Would you like to see this cheatsheet in your native language? You can help us translating it on GitHub!

English

☆ Star 11,859

Tips and tricks

Introduction Type of prediction Type of model

Loss function Gradient descent Likelihood Linear models ↓ Linear regression

Notations and general concepts

 ↓ Logisitic regression Generalized linear models

Support Vector Machines

Optimal margin classifier Hinge loss

Kernel

Generative learning ♦ Naive Bayes

Trees and ensemble methods

Gaussian Discriminant Analysis

CART Random forest Boosting Other methods ♦ k-NN Learning Theory

Hoeffding inequality

VC dimension

PAC

Supervised Learning Unsupervised Learning **Supervised Learning cheatsheet**

By Afshine Amidi and Shervine Amidi

CS 229 - Machine Learning

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

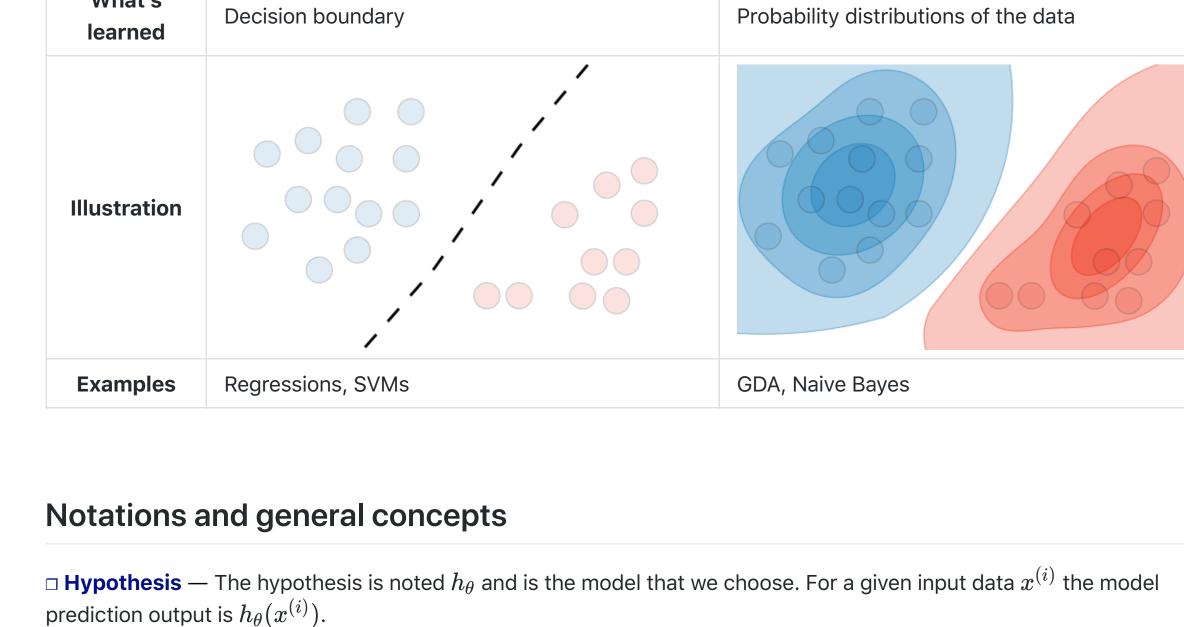
that learns how to predict y from x. □ **Type of prediction** — The different types of predictive models are summed up in the table below: Classification Regression

Deep Learning

Class Continuous Outcome Linear regression Logistic regression, SVM, Naive Bayes

□ **Type of model** — The different models are summed up in the table below: **Generative model Discriminative model** Directly estimate P(y|x)Estimate P(x|y) to then deduce P(y|x)Goal

