

# Supervised Learning cheatsheet

[stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning](https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning)

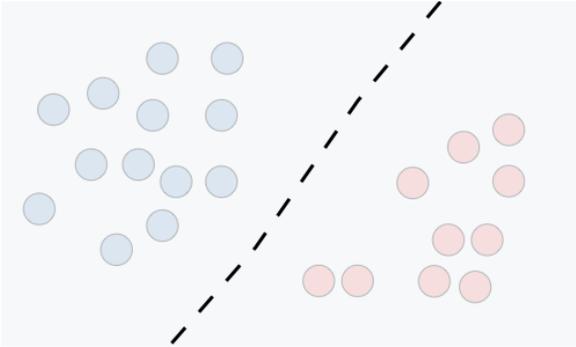
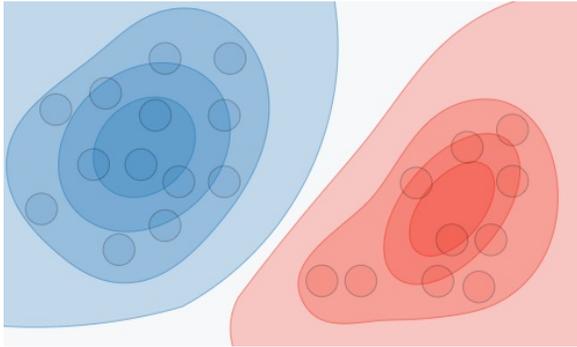
## Introduction to Supervised Learning

Given a set of data points  $\{x(1), \dots, x(m)\}$  associated to a set of outcomes  $\{y(1), \dots, y(m)\}$ , 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
<b>Outcome</b>	Continuous	Class
<b>Examples</b>	Linear regression	Logistic regression, SVM, Naive Bayes

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

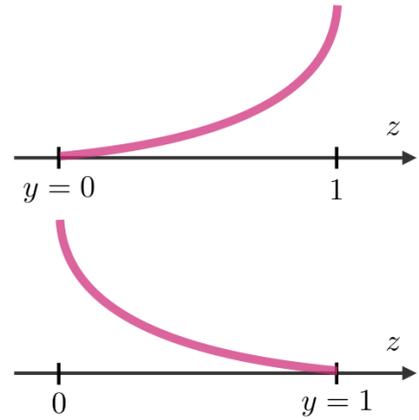
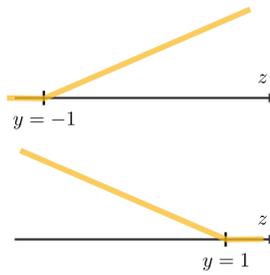
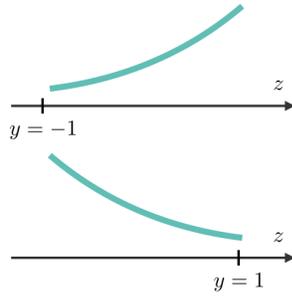
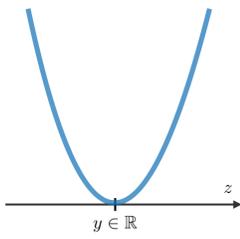
	Discriminative model	Generative model
<b>Goal</b>	Directly estimate $P(y x)P(y x)$	Estimate $P(x y)P(x y)$ to then deduce $P(y x)P(y x)$
<b>What's learned</b>	Decision boundary	Probability distributions of the data
<b>Illustration</b>		
<b>Examples</b>	Regressions, SVMs	GDA, Naive Bayes

## 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  $L: (z, y) \in \mathbb{R} \times \mathbb{Y} \rightarrow \mathbb{R}$  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 error	Logistic loss	Hinge loss	Cross-entropy
$\frac{1}{2}(y-z)^2$	$-\log(1+\exp(-yz))$	$\max(0, 1-yz)$	$-\log(z) - (1-y)\log(1-z)$



