# 1 Supervised Learning

## 1.1 Introduction to Supervised Learning

Given a set of data points $\{x^{(1)}, ..., x^{(m)}\}$ associated to a set of outcomes $\{y^{(1)}, ..., y^{(m)}\}$, we want to build a classifier that learns how to predict $y$ from $x$.

❏ **Type of prediction** – The different types of predictive models are summed up in the table below:

| | Regression | Classifier |
|---|---|---|
| **Outcome** | Continuous | Class |
| **Examples** | Linear regression | Logistic regression, SVM, Naive Bayes |

❏ **Type of model** – The different models are summed up in the table below:

| | Discriminative model | Generative model |
|---|---|---|
| **Goal** | Directly estimate $P(y|x)$ | Estimate $P(x|y)$ to deduce $P(y|x)$ |
| **What's learned** | Decision boundary | Probability distributions of the data |
| **Illustration** | | |
| **Examples** | Regressions, SVMs | GDA, Naive Bayes |

## 1.2 Notations and general concepts

❏ **Hypothesis** – The hypothesis is noted $h_\theta$ and is the model that we choose. For a given input data $x^{(i)}$, the model prediction output is $h_\theta(x^{(i)})$.

❏ **Loss function** – A loss function is a function $L : (z,y) \in \mathbb{R} \times Y \longmapsto L(z,y) \in \mathbb{R}$ that takes as inputs the predicted value $z$ corresponding to the real data value $y$ and outputs how different they are. The common loss functions are summed up in the table below:
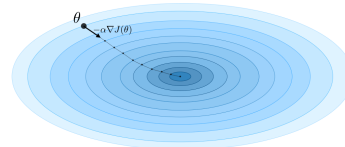
| Least squared | Logistic | Hinge | Cross-entropy |
|---|---|---|---|
| $\frac{1}{2}(y - z)^2$ | $\log(1 + \exp(-yz))$ | $\max(0, 1 - yz)$ | $-\left[y\log(z) + (1 - y)\log(1 - z)\right]$ |
| Linear regression | Logistic regression | SVM | Neural Network |

❏ **Cost function** – The cost function $J$ is commonly used to assess the performance of a model, and is defined with the loss function $L$ as follows:

$$J(\theta) = \sum_{i=1}^{m} L(h_\theta(x^{(i)}), y^{(i)})$$

❏ **Gradient descent** – By noting $\alpha \in \mathbb{R}$ the learning rate, the update rule for gradient descent is expressed with the learning rate and the cost function $J$ as follows:

$$\theta \longleftarrow \theta - \alpha \nabla J(\theta)$$



*Remark: Stochastic gradient descent (SGD) is updating the parameter based on each training example, and batch gradient descent is on a batch of training examples.*

❏ **Likelihood** – The likelihood of a model $L(\theta)$ given parameters $\theta$ is used to find the optimal parameters $\theta$ through maximizing the likelihood. In practice, we use the log-likelihood $\ell(\theta) = \log(L(\theta))$ which is easier to optimize. We have:

$$\theta^{\text{opt}} = \arg \max_\theta L(\theta)$$

❏ **Newton's algorithm** – The Newton's algorithm is a numerical method that finds $\theta$ such that $\ell'(\theta) = 0$. Its update rule is as follows:

$$\theta \leftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

*Remark: the multidimensional generalization, also known as the Newton-Raphson method, has the following update rule:*

$$\theta \leftarrow \theta - \left(\nabla_\theta^2 \ell(\theta)\right)^{-1} \nabla_\theta \ell(\theta)$$

## 1.3 Linear models

### 1.3.1 Linear regression

We assume here that $y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$

❏ **Normal equations** – By noting $X$ the matrix design, the value of $\theta$ that minimizes the cost function is a closed-form solution such that:

$$\theta = (X^T X)^{-1} X^T y$$

❏ **LMS algorithm** – By noting $\alpha$ the learning rate, the update rule of the Least Mean Squares (LMS) algorithm for a training set of $m$ data points, which is also known as the Widrow-Hoff learning rule, is as follows:

$$\forall j, \quad \theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^{m} \left[y^{(i)} - h_\theta(x^{(i)})\right] x_j^{(i)}$$

*Remark: the update rule is a particular case of the gradient ascent.*

❏ **LWR** – Locally Weighted Regression, also known as LWR, is a variant of linear regression that weights each training example in its cost function by $w^{(i)}(x)$, which is defined with parameter $\tau \in \mathbb{R}$ as:

