

Would you like to see this cheatsheet in your native language? You can help us translating it (<https://github.com/shervinea/cheatsheet-translation>) on GitHub!

(<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#cs-229---machine-learning>) **CS 229 - Machine Learning (teaching/cs-229)**

العربية (/ar/teaching/cs-229/cheatsheet-supervised-learning) **English** Español (/es/teaching/cs-229/hoja-referencia-aprendizaje-supervisado) فارسی (/fa/teaching/cs-229/cheatsheet-supervised-learning) Français (/fr/teaching/cs-229/pense-bete-apprentissage-supervise) 한국어 (/ko/teaching/cs-229/cheatsheet-supervised-learning) Português (/pt/teaching/cs-229/dicas-aprendizado-supervisionado) Türkçe (/tr/teaching/cs-229/cheatsheet-supervised-learning) 簡中 (/zh/teaching/cs-229/cheatsheet-supervised-learning) 繁中 (/zh-tw/teaching/cs-229/cheatsheet-supervised-learning)

(<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#cheatsheet>) Supervised Learning cheatsheet

★ Star 10,118

By Afshine Amidi (<https://twitter.com/afshinea>) and Shervine Amidi (<https://twitter.com/shervinea>)

(<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#introduction>) Introduction to Supervised Learning

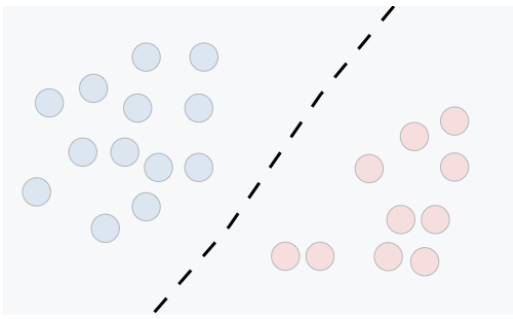
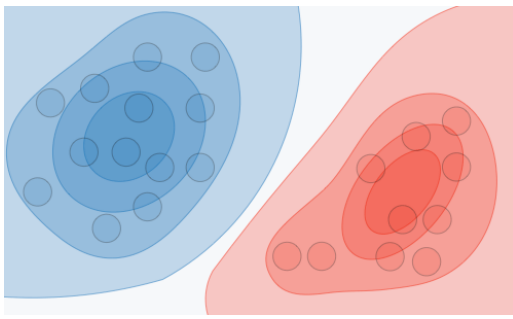
Given a set of data points $\{x^{(1)}, \dots, x^{(m)}\}$ associated to a set of outcomes $\{y^{(1)}, \dots, y^{(m)}\}$, we want to build a classifier that learns how to predict y from x .

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

	Regression	Classifier
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 then deduce $P(y x)$
What's learned	Decision boundary	Probability distributions of the data

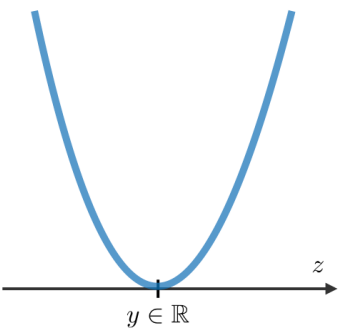
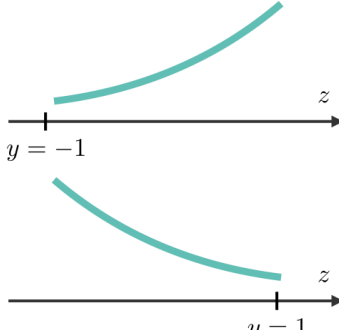
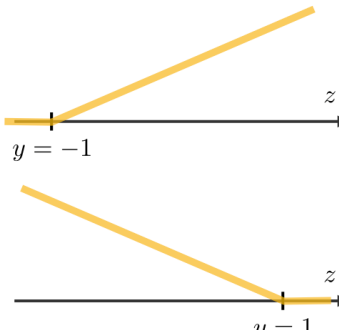
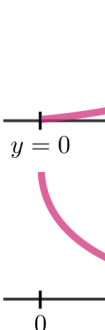
Illustration		
Examples	Regressions, SVMs	GDA, Naive Bayes

(<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#notations>)

Notations and general concepts

Hypothesis — The hypothesis is noted h_θ 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 $L : (z, y) \in \mathbb{R} \times Y \mapsto L(z, y) \in \mathbb{R}$ 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:

