

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

المرجع السريع في علم تعلّم الآلة

30 جوان 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

| | | | |
|----------|----|-----|-----------|
| Notation | ix | iii | المحتويات |
|----------|----|-----|-----------|

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.

Acronyms

Use the template `acronym.tex` together with the Springer document class `SVMono` (monograph-type books) or `SVMult` (edited books) to style your list(s) of abbreviations or symbols in the Springer layout.

Lists of abbreviations, symbols and the like are easily formatted with the help of the Springer-enhanced description environment.

| | |
|------|---|
| ABC | Spelled-out abbreviation and definition |
| BABI | Spelled-out abbreviation and definition |
| CABR | Spelled-out abbreviation and definition |

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 ينفي 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

| 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

We use boldface lower-case to denote vectors, such as \mathbf{x} , and boldface upper-case to denote matrices, such as \mathbf{X} . We denote entries in a matrix by non-bold upper case letters, such as X_{ij} .

Vectors are assumed to be column vectors, unless noted otherwise. We use (x_1, \dots, x_D) to denote a column vector created by stacking D scalars. If we write $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, where the left hand side is a matrix, we mean to stack the \mathbf{x}_i along the columns, creating a matrix.

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

Probability notation

We denote random and fixed scalars by lower case, random and fixed vectors by bold lower case, and random and fixed matrices by bold upper case. Occasionally we use non-bold upper case to denote scalar random variables. Also, we use $p()$ for both 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 $\Lambda = \Sigma^{-1}$ |
| $\text{mode}[X]$ | Most probable value of 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

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 ε 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 |
