

<https://github.com/soulmachine/machine-learning-cheat-sheet>
soulmachine@gmail.com

المرجع السريع في علم تعلم الآلة

2 جويلية 2017

©2013 soulmachine

Except where otherwise noted, This document is licensed under a Creative Commons Attribution-ShareAlike 3.0 Unported (CC BY-SA3.0) license (<http://creativecommons.org/licenses/by/3.0/>).

المحتويات

1.3.2	A	3 المحتويات
	sim- ple		Notation 7
	non- parametric clas-	1	تمهيد 1
	si- fier:		1.1 أنواع تعلم الآلة
	K- nearest	1	Types of machine learning
	neigh- bours	2	1.2 المكونات الثلاثة
1.3.3	Overfitting	2	نماذج تعلم الآلة
1.3.4	Cross val-	1	1.1.2 التمثيل
	i- da- tion	2	Representation 1.2.2 Evaluation 1
			1.2.3 Optimization 2
1.3.5	Model	1.3	Some basic con-
	se- lec- tion	2	cepts 2
			1.3.1 Parametric vs non-parametric models 2

List of Contributors

Wei Zhang

PhD candidate at the Institute of Software, Chinese Academy of Sciences (ISCAS), Beijing, P.R.CHINA, e-mail: zh3feng@gmail.com, has written chapters of Naive Bayes and SVM.

Fei Pan

Master at Beijing University of Technology, Beijing, P.R.CHINA, e-mail: example@gmail.com, has written chapters of KMeans, AdaBoost.

Yong Li

PhD candidate at the Institute of Automation of the Chinese Academy of Sciences (CASIA), Beijing, P.R.CHINA, e-mail: liyong3forever@gmail.com, has written chapters of Logistic Regression.

Jiankou Li

PhD candidate at the Institute of Software, Chinese Academy of Sciences (ISCAS), Beijing, P.R.CHINA, e-mail: lijiankoucoco@163.com, has written chapters of BayesNet.

Notation

مجموع الرموز

Introduction

تمهيد

من الصعوبة التوصل إلى مجموعة وحيدة وثابتة من الرموز لتغطية المجال الشاسع من البيانات (data) و النماذج (models) والخوارزميات (algorithms) التي نناقشها في هذا الكتيب. علاوة على ذلك، العلامات الرياضية المتفق عليها تختلف بين علم تعلم الآلة (machine learning) و علم الإحصاءات (statistics)، و بين الكتب والأوراق العلمية المختلفة. مع ذلك، فقد حاولنا أن تكون الرموز المستعملة متسقة قدر الإمكان. فيما يلي نلخص معظم الرموز المستخدمة، هذا لا يعني أن بعض المقاطع الفردية في الكتيب قد تعرض رموزا جديدة. إعلم أيضا أن بعض الرموز قد يكون لها معان مختلفة تبعا للسياق، رغم أننا سنحرص على تجنب ذلك قدر الإمكان.

General math notation

مجموع الرموز الرياضية

Symbol	Meaning
$\lfloor x \rfloor$	Floor of x , i.e., round down to nearest integer
$\lceil x \rceil$	Ceiling of x , i.e., round up to nearest integer
$\mathbf{x} \otimes \mathbf{y}$	Convolution of \mathbf{x} and \mathbf{y}
$\mathbf{x} \odot \mathbf{y}$	Hadamard (elementwise) product of \mathbf{x} and \mathbf{y}
$a \wedge b$	logical AND
$a \vee b$	logical OR
$\neg a$	logical NOT
$\mathbb{I}(x)$	Indicator function, $\mathbb{I}(x) = 1$ if x is true, else $\mathbb{I}(x) = 0$
∞	Infinity
\rightarrow	Tends towards, e.g., $n \rightarrow \infty$
\propto	Proportional to, so $y = ax$ can be written as $y \propto x$
$ x $	Absolute value
$ \mathcal{S} $	Size (cardinality) of a set
$n!$	Factorial function
∇	Vector of first derivatives
∇^2	Hessian matrix of second derivatives
\triangleq	Defined as
$O(\cdot)$	Big-O: roughly means order of magnitude
\mathbb{R}	The real numbers
$1:n$	Range (Matlab convention): $1:n = 1, 2, \dots, n$
\approx	Approximately equal to
$\arg \max_x f(x)$	Argmax: the value x that maximizes f
$B(a, b)$	Beta function, $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$
$B(\alpha)$	Multivariate beta function, $\frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)}$
$\binom{n}{k}$	n choose k , equal to $n!/(k!(n-k)!)$
$\delta(x)$	Dirac delta function, $\delta(x) = \infty$ if $x = 0$, else $\delta(x) = 0$