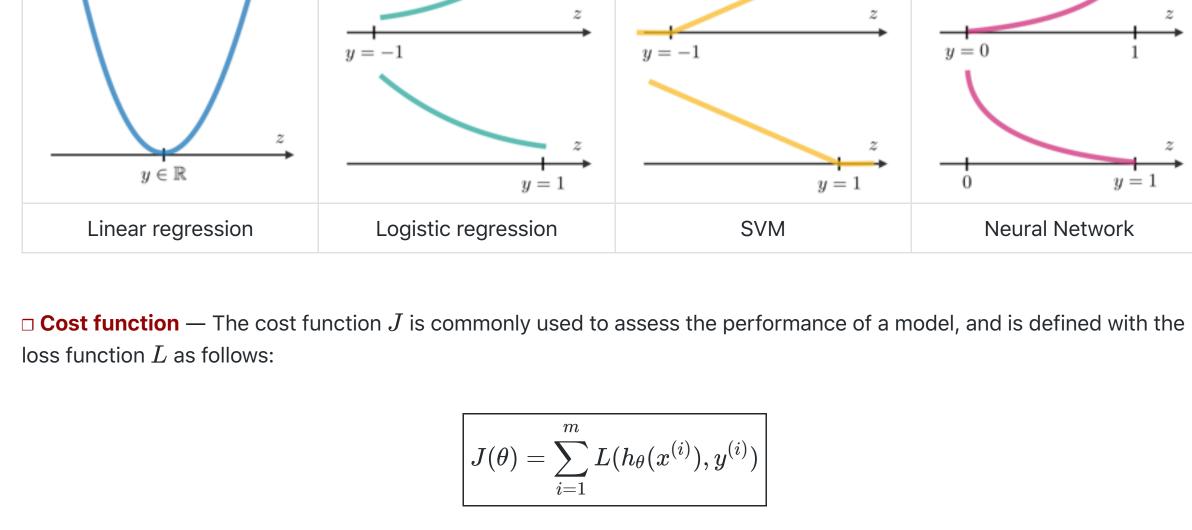
What's **Decision boundary** learned



$exttt{$\square$}$ Loss function — A loss function is a function $L:(z,y)\in\mathbb{R} imes Y\longmapsto L(z,y)\in\mathbb{R}$ 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 error **Logistic loss** Hinge loss **Cross-entropy**

 $- \Big| y \log(z) + (1-y) \log(1 \frac{1}{2}(y-z)^2$ $\max(0, 1 - yz)$ $\log(1+\exp(-yz))$



 \square **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:

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

 \Box **Likelihood** — The likelihood of a model $L(\theta)$ given parameters θ is used to find the optimal parameters θ through

 $ig| heta^{ ext{opt}} = rg \max L(heta)$

 \square **Newton's algorithm** — Newton's algorithm is a numerical method that finds heta such that $\ell'(heta)=0$. Its update rule is as follows:

Linear regression We assume here that $y|x; heta \sim \mathcal{N}(\mu,\sigma^2)$ \Box Normal equations — By noting X the design matrix, the value of θ that minimizes the cost function is a closed-form

 $\left| heta = (X^T X)^{-1} X^T y
ight|$

 \Box LMS algorithm — By noting α the learning rate, the update rule of the Least Mean Squares (LMS) algorithm for a

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 $au \in \mathbb{R}$ as:

 $\left|w^{(i)}(x)=\exp\left(-rac{(x^{(i)}-x)^2}{2 au^2}
ight)
ight|$

 $oxed{ orall j, \quad heta_j \leftarrow heta_j + lpha \sum_{i=1}^m \left[y^{(i)} - h_ heta(x^{(i)})
ight] x_j^{(i)} }$

□ 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 $heta_K=0$, which makes the Bernoulli parameter ϕ_i of each class i be such that:

 $\phi = \overline{p(y=1|x; heta)} = rac{1}{1+\exp(- heta^Tx)} = g(heta^Tx)$

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

The most common exponential distributions are summed up in the following table: T(y) $a(\eta)$

 $\log(1+\exp(\eta))$

b(y)

 $\frac{1}{\sqrt{2\pi}}\exp\left(-\frac{y^2}{2}\right)$ Gaussian $\log(\lambda)$ Poisson

 \mathbb{R}^{n+1} and rely on the following 3 assumptions: $ig|y|x; heta \sim ext{ExpFamily}(\eta)ig|$ $ig| h_{ heta}(x) = E[y|x; heta]$ Remark: ordinary least squares and logistic regression are special cases of generalized linear models.