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$

### 1.3.2 Classification and logistic regression

❏ **Sigmoid function** – The sigmoid function $g$, also known as the logistic function, is defined as follows:

$$\forall z \in \mathbb{R}, \quad g(z) = \frac{1}{1 + e^{-z}} \in ]0,1[$$

❏ **Logistic regression** – We assume here that $y|x; \theta \sim \text{Bernoulli}(\phi)$. We have the following form:

$$\phi = p(y = 1|x; \theta) = \frac{1}{1 + \exp(-\theta^T x)} = g(\theta^T x)$$

*Remark: there is no closed form solution for the case of logistic regressions.*

❏ **Softmax regression** – A softmax regression, also called a multiclass logistic regression, is used to generalize logistic regression when there are more than 2 outcome classes. By convention, we set $\theta_K = 0$, which makes the Bernoulli parameter $\phi_i$ of each class $i$ equal to:

$$\phi_i = \frac{\exp(\theta_i^T x)}{\sum_{j=1}^{K} \exp(\theta_j^T x)}$$

### 1.3.3 Generalized Linear Models

❏ **Exponential family** – A class of distributions is said to be in the exponential family if it can be written in terms of a natural parameter, also called the canonical parameter or link function, $\eta$, a sufficient statistic $T(y)$ and a log-partition function $a(\eta)$ as follows:

$$p(y; \eta) = b(y) \exp(\eta T(y) - a(\eta))$$

*Remark: we will often have $T(y) = y$. Also, $\exp(-a(\eta))$ can be seen as a normalization parameter that will make sure that the probabilities sum to one.*

Here are the most common exponential distributions summed up in the following table:

| Distribution | $\eta$ | $T(y)$ | $a(\eta)$ | $b(y)$ |
|---|---|---|---|---|
| Bernoulli | $\log\left(\frac{\phi}{1-\phi}\right)$ | $y$ | $\log(1 + \exp(\eta))$ | $1$ |
| Gaussian | $\mu$ | $y$ | $\frac{\eta^2}{2}$ | $\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$ |
| Poisson | $\log(\lambda)$ | $y$ | $e^\eta$ | $\frac{1}{y!}$ |
| Geometric | $\log(1 - \phi)$ | $y$ | $\log\left(\frac{e^\eta}{1-e^\eta}\right)$ | $1$ |

❏ **Assumptions of GLMs** – Generalized Linear Models (GLM) aim at predicting a random variable $y$ as a function fo $x \in \mathbb{R}^{n+1}$ and rely on the following 3 assumptions:

$$(1) \quad \boxed{y|x; \theta \sim \text{ExpFamily}(\eta)} \qquad (2) \quad \boxed{h_\theta(x) = E[y|x; \theta]} \qquad (3) \quad \boxed{\eta = \theta^T x}$$

*Remark: ordinary least squares and logistic regression are special cases of generalized linear models.*
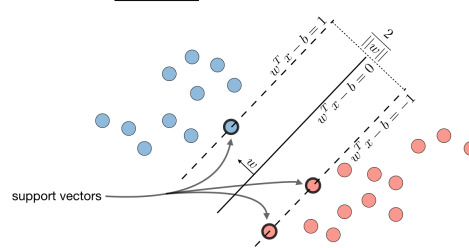
## 1.4 Support Vector Machines

The goal of support vector machines is to find the line that maximizes the minimum distance to the line.

❏ **Optimal margin classifier** – The optimal margin classifier $h$ is such that:

$$h(x) = \text{sign}(w^T x - b)$$

where $(w, b) \in \mathbb{R}^n \times \mathbb{R}$ is the solution of the following optimization problem:

$$\min \frac{1}{2}||w||^2 \quad \text{such that} \quad y^{(i)}(w^T x^{(i)} - b) \geqslant 1$$
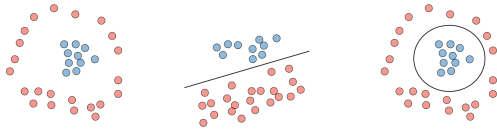


*Remark: the line is defined as $w^T x - b = 0$.*

❏ **Hinge loss** – The hinge loss is used in the setting of SVMs and is defined as follows:

$$L(z, y) = [1 - yz]_+ = \max(0, 1 - yz)$$

❒ **Kernel** – Given a feature mapping $\phi$, we define the kernel $K$ to be defined as:

$$K(x,z) = \phi(x)^T \phi(z)$$

In practice, the kernel $K$ defined by $K(x,z) = \exp\left(-\frac{||x-z||^2}{2\sigma^2}\right)$ is called the Gaussian kernel and is commonly used.



