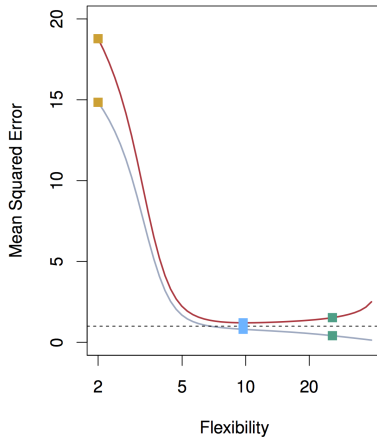
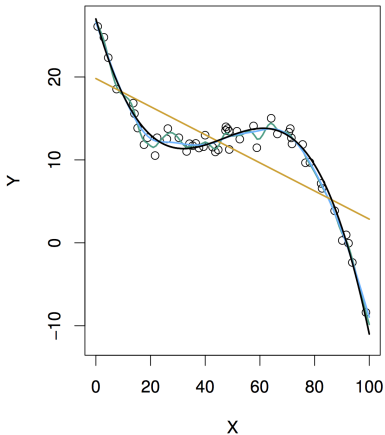
Models: 1) true model — 2) linear regression fit —■— 3) smoothing spline fit with medium flexibility —■— 4) smoothing spline fit with high flexibility —■—

Accuracy measure: (i) training MSE — (ii) test MSE — (iii) irreducible MSE —



Models: 1) true model — 2) linear regression fit —■— 3) smoothing spline fit with medium flexibility —■— 4) smoothing spline fit with high flexibility —■—

Accuracy measure: (i) training MSE — (ii) test MSE — (iii) irreducible MSE —



Models: 1) true model — 2) linear regression fit —■— 3) smoothing spline fit with medium flexibility —■— 4) smoothing spline fit with high flexibility —■—

Accuracy measure: (i) training MSE — (ii) test MSE — (iii) irreducible MSE ----

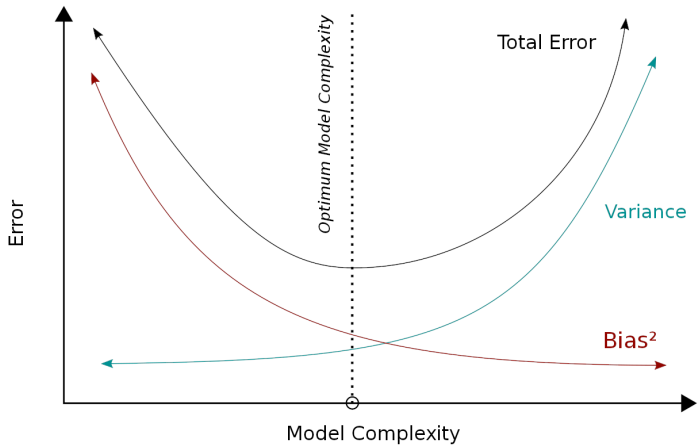


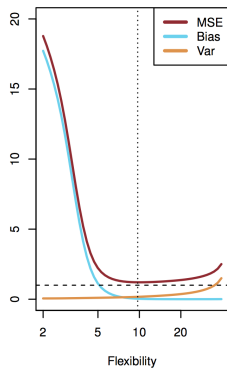
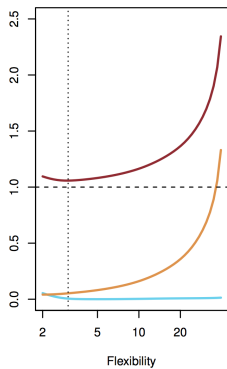
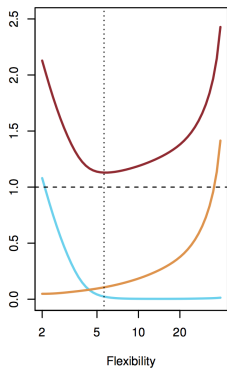
- The expected test error for a new observation value  $\mathbf{x}_0$  can always be decomposed as

