
Mathematical Notations for Course LV 185.A83

Machine Learning for Health Informatics

Andreas Holzinger
Vienna University of Technology, Austria
andreas.holzinger@univie.ac.at

First version April, 12, 2016; Current as of March, 31, 2017

Abstract

Machine learning (ML) is a very practical field. Consequently, the theoretical-mathematical content of the course 185.A83 "Machine Learning for Health Informatics" is kept to a minimum, but for understanding the basics we need this minimum maths skill set. After understanding some principles of practical ML, students will rapidly become interested in the underlying theoretical principles and then the mathematical interest will also raise. This manual shall provide a notation to foster a consistent notation throughout the class to cover the extremely wide variety of data, models and algorithms discussed in this course. Definitions, conventions and usage of one and the same expression can be extremely different in mathematics, in statistics, and in computer science. Consequently, the goal of this short manual is to provide a help for students to some of the most used ML notations. Always consider that one and the same symbol may have different meaning in different context.

1 Introduction and Motivation

Machine learning heavily builds on the three pillars of linear algebra, optimization and probability/statistics, although many other mathematical areas are involved.

The typical data organization is in form of a 2D array, where the rows represent the samples (data items) and the columns represent the attributes (features), which can be seen as a vector of attributes and the array as a matrix.

In this short document the student will find some notations of mathematical foundations, linear algebra, probability, specific ML notations, and graphical model notations, finally some recommendations to mathematical literature what we need for the courses LV 185.A83, LV 706.315, LV 709.049. A very brief introduction to the application area health informatics can be found in [1], and some recent research trends can be found in [2].

2 Most used - at a glance

s is a scalar, \vec{a} a vector, \mathbf{A} a matrix, \mathbf{A} a tensor, \mathbb{A} a set, \mathcal{G} a graph;

$\mathbb{E}[f(x)]$ Expected value of function $f(x)$ with respect to $p(x)$

$\vec{\theta}$ Parameter vector (set of parameters that generated (x, y) ; the goal of ML is to estimate theta from given x s and y s

\mathcal{GP} Gaussian process $f \sim \mathcal{GP}(m(\vec{x}), k(\vec{x}, \vec{x}'))$; where the function f is distributed as a Gaussian process with the mean function $m(\vec{x})$ and the covariance function $k(x, x')$

3 Mathematical Foundations

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
$\vec{x} \otimes \vec{y}$	Convolution of \vec{x} and \vec{y}
$\vec{x} \odot \vec{y}$	Hadamard (elementwise) product of \vec{x} and \vec{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., $n \rightarrow \infty$
\leftarrow	in an algorithm: assign to variable t the new value $t + 1$, e.g., $t \leftarrow t + 1$
\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
$\arg \min_x f(x)$	Argmin: the value x that minimizes f
$B(a, b)$	Beta function, $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$
$B(\vec{\alpha})$	Multivariate beta function, $B(\vec{\alpha}) = \frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)}$
$n!$	n factorial = $n * (n - 1) * (n - 2) * \dots * 1$
$\binom{n}{k} = \frac{n!}{k!(n-k)!}$	n choose k , equal to $n!/(k!(n-k)!)$
$\delta(x)$	Dirac delta function, $\delta(x) = \infty$ if $x = 0$, else $\delta(x) = 0$
$\Gamma(x)$	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., $\mathcal{X} = \mathbb{R}^D$)
\equiv	equivalent to (or defined to be)
$\lim_{a \rightarrow \infty} f(x)$	the value of $f(x)$ in the limit as x approaches a
$m \bmod n$	m modulo n , the remainder when m is divided by n (e.g. $7 \bmod 5 = 2$)
\ln	logarithm base e , or natural logarithm of x
\log	logarithm base 10 of x
\log_2	logarithm base 2 of x
$\exp(x)$ or e^x	exponential of x , i.e., e raised the power of x
$\partial f(x)/\partial x$	partial derivative of f with respect to x
$\int_a^b f(x) dx$	the integral of $f(x)$ between a and b . If no limits are written, the full space is assumed.
$F(X; \theta)$	function of x , with implied dependence upon θ
$\langle x \rangle$	expected value of random variable x
\bar{x}	mean or average value of x
$\mathcal{E}[f(x)]$	the expected value of function $f(x)$ where x is a random variable
$\mathcal{E}_y[f(x, y)]$	the expected value of function over several variables, $f(x)$, taken over a subset y of them
$\sum_{i=1}^n a_i$	the sum from $i = 1$ to n : $a_1 + a_2 + \dots + a_n$
$\prod_{i=1}^n a_i$	the product from $i = 1$ to n : $a_1 * a_2 * \dots * a_n$