Non-linear separability $\Longrightarrow$ Use of a kernel mapping $\phi$ $\Longrightarrow$ Decision boundary in the original space

*Remark: we say that we use the "kernel trick" to compute the cost function using the kernel because we actually don't need to know the explicit mapping $\phi$, which is often very complicated. Instead, only the values $K(x,z)$ are needed.*

❒ **Lagrangian** – We define the Lagrangian $\mathcal{L}(w,b)$ as follows:

$$\mathcal{L}(w,b) = f(w) + \sum_{i=1}^{l} \beta_i h_i(w)$$

*Remark: the coefficients $\beta_i$ are called the Lagrange multipliers.*

### 1.5   Generative Learning

A generative model first tries to learn how the data is generated by estimating $P(x|y)$, which we can then use to estimate $P(y|x)$ by using Bayes' rule.

#### 1.5.1   Gaussian Discriminant Analysis

❒ **Setting** – The Gaussian Discriminant Analysis assumes that $y$ and $x|y = 0$ and $x|y = 1$ are such that:

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \quad \text{and} \quad x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$$

❒ **Estimation** – The following table sums up the estimates that we find when maximizing the likelihood:

| $\widehat{\phi}$ | $\widehat{\mu_j} \quad (j = 0,1)$ | $\widehat{\Sigma}$ |
|---|---|---|
| $\frac{1}{m}\sum_{i=1}^{m} 1_{\{y^{(i)}=1\}}$ | $\frac{\sum_{i=1}^{m} 1_{\{y^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^{m} 1_{\{y^{(i)}=j\}}}$ | $\frac{1}{m}\sum_{i=1}^{m} (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$ |

#### 1.5.2   Naive Bayes

❒ **Assumption** – The Naive Bayes model supposes that the features of each data point are all independent:

$$P(x|y) = P(x_1, x_2, \dots | y) = P(x_1|y)P(x_2|y)\dots = \prod_{i=1}^{n} P(x_i|y)$$

❒ **Solutions** – Maximizing the log-likelihood gives the following solutions, with $k \in \{0,1\}$, $l \in [\![1,L]\!]$

$$P(y = k) = \frac{1}{m} \times \#\{j | y^{(j)} = k\} \quad \text{and} \quad P(x_i = l|y = k) = \frac{\#\{j | y^{(j)} = k \text{ and } x_i^{(j)} = l\}}{\#\{j | y^{(j)} = k\}}$$

*Remark: Naive Bayes is widely used for text classification and spam detection.*

### 1.6   Tree-based and ensemble methods

These methods can be used for both regression and classification problems.

❒ **CART** – Classification and Regression Trees (CART), commonly known as decision trees, can be represented as binary trees. They have the advantage to be very interpretable.

❒ **Random forest** – It is a tree-based technique that uses a high number of decision trees built out of randomly selected sets of features. Contrary to the simple decision tree, it is highly uninterpretable but its generally good performance makes it a popular algorithm.

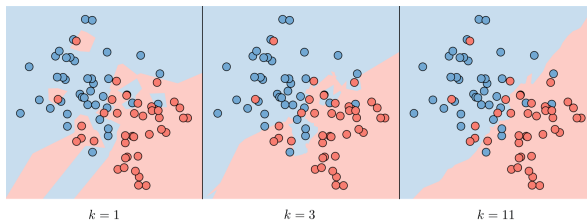*Remark: random forests are a type of ensemble methods.*

❒ **Boosting** – The idea of boosting methods is to combine several weak learners to form a stronger one. The main ones are summed up in the table below:

| Adaptive boosting | Gradient boosting |
|---|---|
| - High weights are put on errors to improve at the next boosting step <br> - Known as Adaboost | - Weak learners trained on remaining errors |

### 1.7   Other non-parametric approaches

❒ **$k$-nearest neighbors** – The $k$-nearest neighbors algorithm, commonly known as $k$-NN, is a non-parametric approach where the response of a data point is determined by the nature of its $k$ neighbors from the training set. It can be used in both classification and regression settings.
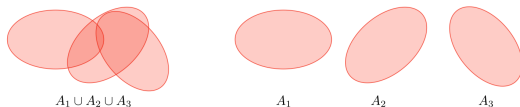
*Remark: The higher the parameter $k$, the higher the bias, and the lower the parameter $k$, the higher the variance.*

