

Machine Learning Notation

Shan-Hung Wu

1 Numbers & Arrays

a	A scalar (integer or real)
A	A scalar constant
\mathbf{a}	A vector
\mathbf{A}	A matrix
\mathbf{A}	A tensor
\mathbf{I}_n	The $n \times n$ identity matrix
\mathbf{D}	A diagonal matrix
$\text{diag}(\mathbf{a})$	A square, diagonal matrix with diagonal entries given by \mathbf{a}
a	A scalar random variable
\mathbf{a}	A vector-valued random variable
\mathbf{A}	A matrix-valued random variable

2 Sets & Graphs

\mathbb{A}	A set
\mathbb{R}	The set of real numbers
$\{0, 1\}$	The set containing 0 and 1
$\{0, 1, \dots, n\}$	The set of all integers between 0 and n
$[a, b]$	The real interval including a and b
$(a, b]$	The real interval excluding a but including b
$\mathbb{A} \setminus \mathbb{B}$	Set subtraction, i.e., the set containing the elements of \mathbb{A} that are not in \mathbb{B}
\mathcal{G}	A graph whose each vertex $\mathbf{x}^{(i)}$ denotes a random variable and edge denotes conditional dependency (directed) or correlation (undirected)
$\text{Pa}(\mathbf{x}^{(i)})$	The parents of a vertex $\mathbf{x}^{(i)}$ in \mathcal{G}

3 Indexing

a_i	Element i of vector \mathbf{a} , with indexing starting at 1
a_{-i}	All elements of vector \mathbf{a} except for element i
$A_{i,j}$	Element (i, j) of matrix \mathbf{A}
$\mathbf{A}_{i,:}$	Row i of matrix \mathbf{A}
$\mathbf{A}_{:,i}$	Column i of matrix \mathbf{A}
$\mathbf{A}_{i,j,k}$	Element (i, j, k) of a 3-D tensor \mathbf{A}
$\mathbf{A}_{:, :, i}$	2-D slice of a 3-D tensor
\mathbf{a}_i	Element i of the random vector \mathbf{a}

4 Functions

$f : \mathbb{A} \rightarrow \mathbb{B}$	A function f with domain \mathbb{A} and range \mathbb{B}
$f \circ g$	Composition of functions f and g
$f(\mathbf{x}; \boldsymbol{\theta})$	A function of \mathbf{x} parametrized by $\boldsymbol{\theta}$ (with $\boldsymbol{\theta}$ omitted sometimes)
$\ln x$	Natural logarithm of x
$\sigma(x)$	Logistic sigmoid, i.e., $(1 + \exp(-x))^{-1}$
$\zeta(x)$	Softplus, $\ln(1 + \exp(x))$
$\ \mathbf{x}\ _p$	L^p norm of \mathbf{x}
$\ \mathbf{x}\ $	L^2 norm of \mathbf{x}
x^+	Positive part of x , i.e., $\max(0, x)$
$1(x; \text{cond})$	The indicator function of x : 1 if the condition is true, 0 otherwise
$g[f; x]$	A functional that maps f to $f(x)$

Sometimes we use a function f whose argument is a scalar, but apply it to a vector, matrix, or tensor: $f(\mathbf{x})$, $f(\mathbf{X})$, or $f(\mathbf{X})$. This means to apply f to the array element-wise. For example, if $\mathbf{C} = \sigma(\mathbf{X})$, then $C_{i,j,k} = \sigma(X_{i,j,k})$ for all i, j and k .

5 Calculus

$f'(a)$ or $\frac{df}{dx}(a)$	Derivative of $f : \mathbb{R} \rightarrow \mathbb{R}$ at input point a
$\frac{\partial f}{\partial x_i}(\mathbf{a})$	Partial derivative of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with respect to x_i at input \mathbf{a}
$\nabla f(\mathbf{a}) \in \mathbb{R}^n$	Gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at input \mathbf{a}
$\nabla f(\mathbf{A}) \in \mathbb{R}^{m \times n}$	Matrix derivatives of $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ at input \mathbf{A}
$\nabla f(\mathbf{A})$	Tensor derivatives of f at input \mathbf{A}
$\mathbf{J}(f)(\mathbf{a}) \in \mathbb{R}^{m \times n}$	The Jacobian matrix of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at input \mathbf{a}
$\nabla^2 f(\mathbf{a})$ or $\mathbf{H}(f)(\mathbf{a}) \in \mathbb{R}^{n \times n}$	The Hessian matrix of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at input point \mathbf{a}
$\int f(\mathbf{x}) d\mathbf{x}$	Definite integral over the entire domain of \mathbf{x}
$\int_{\mathbb{S}} f(\mathbf{x}) d\mathbf{x}$	Definite integral with respect to \mathbf{x} over the set \mathbb{S}