$f(x) * g(x)$	convolution of $f(x)$ with $g(x)$
$\mathcal{A}, \mathcal{B}, \mathcal{C}, \mathcal{D}, \dots$	"Calligraphic" font generally denotes sets or lists, e.g., data set $\mathcal{D} = x_1, \dots, x_n$
$x \in \mathcal{D}$	x is an element of set \mathcal{D}
$x \notin \mathcal{D}$	x is not an element of set \mathcal{D}
$\mathcal{D} \cup \mathcal{D}$	x union of two sets, i.e., the set containing all elements of \mathcal{D} and \mathcal{D}
$ \mathcal{D} $	cardinality of set \mathcal{D} , i.e., the number of (possibly non-distinct) elements in it
$\max_x[\mathcal{D}]$	the maximum x value in set \mathcal{D}
$dom(x)$	Domain of variable x
$x = x$	The variable x is in the state x
$dim(x)$	For a discrete variable x , this denotes the number of states x can take
$x_{a:b}$	x_a, x_{a+1}, \dots, x_b
\nless, \ngtr	not less than; not greater than
\neq	not equal to
\ll, \gg	much less than; much greater than
d/dx	the derivative with respect to x
$\mathcal{M} \subset \mathcal{N}$	\mathcal{M} is a subset of \mathcal{N}
$\mathcal{M} \supset \mathcal{N}$	\mathcal{M} contains \mathcal{N}
$\mathcal{M} \cap \mathcal{N}$	intersection of \mathcal{M} and \mathcal{N}
\implies	implies
\iff	equivalent to
\exists	there exists
\forall	for every

4 Linear algebra notations

We use boldface lower-case to denote vectors, such as \vec{x} , and boldface upper-case to denote matrices, such as \vec{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 $\vec{X} = (\vec{x}_1, \dots, \vec{x}_n)$, where the left hand side is a matrix, we mean to stack the \vec{x}_i along the columns, creating a matrix.

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

$\vec{x} \otimes \vec{y}$	Tensor product of \vec{x} and \vec{y}
\mathbb{R}^d	d-dimensional Euclidean space
$\mathbf{x}, \mathbf{A}, \dots$	boldface is used for (column) vectors and matrices
$f(x)$	vector-valued function (note the boldface) of a scalar
$f(\chi)$	vector-valued function (note the boldface) of a vector
I	identity matrix, square matrix having 1s on the diagonal and 0 everywhere else
Σ	covariance matrix
λ	eigenvalue
\mathbf{e}	eigenvector
\mathbf{u}_i	unit vector in the i th direction in Euclidean space
$\dim x$	The dimension of vector/matrix x

5 Probability notations

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 $p(\cdot)$ for both discrete and continuous random variables