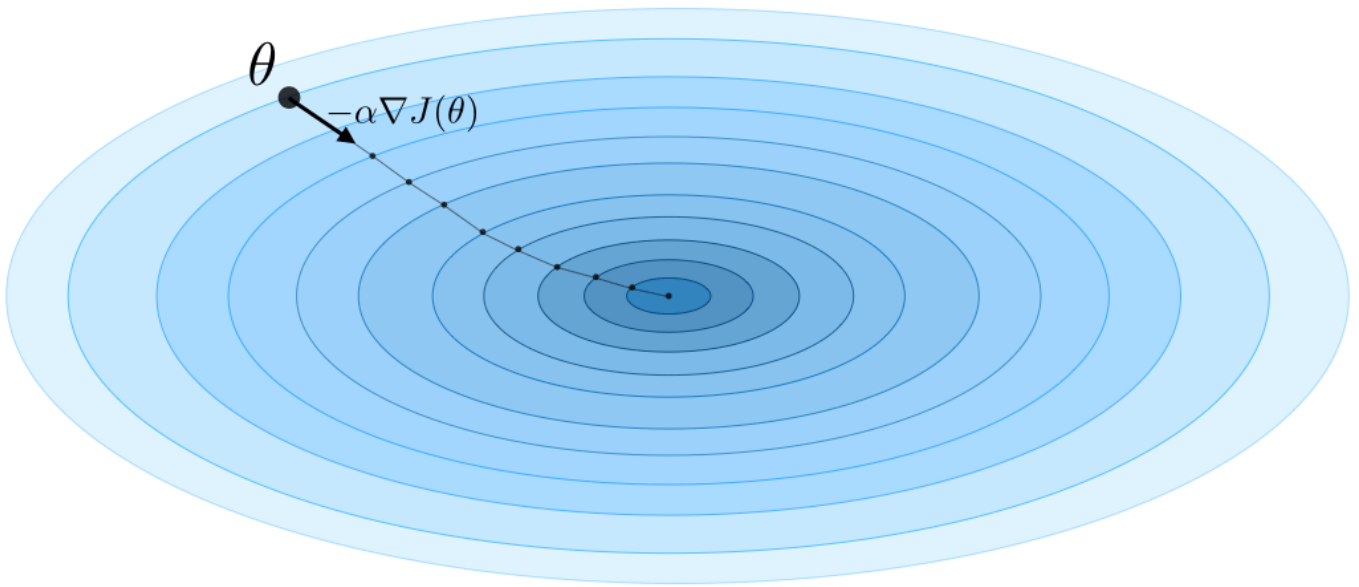
Least squared error	Logistic loss	Hinge loss	Cross entropy
$\frac{1}{2}(y - z)^2$	$\log(1 + \exp(-yz))$	$\max(0, 1 - yz)$	$-[y \log(z) + (1 - y) \log(1 - z)]$
			
Linear regression	Logistic regression	SVM	Neural networks

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:

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



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 θ is used to find the optimal parameters θ 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 θ 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:

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

(<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#linear-models>)

Linear models

Linear regression

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

Normal equations — By noting X the design matrix, the value of θ that minimizes the cost function is a closed-form solution such that:

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

LMS algorithm — By noting α 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:

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

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 ϕ_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, η , 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	η	$T(y)$	$a(\eta)$	$b(y)$
Bernoulli	$\log\left(\frac{\phi}{1-\phi}\right)$	y	$\log(1 + \exp(\eta))$	1

Gaussian	μ	y	$\frac{\eta^2}{2}$	$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$
Poisson	$\log(\lambda)$	y	e^η	$\frac{1}{y!}$
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 of $x \in \mathbb{R}^{n+1}$ and rely on the following 3 assumptions:

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

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

<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#svm> Support Vector Machines