---

$k = 1$      $k = 3$      $k = 11$

### 1.8   Learning Theory

❒ **Union bound** – Let $A_1, \dots, A_k$ be $k$ events. We have:

$$P(A_1 \cup \dots \cup A_k) \leqslant P(A_1) + \dots + P(A_k)$$



$A_1 \cup A_2 \cup A_3$      $A_1$    $A_2$    $A_3$

❒ **Hoeffding inequality** – Let $Z_1, \dots, Z_m$ be $m$ iid variables drawn from a Bernoulli distribution of parameter $\phi$. Let $\widehat{\phi}$ be their sample mean and $\gamma > 0$ fixed. We have:

$$P(|\phi - \widehat{\phi}| > \gamma) \leqslant 2 \exp(-2\gamma^2 m)$$

*Remark: this inequality is also known as the Chernoff bound.*

❒ **Training error** – For a given classifier $h$, we define the training error $\widehat{\epsilon}(h)$, also known as the empirical risk or empirical error, to be as follows:

$$\widehat{\epsilon}(h) = \frac{1}{m}\sum_{i=1}^{m} 1_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

❒ **Probably Approximately Correct (PAC)** – PAC is a framework under which numerous results on learning theory were proved, and has the following set of assumptions:

- the training and testing sets follow the same distribution

- the training examples are drawn independently

❒ **Shattering** – Given a set $S = \{x^{(1)}, \dots, x^{(d)}\}$, and a set of classifiers $\mathcal{H}$, we say that $\mathcal{H}$ shatters $S$ if for any set of labels $\{y^{(1)}, \dots, y^{(d)}\}$, we have:

$$\exists h \in \mathcal{H}, \quad \forall i \in [\![1,d]\!], \quad h(x^{(i)}) = y^{(i)}$$

❒ **Upper bound theorem** – Let $\mathcal{H}$ be a finite hypothesis class such that $|\mathcal{H}| = k$ and let $\delta$ and the sample size $m$ be fixed. Then, with probability of at least $1 - \delta$, we have:

$$\epsilon(\widehat{h}) \leqslant \left(\min_{h \in \mathcal{H}} \epsilon(h)\right) + 2\sqrt{\frac{1}{2m}\log\left(\frac{2k}{\delta}\right)}$$

❒ **VC dimension** – The Vapnik-Chervonenkis (VC) dimension of a given infinite hypothesis class $\mathcal{H}$, noted $\text{VC}(\mathcal{H})$ is the size of the largest set that is shattered by $\mathcal{H}$.

*Remark: the VC dimension of $\mathcal{H} = \{$set of linear classifiers in 2 dimensions$\}$ is 3.*



❒ **Theorem (Vapnik)** – Let $\mathcal{H}$ be given, with $\text{VC}(\mathcal{H}) = d$ and $m$ the number of training examples. With probability at least $1 - \delta$, we have:
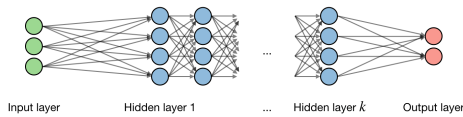
$$\epsilon(\widehat{h}) \leqslant \left(\min_{h \in \mathcal{H}} \epsilon(h)\right) + O\left(\sqrt{\frac{d}{m}\log\left(\frac{m}{d}\right) + \frac{1}{m}\log\left(\frac{1}{\delta}\right)}\right)$$

## 3   Deep Learning

### 3.1   Neural Networks

Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks.

❏ **Architecture** – The vocabulary around neural networks architectures is described in the figure below:



Input layer    Hidden layer 1    ...    Hidden layer $k$    Output layer

By noting $i$ the $i^{th}$ layer of the network and $j$ the $j^{th}$ hidden unit of the layer, we have:

$$z_j^{[i]} = w_j^{[i]T} x + b_j^{[i]}$$

where we note $w, b, z$ the weight, bias and output respectively.