Symbol	Meaning
X, Y	Random variable
$P(\cdot)$	Probability of a random event
$F(\cdot)$	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
$\vec{\alpha}$	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 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
$\ell(\vec{\theta})$	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 $\Lambda = \Sigma^{-1}$
$\text{mode}[\vec{X}]$	Most probable value of \vec{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
$\text{sigm}(x)$	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
\sim	has the distribution, e.g., $p(x) \sim N(\mu, \sigma^2)$
$N(\mu, \sigma^2)$	multidimensional normal or Gaussian distribution with mean μ and variance σ^2
$O(h(x))$	big oh order of $h(x)$
$\Theta(h(x))$	big theta order of $h(x)$
$\Omega(h(x))$	big omega order of $h(x)$
$\sup f(x)$	the supremum value of $f(x)$ -the global maximum of $f(x)$ over all values of x
$p(x = tr)$	Probability of variable x being in the state true
$p(x = fa)$	Probability of variable x being in the state false
$p(x \cap y)$	Probability of x and y
$p(x \cup y)$	Probability of x or y
$p(x y)$	Probability of x conditioned on y
$\langle f(x) \rangle_{g(x)}$	The average of the function $f(x)$ with respect to the distribution $g(x)$
$\sigma(x)$	The logistic sigmoid $\frac{1}{1+\exp(-x)}$
$\text{erf}(x)$	The (Gaussian) error function

6 Specific Machine learning notations

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 \vec{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.

We use uppercase bold roman letters to denote matrices \mathbf{M}

Symbol	Meaning
C	Number of classes
D	Dimensionality of data vector (number of features - of a feature vector gained)
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)
$k(x, y)$	Kernel function
\vec{K}	Kernel matrix
\mathcal{H}	Hypothesis space
L	Loss function
$J(\vec{\theta})$	Cost function
$f(\vec{x})$	Decision function
$P(y \vec{x})$	TODO
λ	Strength of ℓ_2 or ℓ_1 regularizer
$\phi(x)$	Basis function expansion of feature vector \vec{x}
Φ	Basis function expansion of design matrix \vec{X}
$q()$	Approximate or proposal distribution
$Q(\vec{\theta}, \vec{\theta}_{old})$	Auxiliary function in EM
T	Length of a sequence
$T(\mathcal{D})$	Test statistic for data

\vec{T}	Transition matrix of Markov chain
$\vec{\theta}$	Parameter vector
$\vec{\theta}^{(s)}$	s 'th sample of parameter vector
$\hat{\theta}$	Estimate (usually MLE or MAP) of $\vec{\theta}$
$\hat{\theta}_{MLE}$	Maximum likelihood estimate of $\vec{\theta}$
$\hat{\theta}_{MAP}$	MAP estimate of $\vec{\theta}$
$\bar{\theta}$	Estimate (usually posterior mean) of $\vec{\theta}$
\vec{w}	Vector of regression weights (called $\vec{\beta}$ in statistics)
\mathbf{b}	intercept (called ε in statistics)
\vec{W}	Matrix of regression weights
x_{ij}	Component (i.e., feature) j of data case i , for $i = 1 : N, j = 1 : D$
\vec{x}_i	Training case, $i = 1 : N$
\vec{X}	Design matrix of size $N \times D$
\bar{x}	Empirical mean $\bar{x} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i$
\tilde{x}	Future test case
\vec{x}_*	Feature test case
\vec{y}	Vector of all training labels $\vec{y} = (y_1, \dots, y_N)$
z_{ij}	Latent component j for case i
S	Number of samples

7 Graphical model notations

In graphical models, we index nodes by $s, t, u \in V$, and states by $i, j, k \in \mathcal{X}$.

Symbol	Meaning
$\tilde{s}t$	Node s is connected to node t
bel	Belief function
\mathcal{C}	Cliques of a graph
ch_j	Child of node j in a DAG (directed acyclic graph)
$desc_j$	Descendants of node j in a DAG
G	A graph
\mathcal{E}	Edges of a graph
mb_t	Markov blanket of node t
nbd_t	Neighborhood of node t
pa_t	Parents of node t in a DAG
$pred_t$	Predecessors of node t in a Direct Acyclic Graph (DAG) with respect to some ordering
$\psi_c(x_c)$	Potential function for clique c
\mathcal{S}	Separators of a graph
θ_{sjk}	prob. node s is in state k given its parents are in states j
\mathcal{V}	Nodes of a graph
$pa(x)$	The parents of x
$ch(x)$	The children of x
$ne(x)$	The neighbours of x
$\tilde{i}j$	The set of unique neighbouring edges on a graph