$$E[y_0 - \hat{f}(\mathbf{x}_0)]^2 = V[\hat{f}(\mathbf{x}_0)] + B[\hat{f}(\mathbf{x}_0)]^2 + V[\epsilon]$$

where

- $V[\hat{f}(\mathbf{x}_0)]$  is the variance expressing the amount by which  $\hat{f}(\cdot)$  would change if we estimated it using a different training dataset;
  - $B[\hat{f}(\mathbf{x}_0)]$  is the bias expressing the error that is introduced by approximating the data distribution by a statistical model;
  - $V[\epsilon]$  is the irreducible error.
- The expected test error can never be smaller than the irreducible error.
  - In general, more flexible statistical methods have higher variance and smaller bias. On the contrary, simpler models have smaller variance but higher bias.
  - To minimize the expected test error, we need to select a statistical/machine learning method that simultaneously achieves low variance and low bias.





MSE and the bias-variance trade-off

- Suppose that we seek to estimate  $f()$  on the basis of the training observations  $\mathcal{D}_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , where in this case  $y_i \in \{C_1, C_2, \dots, C_K\}$  is the class or label associated with the  $i$ th observation.
- **Training classification error rate** is the proportion of misclassified observations, i.e.

$$\text{CE}_{\text{train}} = \frac{1}{n} \sum_{i \in \mathcal{D}_{\text{train}}} \mathbb{1}(y_i \neq \hat{y}_i)$$

where

- $\hat{y}_i$  is the predicted class label for the  $i$ th observation using  $\hat{f}()$ ;
  - $\mathbb{1}(y_i \neq \hat{y}_i)$  is the indicator function that returns 1 if  $y_i \neq \hat{y}_i$  and 0 otherwise.
- The **test classification error rate** associated with a set of test observations  $\mathcal{D}_{\text{test}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$  is given by

$$\text{CE}_{\text{test}} = \frac{1}{m} \sum_{i \in \mathcal{D}_{\text{test}}} \mathbb{1}(y_i \neq \hat{y}_i)$$

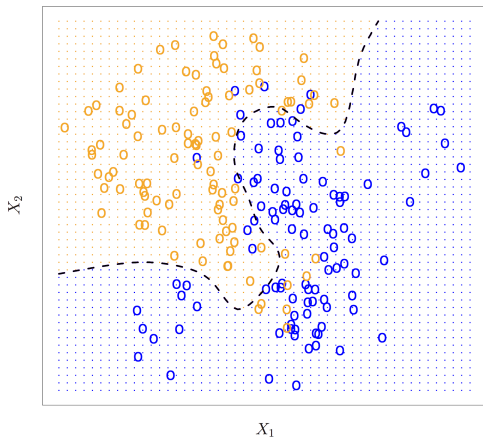
- The test error rate  $CE_{\text{test}}$  is minimized, on average, by a very simple classifier that assigns each observation to the most likely class, given its predictor values.
- According to the **Bayes classifier**, a test observation with predictor vector  $\mathbf{x}_0$  should be assigned to the class  $C_k$  (with  $k = 1, \dots, K$ ) for which

$$\Pr[Y = C_k | \mathbf{x}_0] \quad \text{is maximum}$$

- The Bayes error rate is the lowest possible test error rate produced by the Bayes classifier:
  - error rate at  $\mathbf{x}_0$  is  $1 - \max_k \Pr[Y = C_k | \mathbf{x}_0]$ .
  - overall Bayes error rate is  $1 - E [\max_k \Pr[Y = C_k | X]]$

The Bayes error rate is analogous to the irreducible error for classification tasks.

- The Bayes decision boundary defines the regions in which a test observation will be assigned to one of the  $K$  classes.



Two-class simulated dataset. The dashed line represents the Bayes decision boundary with Bayes error rate  $\approx 13\%$

- For real data, we do not know the conditional distribution of  $Y$  given  $X$ , and so computing the Bayes classifier is impossible.