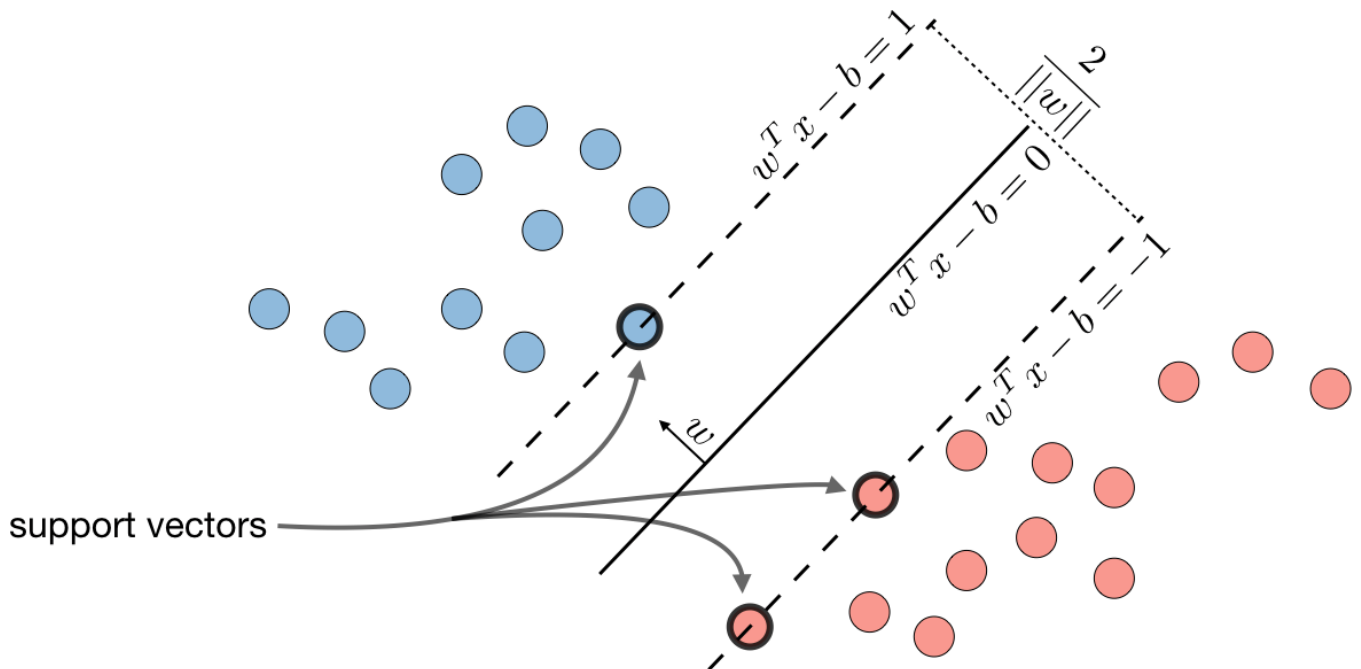
The goal of support vector machines is to find the line that maximizes the minimum distance to the line.

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:

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



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

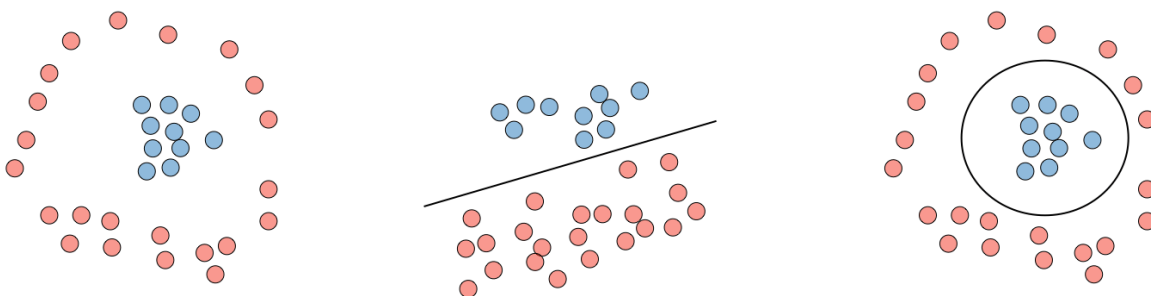
Hinge loss — The hinge loss is used in the setting of SVMs and is defined as follows:

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

Kernel — Given a feature mapping ϕ , 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 ϕ \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 ϕ , 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 β_i are called the Lagrange multipliers.