$\exp(x)$	Exponential function e^x
$\Gamma(x)$	Gamma function, $\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$
$\Psi(x)$	Digamma function, $\Psi(x) = \frac{d}{dx} \log \Gamma(x)$
\mathcal{X}	A set from which values are drawn (e.g., $\mathcal{X} = \mathbb{R}^D$)

Linear algebra notation

رموز علم الجبر الخطي

خلال هذا الكتاب، سنستخدم أحرف **boldface** الصغيرة للدلالة على الناقلات (vectors) مثل x و أحرف **boldface** العلوية للدلالة على المصفوفات مثل X . نشير إلى إدخالات مصفوفة (matrix entries) بأحرف كبيرة غير جريئة، مثل X_{ij} . سنعتبر كل الناقلات (vectors) ناقلات عمود (column vectors)، ما لم يذكر خلاف ذلك في السياق. نستخدم (x_1, \dots, x_D) للدلالة على متجه عمود تم إنشاؤه بواسطة تكديس D أعداد (stacking D scalars). إذا كتبنا $X = (x_1, \dots, x_n)$ ، حيث الجانب الأيسر هو مصفوفة (matrix)، فإننا نقصد تكديس x_i على طول الأعمدة لخلق مصفوفة (matrix).

Symbol	Meaning
$X \succ 0$	X is a positive definite matrix
$tr(X)$	Trace of a matrix
$det(X)$	Determinant of matrix X
$ X $	Determinant of matrix X
X^{-1}	Inverse of a matrix
X^\dagger	Pseudo-inverse of a matrix
X^T	Transpose of a matrix
x^T	Transpose of a vector
$diag(x)$	Diagonal matrix made from vector x
$diag(X)$	Diagonal vector extracted from matrix X
I or I_d	Identity matrix of size $d \times d$ (ones on diagonal, zeros of)
$\mathbf{1}$ or $\mathbf{1}_d$	Vector of ones (of length d)
$\mathbf{0}$ or $\mathbf{0}_d$	Vector of zeros (of length d)
$\ x\ = \ x\ _2$	Euclidean or ℓ_2 norm $\sqrt{\sum_{j=1}^d x_j^2}$
$\ x\ _1$	ℓ_1 norm $\sum_{j=1}^d x_j $
$X_{:,j}$	j 'th column of matrix
$X_{i,:}$	transpose of i 'th row of matrix (a column vector)
$X_{i,j}$	Element (i, j) of matrix X
$x \otimes y$	Tensor product of x and y

Probability notation

رموز علم الإحتمال

نرمز إلى الأعداد العشوائية و الثابتة (random and fixed scalars) بخط صغير (lower case)، و الناقلات العشوائية و الثابتة (random and fixed vectors) بأحرف الصغيرة الجريئة (bold lower case) و المصفوفات العشوائية و الثابتة (fixed matrices) بأحرف الجريئة العلوية (bold upper case). أحيانا نستخدم الأحرف العلوية غير الجريئة (non-bold upper case) للدلالة على المتغيرات العددية العشوائية (scalar random variables). نستخدم، أيضا، $p()$ لكل من المتغيرات العشوائية المنفصلة و المستمرة (discrete and continuous random variables).

Symbol	Meaning
X, Y	Random variable
$P()$	Probability of a random event
$F()$	Cumulative distribution function(CDF), also called distribution function
$p(x)$	Probability mass function(PMF)
$f(x)$	probability density function(PDF)
$F(x, y)$	Joint CDF
$p(x, y)$	Joint PMF
$f(x, y)$	Joint PDF
$p(X Y)$	Conditional PMF, also called conditional probability
$f_{X Y}(x y)$	Conditional PDF
$X \perp Y$	X is independent of Y
$X \not\perp Y$	X is not independent of Y
$X \perp Y Z$	X is conditionally independent of Y given Z
$X \not\perp Y Z$	X is not conditionally independent of Y given Z
$X \sim p$	X is distributed according to distribution p
α	Parameters of a Beta or Dirichlet distribution
$\text{cov}[X]$	Covariance of X
$\mathbb{E}[X]$	Expected value of X
$\mathbb{E}_q[X]$	Expected value of X wrt distribution q
$\mathbb{H}(X)$ or $\mathbb{H}(p)$	Entropy of distribution $p(X)$
$\mathbb{I}(X; Y)$	Mutual information between X and Y
$\mathbb{KL}(p q)$	KL divergence from distribution p to q
$\ell(\theta)$	Log-likelihood function
$L(\theta, a)$	Loss function for taking action a when true state of nature is θ
λ	Precision (inverse variance) $\lambda = 1/\sigma^2$
Λ	Precision matrix $\Lambda = \Sigma^{-1}$
$\text{mode}[\mathbf{X}]$	Most probable value of \mathbf{X}
μ	Mean of a scalar distribution
$\boldsymbol{\mu}$	Mean of a multivariate distribution
Φ	cdf of standard normal
ϕ	pdf of standard normal
π	multinomial parameter vector, Stationary distribution of Markov chain
ρ	Correlation coefficient
$\text{sigm}(x)$	Sigmoid (logistic) function, $\frac{1}{1+e^{-x}}$
σ^2	Variance
Σ	Covariance matrix
$\text{var}[x]$	Variance of x
ν	Degrees of freedom parameter
Z	Normalization constant of a probability distribution