- In supervised machine learning problems, a training set of  $n$  data points is available,  $\mathcal{D}_{\text{train}} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ , where  $\mathbf{x}_i$  represents the  $p$  features on the  $i$ th observation and  $y_i$  represents the value of the response variable.
- The main goal is to build a model whose predictions  $\hat{y}_i$  are as close as possible to the true response values  $y_i$ .
- A **loss function** aims at measuring model's prediction error:

$$\mathcal{L} = \sum_{i=1}^n \mathcal{L}(y_i, \hat{y}_i)$$

- Properties of a loss function:
  - Continuous and differentiable everywhere.
  - Convex, i.e. only one global minimum point exists, so optimization methods like gradient descent are guaranteed to return the globally optimal solution. In practice, this is hard to achieve, and most loss functions are non-convex (i.e. they have multiple local minima).
  - Symmetric, i.e. the error above the target should cause the same loss as the same error below the target.
  - Computationally efficient, i.e. fast and scalable.

## LOSS FUNCTIONS AND MAXIMUM LIKELIHOOD

---

- Many of the loss functions used in ML are derived from the maximum likelihood principle.
- In maximum likelihood estimation (MLE) we are trying to fit a model with parameters  $\theta$  that maximizes the probability of the observed data given the model:  $\Pr(\mathcal{D}|\theta)$ .

- MLEs are computed as

$$\hat{\theta} = \arg \max_{\theta} \log \Pr(\mathcal{D}|\theta)$$

- Thus, the loss function for a random sample  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$  can be defined as

$$\mathcal{L} = -\log \Pr(\mathcal{D}|\theta) = -\sum_{i=1}^n \log \Pr(\mathbf{x}_i|\theta)$$

so  $\mathcal{L}(y_i, \hat{y}_i) = -\log \Pr(\mathbf{x}_i|\theta)$ .

- Because negative logarithm is a monotonically decreasing function, maximizing the likelihood is equivalent to minimizing the loss.



### SQUARED LOSS

---

- The squared loss is defined as

$$\mathcal{L}_{\text{sq}}(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$$

- This is the loss function used in ordinary least squares (OLS), the most common method for solving linear regression problems.
- Pros:
  - Continuous and differentiable everywhere.
  - Convex (has only one global minimum).
  - Easy to compute.
  - Obtained assuming a Gaussian distribution for the errors.

Cons:

- Sensitive to outliers.

## ABSOLUTE LOSS

---

- The absolute loss is defined as

$$\mathcal{L}_{\text{abs}}(y_i, \hat{y}_i) = |y_i - \hat{y}_i|$$

- Pros:
  - Not overly affected by outliers.
  - Easy to compute.
  - Obtained assuming a Laplace distribution for the errors.

Cons:

- Non-differentiable at 0, which makes it hard to be used by derivative-based optimization methods, such as gradient descent.

## HUBER LOSS

---

- The Huber loss is a combination of squared loss and absolute loss and it is defined as

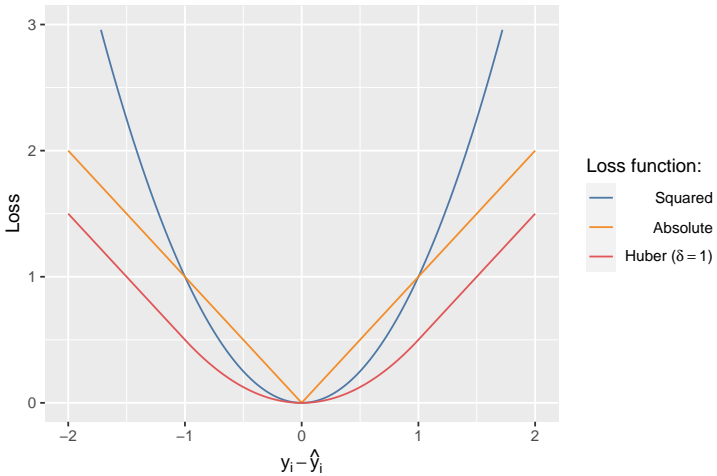
$$\mathcal{L}_{\text{Huber}}(y_i, \hat{y}_i) = \begin{cases} \frac{1}{2}(y_i - \hat{y}_i)^2 & \text{if } |y_i - \hat{y}_i| \leq \delta \\ \delta(|y_i - \hat{y}_i| - \frac{1}{2}\delta) & \text{if } |y_i - \hat{y}_i| > \delta \end{cases}$$