8 Abbreviations (incomplete)

cdf ... cumulative distribution function DAG ... directed acyclic graph HMM ... Hidden Markov Model iff ... if and only if pmf ... probability mass function

9 Recommended Literature (incomplete)

To strengthen the mathematical understanding the following textbooks can be recommended:

Dan SIMOVICI & Chabane DJERABA (2014) *Mathematical Tools for Data Mining: Set Theory, Partial Orders, Combinatorics*, Second Edition. London, Heidelberg, New York, Dordrecht: Springer [3]. This is a must-have book on every desk, a comprehensive compendium of the maths we need in our daily work, includes topologies and measures in metric spaces.

Keneth H. ROSEN (2013) *Discrete Mathematics and its Applications*. New York: McGraw-Hill [4]. This discrete mathematics course book spans a thread through mathematical reasoning, combinatorial analysis, discrete structures, algorithmic thinking and applications as well as modeling very recommendable.

Richard O. DUDA, Peter E. HART & David G. STORK (2001) *Pattern Classification*. New York: John Wiley [5]. This is THE classic work from Bayesian Decision Theory, Nonparametric Techniques, Linear Discriminant Functions and Stochastic Methods with a useful and applicable mathematical foundation. A must-have for any data scientist.

About the Lecturer

Andreas Holzinger is lead of the Holzinger Group HCI-KDD at the Institute for Medical Informatics, Statistics and Documentation at the Medical University Graz, and Associate Professor of Applied Computer Science at the Institute of Interactive Systems and Data Science at the Faculty of Computer Science and Biomedical Engineering at Graz University of Technology. Andreas Holzinger is Visiting Professor for Machine Learning in Health Informatics at the Faculty of Informatics at Vienna University of Technology. He serves as consultant for the Canadian, US, UK, Swiss, French, Italian and Dutch governments, for the German Excellence Initiative, and as national expert in the European Commission. His research interests are in supporting human intelligence with machine intelligence to help to solve problems in health informatics. Andreas obtained a Ph.D. in Cognitive Science from Graz University in 1998 and his Habilitation (second Ph.D.) in Computer Science from Graz University of Technology in 2003. Andreas was Visiting Professor in Berlin, Innsbruck, London (2 times), and Aachen. Andreas founded the international Expert Network HCI-KDD to foster a synergistic combination of methodologies of two areas that offer ideal conditions towards unraveling problems in understanding complex data: Human-Computer Interaction (HCI) and Knowledge Discovery from Data (KDD), with the goal of supporting human intelligence with machine learning for knowledge discovery. Andreas Holzinger is Associate Editor of Springer Knowledge and Information Systems (KAIS), Section Editor for Machine Learning of BMC Medical Informatics and Decision Making (MIDM) and member of IFIP WG 12.9 Computational Intelligence. Personal homepage: <http://aholzinger.at>

References

- [1] Andreas Holzinger. *Biomedical Informatics: Discovering Knowledge in Big Data*. Springer, New York, 2014. doi: 10.1007/978-3-319-04528-3.
- [2] Andreas Holzinger. Machine learning for health informatics. In Andreas Holzinger, editor, *Machine Learning for Health Informatics: State-of-the-Art and Future Challenges, Lecture Notes in Artificial Intelligence LNAI 9605*, pages 1–24. Springer, Cham, 2016. URL http://dx.doi.org/10.1007/978-3-319-50478-0_1.
- [3] Dan A Simovici and Chabane Djeraba. *Mathematical tools for data mining. Second Edition*. 2014. doi: 10.1007/978-1-4471-6407-4.
- [4] Kenneth H Rosen and Kamala Krithivasan. *Discrete mathematics and its applications. Seventh Edition*. McGraw-Hill, New York, 2012.
- [5] Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern Classification. Second Edition*. Wiley, New York et al., 2000.