

Machine Learning Notation Template

Andreas Holzinger, a.holzinger@hci-kdd.org

Vienna University of Technology, Austria

Vienna, April 19, 2016



Abstract

The mathematical content of the course 185.A83 "Machine Learning for Health Informatics" is kept to a minimum, however, a nonzero level is necessary. It is hard to keep a consistent notation throughout the course, therefore in this short document the used mathematical notation is summarized.

1 Variables, Symbols, and Operations

An excellent compilation can be found in Duda, Hart and Stock (2000): Pattern Classification.

Introduction

It is very difficult to come up with a single, consistent notation to cover the wide variety of data, models and algorithms that we discuss. Furthermore, conventions differ between machine learning and statistics, and between different books and papers. Nevertheless, we have tried to be as consistent as possible. Below we summarize most of the notation used in this book, although individual sections may introduce new notation. Note also that the same symbol may have different meanings depending on the context, although we try to avoid this where possible.

General math notation

Symbol	Meaning
--------	---------

$\lfloor \mathbf{x} \rfloor$	Floor of \mathbf{x} , i.e., round down to nearest integer
$\lceil \mathbf{x} \rceil$	Ceiling of \mathbf{x} , i.e., round up to nearest integer
$\tilde{\mathbf{x}} \otimes \tilde{\mathbf{y}}$	Convolution of $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{y}}$
$\tilde{\mathbf{x}} \odot \tilde{\mathbf{y}}$	Hadamard (elementwise) product of $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{y}}$
$a \wedge b$	logical AND
$a \vee b$	logical OR
$\neg a$	logical NOT
$\mathbb{I}(\mathbf{x})$	Indicator function, $\mathbb{I}(\mathbf{x}) = 1$ if \mathbf{x} is true, else $\mathbb{I}(\mathbf{x}) = 0$
∞	Infinity
\rightarrow	Tends towards, e.g., $n \rightarrow \infty$
\propto	Proportional to, so $\mathbf{y} = \mathbf{a}\mathbf{x}$ can be written as $\mathbf{y} \propto \mathbf{x}$
$ \mathbf{x} $	Absolute value
$ 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
$\arg \max_{\mathbf{x}} f(\mathbf{x})$	Argmax: the value \mathbf{x} that maximizes f
$B(a, b)$	Beta function, $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$
$B(\vec{\alpha})$	Multivariate beta function, $\frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)}$
$\binom{n}{k}$	n choose k , equal to $n!/(k!(n-k)!)$
$\delta(\mathbf{x})$	Dirac delta function, $\delta(\mathbf{x}) = \infty$ if $\mathbf{x} = 0$, else $\delta(\mathbf{x}) = 0$
$\exp(\mathbf{x})$	Exponential function $e^{\mathbf{x}}$
$\Gamma(\mathbf{x})$	Gamma function, $\Gamma(\mathbf{x}) = \int_0^\infty u^{\mathbf{x}-1} e^{-u} du$
$\Psi(\mathbf{x})$	Digamma function, $\Psi(\mathbf{x}) = \frac{d}{d\mathbf{x}} \log \Gamma(\mathbf{x})$
\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 $\vec{\mathbf{x}}$, and boldface upper-case to denote matrices, such as $\vec{\mathbf{X}}$. 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 $(\mathbf{x}_1, \dots, \mathbf{x}_D)$ to denote a column vector created by stacking D scalars. If we write $\vec{\mathbf{X}} = (\vec{\mathbf{x}}_1, \dots, \vec{\mathbf{x}}_n)$, where the left hand side is a matrix, we mean to stack the $\vec{\mathbf{x}}_i$ along the columns, creating a matrix.

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

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 $\boldsymbol{p}()$ 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
$p(\mathbf{x})$	Probability mass function(PMF)
$f(\mathbf{x})$	probability density function(PDF)
$F(\mathbf{x}, \mathbf{y})$	Joint CDF
$p(\mathbf{x}, \mathbf{y})$	Joint PMF
$f(\mathbf{x}, \mathbf{y})$	Joint PDF
$p(X Y)$	Conditional PMF, also called conditional probability
$f_{X Y}(\mathbf{x} \mathbf{y})$	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 \boldsymbol{p}
$\tilde{\alpha}$	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 \boldsymbol{q}
$\mathbb{H}(X)$ or $\mathbb{H}(\boldsymbol{p})$	Entropy of distribution $\boldsymbol{p}(X)$
$\mathbb{I}(X; Y)$	Mutual information between X and Y
$\mathbb{KL}(\boldsymbol{p} \boldsymbol{q})$	KL divergence from distribution \boldsymbol{p} to \boldsymbol{q}
$\ell(\tilde{\theta})$	Log-likelihood function
$L(\boldsymbol{\theta}, \boldsymbol{a})$	Loss function for taking action \boldsymbol{a} when true state of nature is $\boldsymbol{\theta}$
λ	Precision (inverse variance) $\lambda = 1/\sigma^2$
Λ	Precision matrix $\Lambda = \Sigma^{-1}$
$\text{mode}[\tilde{X}]$	Most probable value of \tilde{X}