for some hyperparameter  $\delta > 0$ .

- Pros:
  - Continuous and differentiable everywhere.
  - Less sensitive to outliers than squared loss.

Cons:

- Slower to compute.
- Requires tuning of the hyperparameter  $\delta$ .
- Does not have a maximum likelihood interpretation.



## ZERO-ONE LOSS

---

- The simplest loss function is the zero-one loss function defined as

$$\mathcal{L}_{01}(y_i, \hat{y}_i) = \mathbb{1}(y_i \neq \hat{y}_i)$$

where  $y_i \in \{0,1\}$  is the observed class,  $\hat{y}_i$  the corresponding predicted class, and  $\mathbb{1}()$  the indicator function that returns 1 if its argument is true, and 0 otherwise.

- By encoding  $y_i \in \{-1, +1\}$ , the zero-one loss can also be defined as

$$\mathcal{L}_{01}(y_i, s_i) = \mathbb{1}(y_i s_i < 0)$$

where  $s_i \in \mathbb{R}$  is the linear score s.t.  $\Pr(y_i = +1) = 1/(1 + \exp(-s_i))$ .

- This loss function counts the number of prediction errors made by the classifier ([misclassification error](#)).
- Pros:
  - Easy to compute.

Cons:

- Non-differentiable and non-continuous.

## LOG LOSS OR CROSS-ENTROPY LOSS

---

- Denote the binary response as  $y_i \in \{0,1\}$  and the probability of positive case as  $\Pr(y_i = 1) = p_i$ , then the log loss is defines as

$$\mathcal{L}_{\log}(y_i, p_i) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i)$$

- Equivalently, denoting the binary response as  $y_i \in \{-1, +1\}$  and  $s_i \in \mathbb{R}$  the linear score, the log loss can also be defined as

$$\mathcal{L}_{\log}(y_i, s_i) = \log(1 + \exp(-y_i s_i))$$

- Pros:
  - Continuous and differentiable everywhere.
  - Convex (has only one global minimum).
  - Obtained assuming a Bernoulli distribution for the response variable.
  - Loss function used in logistic regression.
  - Easily extended to multi-class classification problems.

Cons:

- Symmetric.

## HINGE LOSS

---

- For the binary response  $y_i \in \{-1, +1\}$ , the hinge loss is defined as

$$\mathcal{L}_{\text{hinge}}(y_i, s_i) = \max(0, 1 - y_i s_i)$$

- Hinge loss is employed by support vector machines (SVM) to obtain a classifier with “maximal margin”:
  - when  $y_i$  and  $s_i$  have the same sign (a correct prediction) and  $s_i \geq 1$  the loss is 0;
  - when  $y_i$  and  $s_i$  have opposite signs, the loss increases linearly with  $s_i$ , and similarly if  $s_i < 1$ , even if it has the same sign (a correct prediction, but not by enough margin).