 $\log\left(\frac{\phi}{1-\phi}\right)$

Support Vector Machines The goal of support vector machines is to find the line that maximizes the minimum distance to the line. \Box Optimal margin classifier — The optimal margin classifier h is such that: $h(x) = \operatorname{sign}(w^T x - b)$ where $(w,b)\in\mathbb{R}^n imes\mathbb{R}$ is the solution of the following optimization problem:

 $L(z,y) = [1-yz]_+ = \max(0,1-yz)$ \square Kernel — Given a feature mapping ϕ , we define the kernel K as follows: $igg|K(x,z)=\phi(x)^T\phi(z)igg|$ In practice, the kernel K defined by $K(x,z)=\exp\left(-rac{||x-z||^2}{2\sigma^2}
ight)$ is called the Gaussian kernel and is commonly used.

 \square **Setting** — The Gaussian Discriminant Analysis assumes that y and x|y=0 and x|y=1 are such that: $|x|y=0 \sim \mathcal{N}(\mu_0,\Sigma)$ (3) $|x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$ $y \sim \mathrm{Bernoulli}(\phi)$ □ **Estimation** — The following table sums up the estimates that we find when maximizing the likelihood: $\widehat{\mu_j} \quad (j=0,1)$

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

□ **Solutions** — Maximizing the log-likelihood gives the following solutions:

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

Other non-parametric approaches \square 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 = 3

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

 $A_1 \cup A_2 \cup A_3$ A_3 A_1 A_2 \Box Hoeffding inequality — Let $Z_1,..,Z_m$ be m iid variables drawn from a Bernoulli distribution of parameter ϕ . Let $\widehat{\phi}$ be their sample mean and $\gamma>0$ fixed. We have: $\left|P(|\phi-\widehat{\phi}|>\gamma)\leqslant 2\exp(-2\gamma^2m)
ight|$ Remark: this inequality is also known as the Chernoff bound. \Box **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:

• the training examples are drawn independently \Box Shattering — Given a set $S=\{x^{(1)},...,x^{(d)}\}$, and a set of classifiers $\mathcal H$, we say that $\mathcal H$ shatters S if for any set of

 $\left| \epsilon(\widehat{h}) \leqslant \left(\min_{h \in \mathcal{H}} \epsilon(h)
ight) + O\left(\sqrt{rac{d}{m}} \log\left(rac{m}{d}
ight) + rac{1}{m} \log\left(rac{1}{\delta}
ight)
ight)
ight|$

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 maximization. We have:

Remark: in practice, we use the log-likelihood $\ell(\theta) = \log(L(\theta))$ which is easier to optimize.

Remark: the multidimensional generalization, also known as the Newton-Raphson method, has the following update rule: $heta \leftarrow heta - \left(
abla_{ heta}^2 \ell(heta)
ight)^{-1}
abla_{ heta} \ell(heta)$ **Linear models**

training set of m data points, which is also known as the Widrow-Hoff learning rule, is as follows:

Classification and logistic regression

Remark: logistic regressions do not have closed form solutions.

Distribution

Bernoulli

Geometric

support vectors

Remark: the decision boundary is defined as $\left| w^T x - b = 0 \right|$

Remark: the coefficients β_i are called the Lagrange multipliers.

Generative Learning

P(y|x) by using Bayes' rule.

with $k \in \{0,1\}$ and $l \in \llbracket 1,L
rbracket$

are summed up in the table below:

boosting step

Known as Adaboost

Learning Theory

Adaptive boosting

• High weights are put on errors to improve at the next

k = 1

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

Gaussian Discriminant Analysis

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

function $a(\eta)$ as follows:

solution such that:

 \square Sigmoid function — The sigmoid function g, also known as the logistic function, is defined as follows: $orall z \in \mathbb{R}, \quad \left| g(z) = rac{1}{1+e^{-z}} \in]0,1[
ight|$

 \Box Logistic regression — We assume here that $y|x; \theta \sim \mathrm{Bernoulli}(\phi)$. We have the following form:

$$\sum_{j=1}^K \exp(heta_j^T x)$$

 $\left| p(y;\eta) = b(y) \exp(\eta T(y) - a(\eta)) \right|$ 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.