❏ **Activation function** – Activation functions are used at the end of a hidden unit to introduce non-linear complexities to the model. Here are the most common ones:

| Sigmoid | Tanh | ReLU | Leaky ReLU |
|---|---|---|---|
| $g(z) = \dfrac{1}{1 + e^{-z}}$ | $g(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | $g(z) = \max(0, z)$ | $g(z) = \max(\epsilon z, z)$ with $\epsilon \ll 1$ |
| | | | |

❏ **Cross-entropy loss** – In the context of neural networks, the cross-entropy loss $L(z,y)$ is commonly used and is defined as follows:

$$L(z,y) = -\Big[ y \log(z) + (1 - y) \log(1 - z) \Big]$$

❏ **Learning rate** – The learning rate, often noted $\eta$, indicates at which pace the weights get updated. This can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.

❏ **Backpropagation** – Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to weight $w$ is computed using chain rule and is of the following form:

$$\frac{\partial L(z,y)}{\partial w} = \frac{\partial L(z,y)}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}$$

As a result, the weight is updated as follows:

$$w \longleftarrow w - \eta \frac{\partial L(z,y)}{\partial w}$$

❏ **Updating weights** – In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data.

- Step 2: Perform forward propagation to obtain the corresponding loss.

- Step 3: Backpropagate the loss to get the gradients.

- Step 4: Use the gradients to update the weights of the network.

❏ **Dropout** – Dropout is a technique meant at preventing overfitting the training data by dropping out units in a neural network. In practice, neurons are either dropped with probability $p$ or kept with probability $1 - p$.

### 3.2   Convolutional Neural Networks

❏ **Convolutional layer requirement** – By noting $W$ the input volume size, $F$ the size of the convolutional layer neurons, $P$ the amount of zero padding, then the number of neurons $N$ that fit in a given volume is such that:

$$N = \frac{W - F + 2P}{S} + 1$$

❏ **Batch normalization** – It is a step of hyperparameter $\gamma, \beta$ that normalizes the batch $\{x_i\}$. By noting $\mu_B, \sigma_B^2$ the mean and variance of that we want to correct to the batch, it is done as follows:

$$x_i \longleftarrow \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

It is usually done after a fully connected/convolutional layer and before a non-linearity layer and aims at allowing higher learning rates and reducing the strong dependence on initialization.

### 3.3   Recurrent Neural Networks

❏ **Types of gates** – Here are the different types of gates that we encounter in a typical recurrent neural network:

| Input gate | Forget gate | Output gate | Gate |
|---|---|---|---|
| Write to cell or not? | Erase a cell or not? | Reveal a cell or not? | How much writing? |

❏ **LSTM** – A long short-term memory (LSTM) network is a type of RNN model that avoids the vanishing gradient problem by adding 'forget' gates.

## 4   Machine Learning Tips and Tricks

### 4.1   Metrics

Given a set of data points $\{x^{(1)}, ..., x^{(m)}\}$, where each $x^{(i)}$ has $n$ features, associated to a set of outcomes $\{y^{(1)}, ..., y^{(m)}\}$, we want to assess a given classifier that learns how to predict $y$ from $x$.

#### 4.1.1   Classification

In a context of a binary classification, here are the main metrics that are important to track to assess the performance of the model.

❏ **Confusion matrix** – The confusion matrix is used to have a more complete picture when assessing the performance of a model. It is defined as follows:
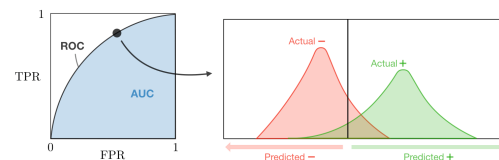


❏ **Main metrics** – The following metrics are commonly used to assess the performance of classification models:

| Metric | Formula | Interpretation |
|---|---|---|
| Accuracy | $\dfrac{TP + TN}{TP + TN + FP + FN}$ | Overall performance of model |
| Precision | $\dfrac{TP}{TP + FP}$ | How accurate the positive predictions are |
| Recall Sensitivity | $\dfrac{TP}{TP + FN}$ | Coverage of actual positive sample |
| Specificity | $\dfrac{TN}{TN + FP}$ | Coverage of actual negative sample |
| F1 score | $\dfrac{2TP}{2TP + FP + FN}$ | Hybrid metric useful for unbalanced classes |

❏ **ROC** – The receiver operating curve, also noted ROC, is the plot of TPR versus FPR by varying the threshold. These metrics are are summed up in the table below:

| Metric | Formula | Equivalent |
|---|---|---|
| True Positive Rate TPR | $\dfrac{TP}{TP + FN}$ | Recall, sensitivity |
| False Positive Rate FPR | $\dfrac{FP}{TN + FP}$ | 1-specificity |