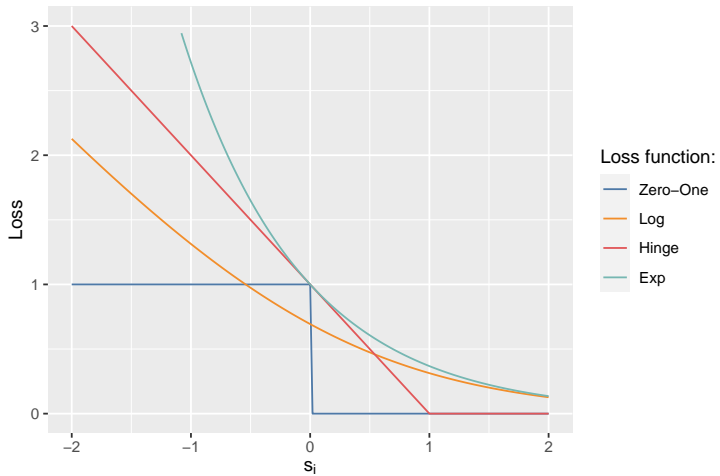
## EXPONENTIAL LOSS

---

- For the binary response  $y_i \in \{-1, +1\}$ , the exponential loss is defined as

$$\mathcal{L}_{\text{exp}}(y_i, s_i) = \exp(-y_i s_i)$$

- A more aggressive loss function which grows exponentially for negative values and is thus very sensitive to wrong predictions.
- Exponential loss is employed by AdaBoost classifier.





- Often models require the tuning of [hyper-parameters](#) ( $k$  in KNN, number of components in GMM, smoothing parameter, lasso and/or ridge parameters, number of hidden layers and number of nodes in NNET, etc.).
- Sometimes we have [no test data](#) available for estimating MSE, classification error, etc.
- In all these cases, a separated [validation dataset](#)  $\mathcal{D}_{\text{val}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^V$  should be used.
- However, instead of setting aside a validation set, it is preferable to use [resampling methods](#).
- [No free lunch theorem](#): *no method/algorithm/model dominates all others over all possible datasets.*
- Realistically, we should decide for any given set of data which method produces the best results.
- This is the most challenging part of statistical/machine learning in practice.

- Resampling methods are a fundamental tool in modern statistics.
- They involve repeatedly drawing samples from a training set and refitting a model of interest on each sample to obtain additional information about the fitted model.
- They can be computationally expensive, because the same statistical model must be fitted multiple times using different subsets of the training data.
- Goals
  - Model assessment (evaluating model's performance)
  - Model selection (selecting the level of flexibility of a model, i.e. hyperparameters tuning)
  - Model inference (provide a measure of accuracy of a parameter estimate or of a given statistical/machine learning method)

- Several possible **performance metrics** can be adopted.
- For regression problems, the error is usually measured by the **root mean square error**:

$$\text{RMSE} = \sqrt{\text{MSE}} = \sqrt{\frac{1}{n} \sum_i (y_i - \hat{f}(x_i))^2}$$

or directly using the MSE.

- For classification problems, the error can be measured by the **classification error**:

$$\text{CE} = \frac{1}{n} \sum_i \mathbb{1}(y_i \neq \hat{y}_i)$$

Many other measures are available: [sensitivity/specificity](#), [ROC-AUC](#), [precision/recal](#), [F-score](#), [log-loss](#) or [cross-entropy](#), [Brier score](#), etc.

- If a validation set is not available, an estimate of the true error must be obtained by [resampling methods](#).

## CROSS-VALIDATION

- Cross-validation is a widely used resampling approach for estimating the performance of a statistical/machine learning model/algorithm.

### V-fold cross-validation

The set of training observations is randomly splitted into  $V$  parts or [folds](#). The model is trained using all but the  $v$ th fold, then the remaining  $v$ th fold is used as validation set. This is done in turn for each fold  $v = 1, \dots, V$ , and then the results are combined.

10-fold cross-validation scheme



- When  $V = n$ , the procedure is called **leave-one-out cross-validation** (LOOCV), because we leave out one data point at a time.