Machine learning/statistics notation

رموز علم تعلم الآلة والإحصاءات

In general, we use upper case letters to denote constants, such as C, K, M, N, T , etc. We use lower case letters as dummy indexes of the appropriate range, such as $c = 1 : C$ to index classes, $i = 1 : M$ to index data cases, $j = 1 : N$ to index input features, $k = 1 : K$ to index states or clusters, $t = 1 : T$ to index time, etc.

We use x to represent an observed data vector. In a supervised problem, we use y or \mathbf{y} to represent the desired output label. We use \mathbf{z} to represent a hidden variable. Sometimes we also use q to represent a hidden discrete variable.

Symbol	Meaning
C	Number of classes
D	Dimensionality of data vector (number of features)
N	Number of data cases
N_c	Number of examples of class c , $N_c = \sum_{i=1}^N \mathbb{I}(y_i = c)$
R	Number of outputs (response variables)
\mathcal{D}	Training data $\mathcal{D} = \{(\mathbf{x}_i, y_i) i = 1 : N\}$
\mathcal{D}_{test}	Test data
\mathcal{X}	Input space
\mathcal{Y}	Output space
K	Number of states or dimensions of a variable (often latent)
$k(x, y)$	Kernel function
\mathbf{K}	Kernel matrix
\mathcal{H}	Hypothesis space
L	Loss function
$J(\boldsymbol{\theta})$	Cost function
$f(\mathbf{x})$	Decision function
$P(y \mathbf{x})$	Conditional probability
λ	Strength of ℓ_2 or ℓ_1 <i>regularizer</i>
$\phi(x)$	Basis function expansion of feature vector \mathbf{x}
Φ	Basis function expansion of design matrix \mathbf{X}
$q()$	Approximate or proposal distribution
$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{old})$	Auxiliary function in EM
T	Length of a sequence
$T(\mathcal{D})$	Test statistic for data
\mathbf{T}	Transition matrix of Markov chain
$\boldsymbol{\theta}$	Parameter vector
$\boldsymbol{\theta}^{(s)}$	s 'th sample of parameter vector
$\hat{\boldsymbol{\theta}}$	Estimate (usually MLE or MAP) of $\boldsymbol{\theta}$
$\hat{\boldsymbol{\theta}}_{MLE}$	Maximum likelihood estimate of $\boldsymbol{\theta}$
$\hat{\boldsymbol{\theta}}_{MAP}$	MAP estimate of $\boldsymbol{\theta}$
$\bar{\boldsymbol{\theta}}$	Estimate (usually posterior mean) of $\boldsymbol{\theta}$
\mathbf{w}	Vector of regression weights (called $\boldsymbol{\beta}$ in statistics)
b	intercept (called $\boldsymbol{\varepsilon}$ in statistics)
\mathbf{W}	Matrix of regression weights
x_{ij}	Component (i.e., feature) j of data case i , for $i = 1 : N, j = 1 : D$
\mathbf{x}_i	Training case, $i = 1 : N$
\mathbf{X}	Design matrix of size $N \times D$
$\bar{\mathbf{x}}$	Empirical mean $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$
$\tilde{\mathbf{x}}$	Future test case
\mathbf{x}_*	Feature test case
\mathbf{y}	Vector of all training labels $\mathbf{y} = (y_1, \dots, y_N)$
z_{ij}	Latent component j for case i

Loss function and risk function

Definition 1.1. In order to measure how well a function fits the training data, a **loss function** $L: Y \times Y \rightarrow R \geq 0$ is defined. For training example (x_i, y_i) , the loss of predicting the value \hat{y} is $L(y_i, \hat{y})$.