❏ **AUC** – The area under the receiving operating curve, also noted AUC or AUROC, is the area below the ROC as shown in the following figure:



#### 4.1.2   Regression

❏ **Basic metrics** – Given a regression model $f$, the following metrics are commonly used to assess the performance of the model:

| Total sum of squares | Explained sum of squares | Residual sum of squares |
|---|---|---|
| $SS_{tot} = \displaystyle\sum_{i=1}^{m}(y_i - \overline{y})^2$ | $SS_{reg} = \displaystyle\sum_{i=1}^{m}(f(x_i) - \overline{y})^2$ | $SS_{res} = \displaystyle\sum_{i=1}^{m}(y_i - f(x_i))^2$ |

❏ **Coefficient of determination** – The coefficient of determination, often noted $R^2$ or $r^2$, provides a measure of how well the observed outcomes are replicated by the model and is defined as follows:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

❏ **Main metrics** – The following metrics are commonly used to assess the performance of regression models, by taking into account the number of variables $n$ that they take into consideration:

| Mallow's Cp | AIC | BIC | Adjusted $R^2$ |
|---|---|---|---|
| $\dfrac{SS_{res} + 2(n+1)\widehat{\sigma}^2}{m}$ | $2\Big[(n+2) - \log(L)\Big]$ | $\log(m)(n+2) - 2\log(L)$ | $1 - \dfrac{(1 - R^2)(m-1)}{m - n - 1}$ |

where $L$ is the likelihood and $\widehat{\sigma}^2$ is an estimate of the variance associated with each response.

### 4.2  Model selection

❑ **Vocabulary** – When selecting a model, we distinguish 3 different parts of the data that we have as follows:

| Training set | Validation set | Testing set |
|---|---|---|
| - Model is trained<br>- Usually 80% of the dataset | - Model is assessed<br>- Usually 20% of the dataset<br>- Also called hold-out<br>or development set | - Model gives predictions<br>- Unseen data |

Once the model has been chosen, it is trained on the entire dataset and tested on the unseen test set. These are represented in the figure below:



❑ **Cross-validation** – Cross-validation, also noted CV, is a method that is used to select a model that does not rely too much on the initial training set. The different types are summed up in the table below:

| $k$-fold | Leave-$p$-out |
|---|---|
| - Training on $k-1$ folds and<br>assessment on the remaining one<br>- Generally $k = 5$ or $10$ | - Training on $n-p$ observations and<br>assessment on the $p$ remaining ones<br>- Case $p = 1$ is called leave-one-out |

The most commonly used method is called $k$-fold cross-validation and splits the training data into $k$ folds to validate the model on one fold while training the model on the $k-1$ other folds, all of this $k$ times. The error is then averaged over the $k$ folds and is named cross-validation error.



❑ **Regularization** – The regularization procedure aims at avoiding the model to overfit the data and thus deals with high variance issues. The following table sums up the different types of commonly used regularization techniques:

| LASSO | Ridge | Elastic Net |
|---|---|---|
| - Shrinks coefficients to 0<br>- Good for variable selection | Makes coefficients smaller | Tradeoff between variable<br>selection and small coefficients |
|  |  |  |
| $... + \lambda\|\|\theta\|\|_1$<br>$\lambda \in \mathbb{R}$ | $... + \lambda\|\|\theta\|\|_2^2$<br>$\lambda \in \mathbb{R}$ | $... + \lambda\left[(1-\alpha)\|\|\theta\|\|_1 + \alpha\|\|\theta\|\|_2^2\right]$<br>$\lambda \in \mathbb{R}, \quad \alpha \in [0,1]$ |

❑ **Model selection** – Train model on training set, then evaluate on the development set, then pick best performance model on the development set, and retrain all of that model on the whole training set.

### 4.3  Diagnostics

❑ **Bias** – The bias of a model is the difference between the expected prediction and the correct model that we try to predict for given data points.

❑ **Variance** – The variance of a model is the variability of the model prediction for given data points.

❑ **Bias/variance tradeoff** – The simpler the model, the higher the bias, and the more complex the model, the higher the variance.

| | Underfitting | Just right | Overfitting |
|---|---|---|---|
| Symptoms | - High training error<br>- Training error close<br>to test error<br>- High bias | - Training error<br>slightly lower than<br>test error | - Low training error<br>- Training error much<br>lower than test error<br>- High variance |
| Regression |  |  |  |