## BOOTSTRAP

---

- The **bootstrap** is a flexible and powerful statistical tool that can be used to quantify the uncertainty associated with a given estimator (e.g. standard errors or confidence intervals for regression coefficients) or the predictions provided by a statistical/machine learning method.
- Bootstrap takes random samples with replacement of the same size as the original data set.
- Since sampling is made with replacement, some observations may be selected more than once and each observation has a 63.2% chance of showing up at least once.

*The probability for an observation of not being selected in any of  $n$  draws from  $n$  samples with replacement is  $(1 - 1/n)^n$ .*

*Then  $\lim_{n \rightarrow \infty} (1 - 1/n)^n = e^{-1} \approx 0.368$ , and the probability of being selected at least once is  $1 - e^{-1} \approx 0.632$ .*

- The observations not selected (approximately 1/3 of the sample) are usually referred to as the out-of-bag observations.

## **BOOTSTRAP ALGORITHM FOR A CLASSIFICATION TASK**

---

1. Create a bootstrap sample by random sampling with replacement;
2. Fit a classifier using the bootstrap sample as training set;
3. Predict out-of-bag observations to get bootstrap classification error;
4. Repeated steps 1-3 multiple times (usually 30 – 100) and then combine the results.

## **BIAS-VARIANCE TRADE-OFF FOR BOOTSTRAP**

---

- The bootstrap estimates of error rate have less variability than V-fold CV, but larger bias (similar to 2-fold CV).  
If the training set size is small, this bias may be problematic, but will decrease as the training set sample size becomes larger.
- The [“632” bootstrap](#) method tries to reduce the bias by creating a performance estimate that is a combination of the simple bootstrap estimate and the estimate from predicting the training set:

$$(0.632 \times \text{bootstrap error rate}) + (0.368 \times \text{training error rate})$$

## ONE STANDARD ERROR RULE

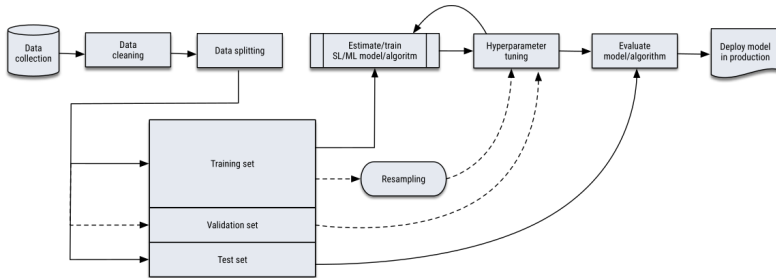
---

- Instead of selecting the model with the “best” tuning parameter value, other schemes for selecting a single model can be used.
- A popular choice is the so-called “one standard error rule”:

*“all else equal (up to one standard error), go for the simpler (more regularized/parsimonious) model”*

- In practice:
  - the model with the best performance value is identified;
  - an estimate of the standard error of performance is computed by a resampling method;
  - the final model is the simplest model whose estimated performance is within one standard error from the best model performance.

# STATISTICAL/MACHINE LEARNING PIPELINE





- Machine Learning and Statistical Learning play a vital role in the broader field of Artificial Intelligence.
- They enable computers to learn from data, make predictions, and solve complex tasks without being explicitly programmed for a specific task.
- These techniques have widespread applications across various domains, revolutionizing industries and improving decision-making processes.

### References (for starters...)

- James G., Witten D., Hastie T. and Tibshirani R. (2021) [An Introduction to Statistical Learning: with Applications in R](https://www.statlearning.com), 2nd edition, Springer-Verlag  
<https://www.statlearning.com>
- Murphy K. P. (2022) [Probabilistic Machine Learning: An Introduction](https://probml.github.io/pml-book/book1.html), MIT Press  
<https://probml.github.io/pml-book/book1.html>
- Murphy K. P. (2023) [Probabilistic Machine Learning: Advanced Topics](https://probml.github.io/pml-book/book2.html), MIT Press  
<https://probml.github.io/pml-book/book2.html>



• **THANK YOU** •

■  
ANY QUESTION?