<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#generative-learning> 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 y and $x|y = 0$ and $x|y = 1$ are such that:

- (1) $y \sim \text{Bernoulli}(\phi)$
- (2) $x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma)$
- (3) $x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$

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

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

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 $k \in \{0, 1\}, l \in \llbracket 1, L \rrbracket$

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

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

<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#tree>

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:

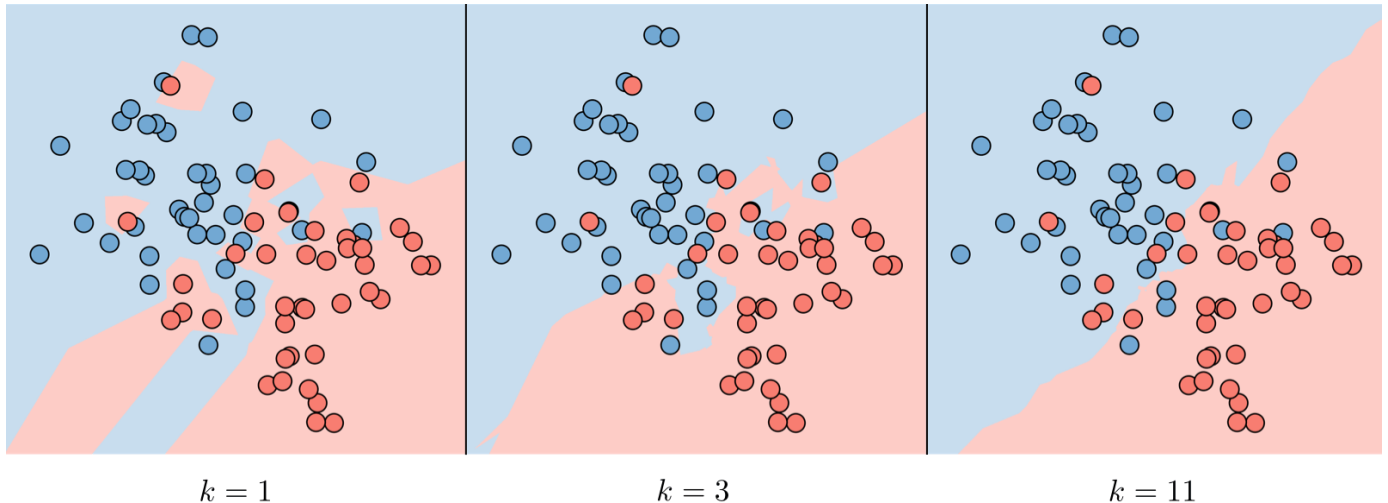
Adaptive boosting	Gradient boosting
<ul style="list-style-type: none"> Known as Adaboost High weights are put on errors to improve at the next boosting step 	<ul style="list-style-type: none"> Weak learners trained on remaining error

(<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#other>)

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.

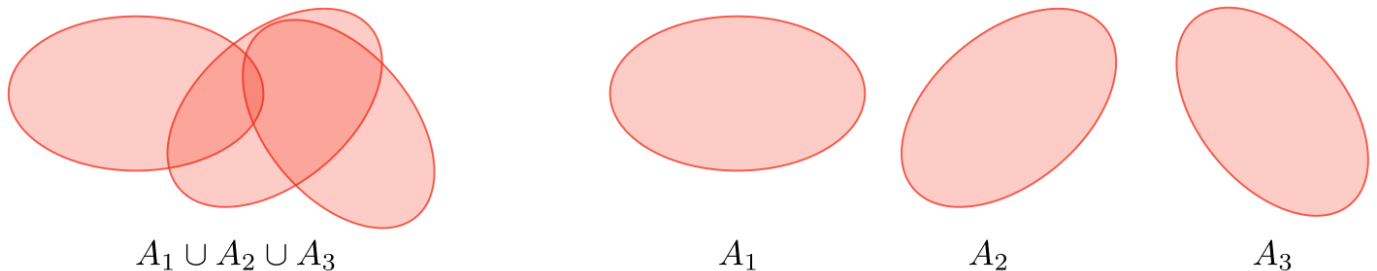


(<https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning#learning-theory>)

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 ϕ . Let $\hat{\phi}$ be their sample mean and $\gamma > 0$ fixed. We have:

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

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:

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

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:

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

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

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

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 $\mathcal{H} = \{\text{set of linear classifiers in 2 dimensions}\}$ is 3.



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

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



(<https://twitter.com/shervinea>)



(<https://linkedin.com/in/shervineamidi>)



(<https://github.com/shervinea>)



(<https://scholar.google.com/citations?user=nMnMTm8AAAAJ>)