Linear regression

Logistic regression

SVM

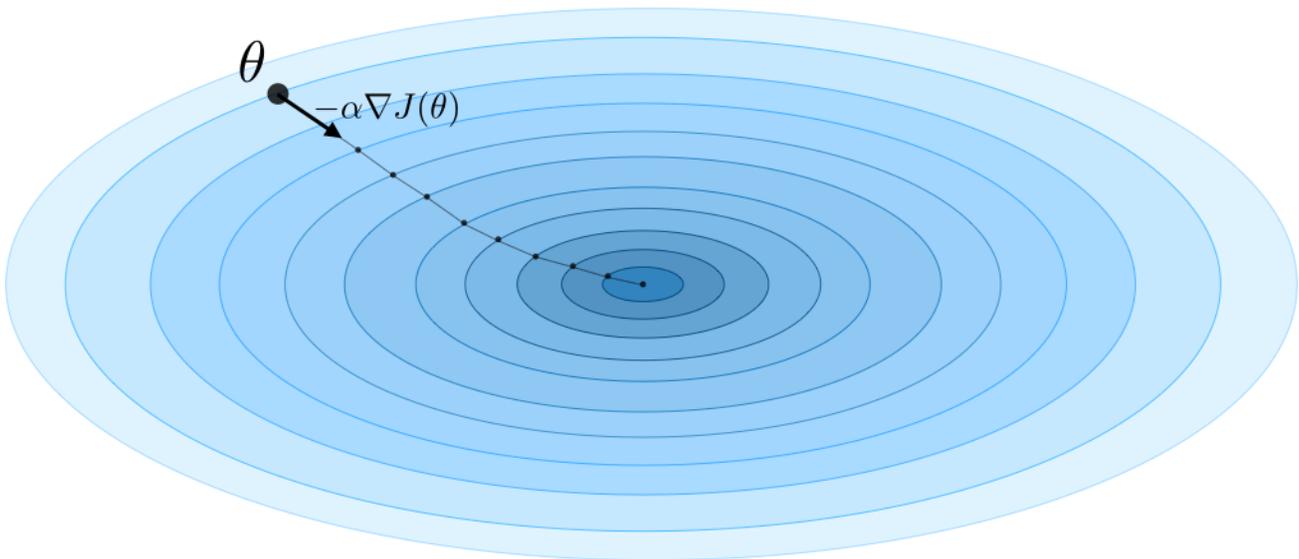
Neural Network

**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:

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

**Gradient descent** — By noting  $\alpha \in \mathbb{R}, \alpha > 0$  the learning rate, the update rule for gradient descent is expressed with the learning rate and the cost function  $J$  as follows:

$$\theta \leftarrow \theta - \alpha \nabla J(\theta) \quad \theta \leftarrow \theta - \alpha \nabla \ell(\theta)$$



*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:

$$\theta_{\text{opt}} = \arg \max_{\theta} \ell(\theta)$$

**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:

$$\theta \leftarrow \theta - \ell''(\theta)^{-1} \ell'(\theta)$$

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

$$\theta \leftarrow \theta - (\nabla^2 \ell(\theta))^{-1} \nabla \ell(\theta)$$

## Linear models

### Linear regression

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

**Normal equations** — By noting  $X$  the design matrix, 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:

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

*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  $w(i)$ , which is defined with parameter  $\tau \in \mathbb{R}$  as:

$$w(i) = \exp(-(x(i) - x)^2 / (2\tau^2))$$

## Classification and logistic regression

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

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

**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:

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

## 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:

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

*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	$\log(\phi) - \phi$	$y$	$\log(1 + \exp(\eta))$	$1$
Gaussian	$\mu$	$y^2$	$\frac{1}{2}\eta^2$	$-\frac{1}{2}y^2$
Poisson	$\log(\lambda)$	$y$	$e^\eta$	$y \log y$
Geometric	$\log(1 - \phi)$	$y$	$-\log(1 - \phi)$	$1 - \phi$

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

$$(1) y | x; \theta \sim \text{ExpFamily}(\eta)$$

$$(2) \eta = \theta^T x$$

$$(3) \eta = \theta^T x$$

*Remark: ordinary least squares and logistic regression are special cases of generalized linear models.*