 $\log\left(\frac{e^{\eta}}{1-e^{\eta}}\right)$ $\log(1-\phi)$ $lue{}$ Assumptions of GLMs — Generalized Linear Models (GLM) aim at predicting a random variable y as a function of $x \in \mathbb{R}$

> $\Big|\minrac{1}{2}||w||^2$ $y^{(i)}(w^Tx^{(i)}-b)\geqslant 1$ such that

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. $oxedsymbol{\square}$ Lagrangian — We define the Lagrangian $\mathcal{L}(w,b)$ as follows:

 $oxedsymbol{\mathcal{L}}(w,b) = f(w) + \sum_{i=1}^{n} eta_i h_i(w)$

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

 $rac{1}{m} \sum_{i=1}^m 1_{\{y^{(i)}=1\}} \quad rac{\sum_{i=1}^m 1_{\{y^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m 1_{\{y^{(i)}=j\}}} \quad rac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T$ **Naive Bayes** □ **Assumption** — The Naive Bayes model supposes that the features of each data point are all independent: $oxed{P(x|y) = P(x_1, x_2, ...|y) = P(x_1|y)P(x_2|y)... = \prod_{i=1}^n P(x_i|y)}$

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

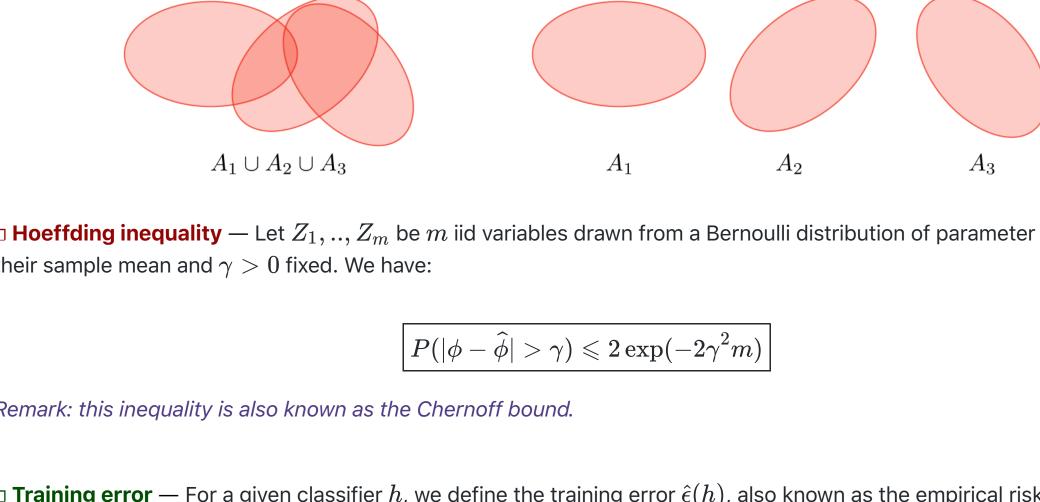
Gradient boosting

k = 11

• Weak learners are trained on residuals

• Examples include XGBoost

 $oxed{P(y=k) = rac{1}{m} imes \#\{j|y^{(j)} = k\}} \quad ext{and} \quad egin{array}{c} P(x_i = l|y = k) = rac{\#\{j|y^{(j)} = k ext{ and } x_i^{(j)} = l\}}{\#\{j|y^{(j)} = k\}} \end{array}$



labels $\{y^{(1)},...,y^{(d)}\}$, we have:

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

size of the largest set that is shattered by \mathcal{H} . least $1-\delta$, we have:

• View PDF version on GitHub

□ 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

 \Box **VC dimension** — The Vapnik-Chervonenkis (VC) dimension of a given infinite hypothesis class \mathcal{H} , noted $VC(\mathcal{H})$ is the

 $\exists h \in \mathcal{H}, \quad orall i \in \llbracket 1, d
rbracket, \quad h(x^{(i)}) = y^{(i)}$ lacktriangle Upper bound theorem — Let ${\cal H}$ be a finite hypothesis class such that $|{\cal H}|=k$ and let δ and the sample size m be fixed. Then, with probability of at least $1-\delta$, we have:

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

 \square **Theorem (Vapnik)** — Let $\mathcal H$ be given, with $\mathrm{VC}(\mathcal H)=d$ and m the number of training examples. With probability at

y in G S