The following is some common loss functions:

1. 0-1 loss function

$$L(Y, f(X)) = \mathbb{I}(Y, f(X)) = \begin{cases} 1, & Y \neq f(X) \\ 0, & Y = f(X) \end{cases}$$

2. Quadratic loss function $L(Y, f(X)) = (Y - f(X))^2$

3. Absolute loss function $L(Y, f(X)) = |Y - f(X)|$

4. Logarithmic loss function

$$L(Y, P(Y|X)) = -\log P(Y|X)$$

Definition 1.2. The risk of function f is defined as the expected loss of f :

$$R_{\text{exp}}(f) = E[L(Y, f(X))] = \int L(y, f(x)) P(x, y) dx dy$$

which is also called expected loss or **risk function**.

Definition 1.3. The risk function $R_{\text{exp}}(f)$ can be estimated from the training data as

$$R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

which is also called empirical loss or **empirical risk**.

You can define your own loss function, but if you're a novice, you're probably better off using one from the literature. There are conditions that loss functions should meet²:

1. They should approximate the actual loss you're trying to minimize. As was said in the other answer, the standard loss functions for classification is zero-one-loss (misclassification rate) and the ones used for training classifiers are approximations of that loss.
2. The loss function should work with your intended optimization algorithm. That's why zero-one-loss is not used directly: it doesn't work with gradient-based optimization methods since it doesn't have a well-defined gradient (or even a subgradient, like the hinge loss for SVMs has).

The main algorithm that optimizes the zero-one-loss directly is the old perceptron algorithm (chapter §??).

¹ Model = Representation + Evaluation + Optimization. Domingos, P. A few useful things to know about machine learning, Commun. ACM, 87-78:(10)55, (2012)

² <http://t.cn/zTrDxLO>

باب 1

Introduction

تمهيد

Types of machine learning أنواع تعلم الآلة

Supervised learning	<div>Classification</div> <div>Regression</div>
Unsupervised learning	<div>Discovering clusters</div> <div>Discovering latent factors</div> <div>Discovering graph structure</div> <div>Matrix completion</div>

المكونات الثلاثة لنماذج تعلم الآلة

النموذج (Model) = التمثيل (Representation) + التقييم (Evaluation) + التحسين (Optimization)¹
في ما يلي سيقع تفسير كل مكون من هذه المكونات الثلاثة على حدة.

Representation

التمثيل

In supervised learning, a model must be represented as a conditional probability distribution $P(y|x)$ (usually we call it classifier) or a decision function $f(x)$. The set of classifiers (or decision functions) is called the hypothesis space of the model. Choosing a representation for a model is tantamount to choosing the hypothesis space that it can possibly learn.

Evaluation

التقييم

In the hypothesis space, an evaluation function (also called objective function or risk function) is needed to distinguish good classifiers (or decision functions) from bad ones.

Evaluation

No training is needed.

Optimization

No training is needed.

Overfitting

Cross validation

Definition 1.7. **Cross validation**, sometimes called *rotation estimation*, is a *model validation* technique for assessing how the results of a statistical analysis will generalize to an independent data set³.

Common types of cross-validation:

1. K-fold cross-validation. In k-fold cross-validation, the original sample is randomly partitioned into k equal size subsamples. Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k - 1 subsamples are used as training data.
2. 2-fold cross-validation. Also, called simple cross-validation or holdout method. This is the simplest variation of k-fold cross-validation, k=2.
3. Leave-one-out cross-validation(LOOCV). k=M, the number of original samples.

Model selection

When we have a variety of models of different complexity (e.g., linear or logistic regression models with different degree polynomials, or KNN classifiers with different values of K), how should we pick the right one? A natural approach is to compute the misclassification rate on the training set for each method.

ERM and SRM

Definition 1.4. ERM(Empirical risk minimization)

$$\min_{f \in \mathcal{F}} R_{\text{emp}}(f) = \min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) \quad (1.3)$$

Definition 1.5. Structural risk

$$R_{\text{smp}}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f) \quad (1.4)$$

Definition 1.6. SRM(Structural risk minimization)

$$\min_{f \in \mathcal{F}} R_{\text{srm}}(f) = \min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f) \quad (1.5)$$

Optimization

التقييم

Finally, we need a training algorithm(also called learning algorithm) to search among the classifiers in the hypothesis space for the highest-scoring one. The choice of optimization technique is key to the efficiency of the model.

Some basic concepts

Parametric vs non-parametric models

A simple non-parametric classifier: K-nearest neighbours

Representation

$$y = f(\mathbf{x}) = \arg \min_c \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} \mathbb{I}(y_i = c) \quad (1.6)$$

where $N_k(\mathbf{x})$ is the set of k points that are closest to point \mathbf{x} .

Usually use k-d tree to accelerate the process of finding k nearest points.

³ [http://en.wikipedia.org/wiki/Cross-validation_\(statistics\)](http://en.wikipedia.org/wiki/Cross-validation_(statistics))