## Support Vector Machines

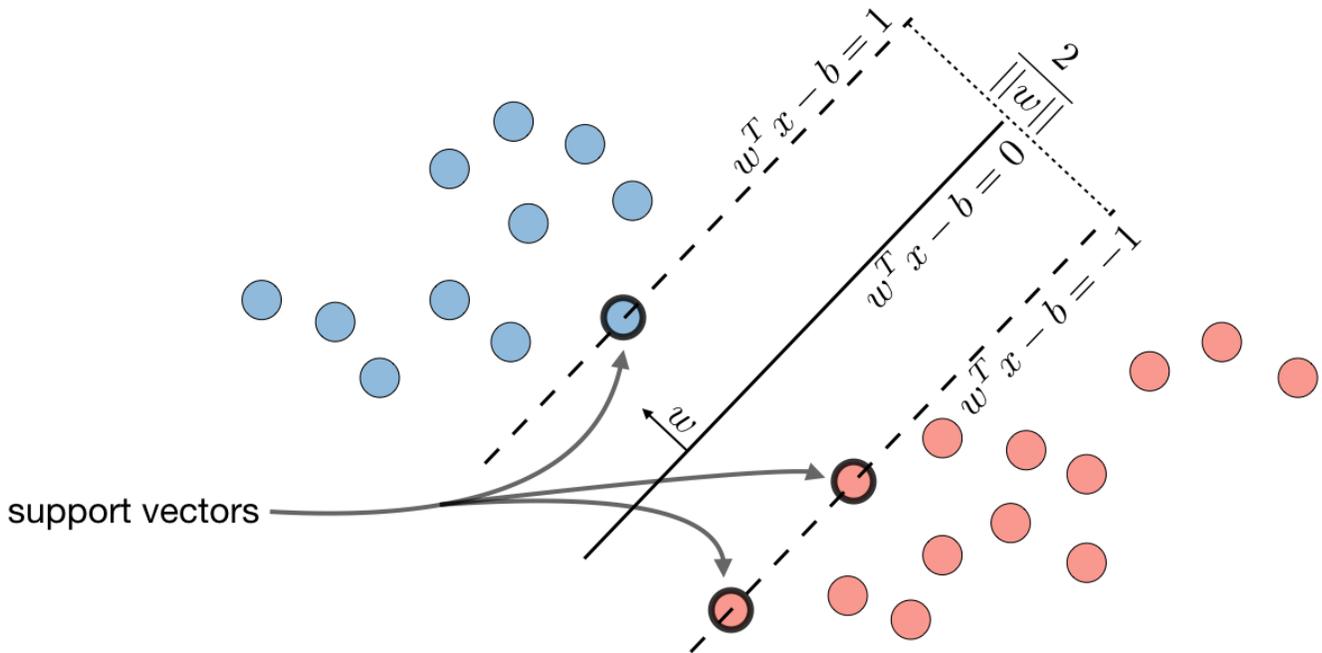
The goal of support vector machines is to find the line that maximizes the minimum distance to the line.

**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:

$$\min_{\|w\|_2} \sum_i \max(0, 1 - y_i(w^T x^{(i)} - b))$$



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

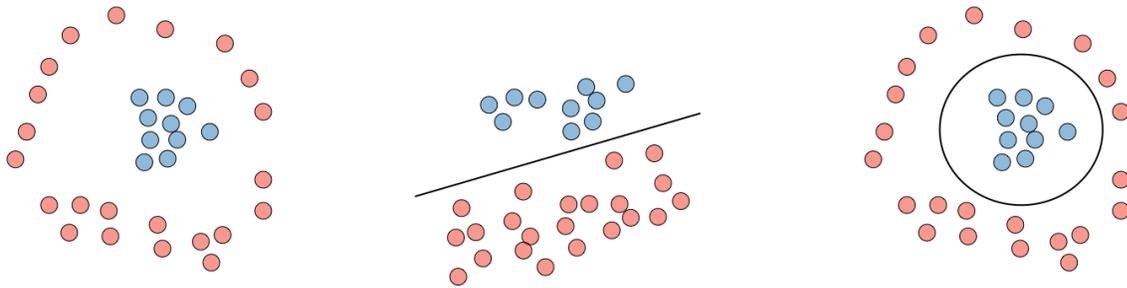
**Hinge loss** — The hinge loss is used in the setting of SVMs and is defined as follows:

$$L(z, y) = \max(0, 1 - yz)$$

**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  $K(x, z) = \exp(-\frac{1}{2\sigma^2} \|x - z\|^2)$  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  $L(w, b)$  as follows:

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

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

## 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.

### Gaussian Discriminant Analysis

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

$$(1) y \sim \text{Bernoulli}(\phi)$$

$$(2) x|y=0 \sim N(\mu_0, \Sigma) \quad x|y=1 \sim N(\mu_1, \Sigma)$$

$$(3) x|y=1 \sim N(\mu_1, \Sigma) \quad x|y=0 \sim N(\mu_0, \Sigma)$$

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

$\hat{\phi}$	$\hat{\mu}_j (j=0, 1)$	$\hat{\Sigma}$
$\frac{1}{m} \sum_{i=1}^m \{y(i)=1\}$	$\frac{1}{m} \sum_{i=1}^m \{y(i)=j\} x(i)$	$\frac{1}{m} \sum_{i=1}^m (x(i) - \mu_{y(i)})(x(i) - \mu_{y(i)})^T$

### Naive Bayes

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

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

**Solutions** — Maximizing the log-likelihood gives the following solutions, with  $k \in \{0, 1\}$ ,  $l \in \{1, L\}$

$P(y=k) = \frac{1}{m} \sum_{j=1}^m \mathbb{1}_{\{y(j)=k\}}$  and  $P(x_i=l | y=k) = \frac{\sum_{j=1}^m \mathbb{1}_{\{y(j)=k \text{ and } x_i(j)=l\}}}{\sum_{j=1}^m \mathbb{1}_{\{y(j)=k\}}}$