μ	Mean of a scalar distribution
$\vec{\mu}$	Mean of a multivariate distribution
Φ	cdf of standard normal
ϕ	pdf of standard normal
$\vec{\pi}$	multinomial parameter vector, Stationary distribution of Markov chain
ρ	Correlation coefficient
$\text{sigm}(\mathbf{x})$	Sigmoid (logistic) function, $\frac{1}{1 + e^{-x}}$
σ^2	Variance
Σ	Covariance matrix
$\text{var}[\mathbf{x}]$	Variance of \mathbf{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 \mathbf{x} to represent an observed data vector. In a supervised problem, we use y or \tilde{y} to represent the desired output label. We use \vec{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 $\mathcal{D} = \{(\vec{x}_i, y_i) i = 1 : N\}$
\mathcal{D}_{test}	Test data
\mathcal{X}	Input space
\mathcal{Y}	Output space
K	Number of states or dimensions of a variable (often latent)

$k(\mathbf{x}, y)$	Kernel function
\tilde{K}	Kernel matrix
\mathcal{H}	Hypothesis space
L	Loss function
$J(\tilde{\theta})$	Cost function
$f(\tilde{\mathbf{x}})$	Decision function
$P(y \tilde{\mathbf{x}})$	TODO
λ	Strength of ℓ_2 or ℓ_1 <i>regularizer</i>
$\phi(\mathbf{x})$	Basis function expansion of feature vector $\tilde{\mathbf{x}}$
Φ	Basis function expansion of design matrix \tilde{X}
$q()$	Approximate or proposal distribution
$Q(\tilde{\theta}, \tilde{\theta}_{old})$	Auxiliary function in EM
T	Length of a sequence
$T(\mathcal{D})$	Test statistic for data
\tilde{T}	Transition matrix of Markov chain
$\tilde{\theta}$	Parameter vector
$\tilde{\theta}^{(s)}$	s 'th sample of parameter vector
$\hat{\tilde{\theta}}$	Estimate (usually MLE or MAP) of $\tilde{\theta}$
$\hat{\tilde{\theta}}_{MLE}$	Maximum likelihood estimate of $\tilde{\theta}$
$\hat{\tilde{\theta}}_{MAP}$	MAP estimate of $\tilde{\theta}$
$\bar{\tilde{\theta}}$	Estimate (usually posterior mean) of $\tilde{\theta}$
$\tilde{\mathbf{w}}$	Vector of regression weights (called $\tilde{\beta}$ in statistics)
b	intercept (called ϵ in statistics)
\tilde{W}	Matrix of regression weights
\mathbf{x}_{ij}	Component (i.e., feature) j of data case i , for $i = 1 : N, j = 1 : D$
$\tilde{\mathbf{x}}_i$	Training case, $i = 1 : N$
\tilde{X}	Design matrix of size $N \times D$
$\bar{\tilde{\mathbf{x}}}$	Empirical mean $\bar{\tilde{\mathbf{x}}} = \frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{x}}_i$
$\tilde{\tilde{\mathbf{x}}}$	Future test case
$\tilde{\mathbf{x}}_*$	Feature test case
$\tilde{\mathbf{y}}$	Vector of all training labels $\tilde{\mathbf{y}} = (y_1, \dots, y_N)$
z_{ij}	Latent component j for case i

2 Future Research

There are many future challenges in machine learning generally and in the application of machine learning to health informatics specifically. The ultimate goal is to design and develop algorithms which can learn from data, hence can improve with experience over time. Interactive Machine Learning (iML) is a relatively new approach and can be of particular important to solve problems in health informatics, where we are lacking big data sets, deal with complex data and/or rare events, where traditional learning algorithms suffer due to insufficient training samples. Here the doctor-in-the-loop can help, where human expertise and long-term experience can assist in solving hard problems.

- TODO
- TODO
- TODO

3 References