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

## 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.

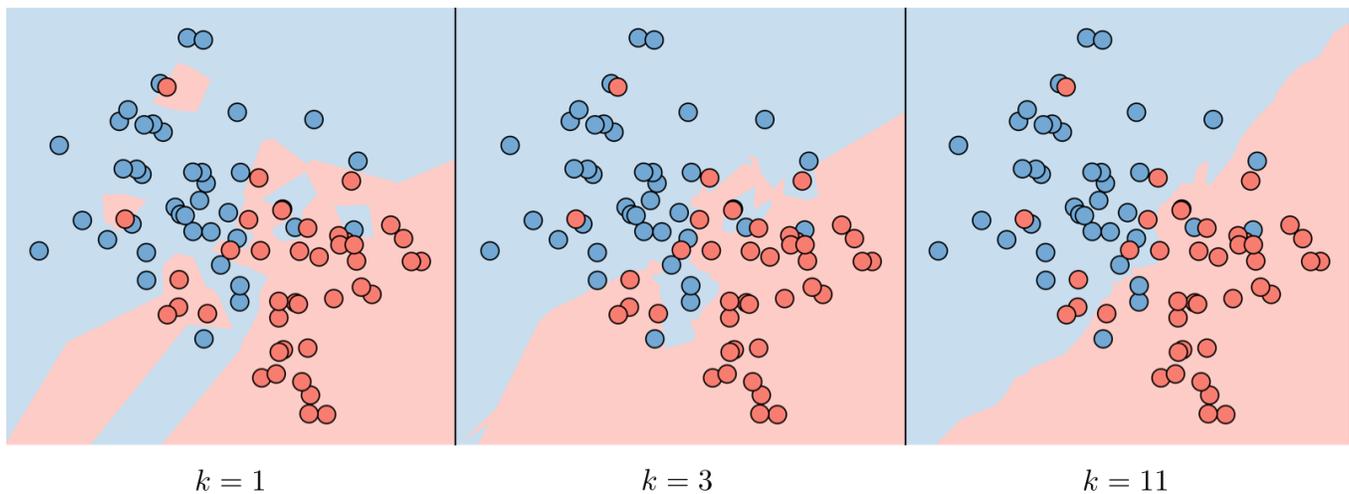
**Boosting** — The idea of boosting methods is to combine several weak learners to form a stronger one. The main ones are summed up in the table below:

Adaptive boosting	Gradient boosting
<ul style="list-style-type: none"> <li>• Known as Adaboost</li> <li>• High weights are put on errors to improve at the next boosting step</li> </ul>	<ul style="list-style-type: none"> <li>• Weak learners trained on remaining errors</li> </ul>

## 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.

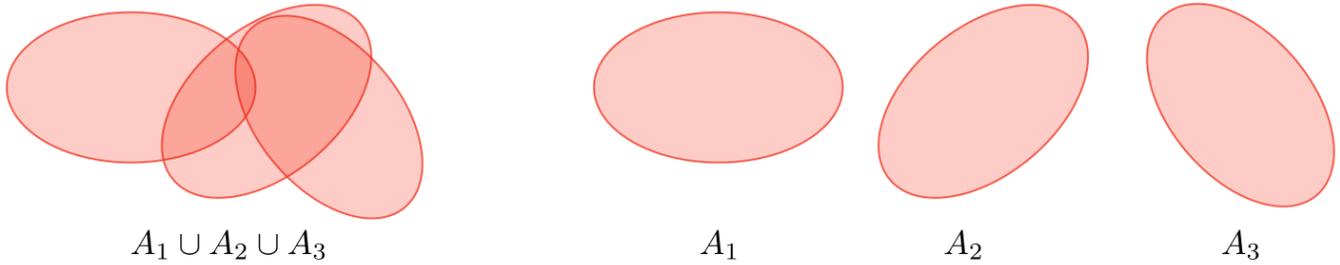
Remark: The higher the parameter k, the higher the bias, and the lower the parameter k, the higher the variance.



## Learning Theory

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

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



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

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

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

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

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

**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  $H$ , we say that  $H$  shatters  $S$  if for any set of labels  $\{y(1), \dots, y(d)\}$ , we have:

$$\exists h \in H, \forall i \in [1, d], h(x(i)) = y(i)$$

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

$$\hat{\epsilon}(h) \leq \min_{h \in H} \epsilon(h) + 2\sqrt{\frac{2 \log(2k\delta)}{m}}$$

**VC dimension** — The Vapnik-Chervonenkis (VC) dimension of a given infinite hypothesis class  $H$ , noted  $VC(H)$  is the size of the largest set that is shattered by  $H$ .

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



**Theorem (Vapnik)** — Let  $H$  be given, with  $VC(H) = d$  and  $m$  the number of training examples. With probability at least  $1 - \delta$ , we have:

$$\hat{\epsilon}(h) \leq \min_{h \in H} \epsilon(h) + O\left(\sqrt{\frac{d \log(md)}{m}} + 1 \log(1/\delta)\right)$$