6 Linear Algebra

\mathbf{A}^\top	Transpose of matrix \mathbf{A}
\mathbf{A}^\dagger	Moore-Penrose pseudo-inverse of \mathbf{A}
$\mathbf{A} \odot \mathbf{B}$	Element-wise (Hadamard) product of \mathbf{A} and \mathbf{B}
$\det(\mathbf{A})$	Determinant of \mathbf{A}
$\text{tr}(\mathbf{A})$	Trace of \mathbf{A}
$\mathbf{e}^{(i)}$	The i -th standard basis vector (a one-hot vector)

7 Probability & Info. Theory

$a \perp b$	Random variables a and b are independent
$a \perp b \mid c$	They are conditionally independent given c
$\Pr(a \mid b)$ or $\Pr(a \mid b)$	Shorthand for the probability $\Pr(a = a \mid b = b)$
$P_a(a)$	A probability mass function of the discrete random variable a
$p_a(a)$	A probability density function of the continuous random variable a
$P(a = a)$	Either $P_a(a)$ or $p_a(a)$
$P(\theta)$	A probability distribution parametrized by θ
$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	The Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$
$x \sim P(\theta)$	Random variable x has distribution P
$E_{x \sim P}[f(x)]$	Expectation of $f(x)$ with respect to P
$\text{Var}[f(x)]$	Variance of $f(x)$
$\text{Cov}[f(x), g(x)]$	Covariance of $f(x)$ and $g(x)$
$H(x)$	Shannon entropy of the random variable x
$D_{\text{KL}}(P \parallel Q)$	Kullback-Leibler (KL) divergence from distribution Q to P

8 Machine Learning

\mathbb{X}	The set of training examples
N	Size of \mathbb{X}
$(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$	The i -th example pair in \mathbb{X} (supervised learning)
$\mathbf{x}^{(i)}$	The i -th example in \mathbb{X} (unsupervised learning)
D	Dimension of a data point $\mathbf{x}^{(i)}$
K	Dimension of a label $\mathbf{y}^{(i)}$
$\mathbf{X} \in \mathbb{R}^{N \times D}$	Design matrix, where $\mathbf{X}_{i,:}$ denotes $\mathbf{x}^{(i)}$
$P(\mathbf{x}, \mathbf{y})$	A data generating distribution
\mathbb{F}	Hypothesis space of functions to be learnt, i.e., a model
$C[f]$	A cost functional of $f \in \mathbb{F}$
$C(\theta)$	A cost function of θ parametrizing $f \in \mathbb{F}$
$(\mathbf{x}', \mathbf{y}')$	A testing pair
$\hat{\mathbf{y}}$	Label predicted by a function f , i.e., $\hat{\mathbf{y}} = f(\mathbf{x}')$ (supervised learning)

9 Typesetting

Section*	Section that can be skipped for the first time reading
Section**	Section for reference only (will not be taught)
[Proof]	Prove it yourself
[Homework]	You have homework

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, $\frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)}$
$n!$	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}
$\text{dom}(x)$	Domain of variable x
$x = x$	The variable x is in the state x
$\text{dim}(x)$	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
\longleftrightarrow	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
$\text{tr}(\vec{X})$	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
$\text{diag}(x)$	Diagonal matrix made from vector \vec{x}
$\text{diag}(\vec{X})$	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}
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 $p()$ 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
$\vec{\alpha}$	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 $h(x)$
$\sup_x f(x)$	the supremum value of $f(x)$ -the global maximum of $f(x)$ over all values of x
$p(x = \text{true})$	Probability of variable x being in the state true
$p(x = \text{false})$	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 $p(x)$
$\sigma(x)$	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 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 **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 $\mathcal{D} = \{(\vec{x}_i, y_i) i = 1 : N\}$
$\mathcal{D}_{\text{test}}$	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
$Q(\vec{\theta}, \vec{\theta}_{\text{old}})$	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
$\hat{\vec{\theta}}$	Estimate (usually MLE or MAP) of $\vec{\theta}$
$\hat{\vec{\theta}}_{MLE}$	Maximum likelihood estimate of $\vec{\theta}$
$\hat{\vec{\theta}}_{MAP}$	MAP estimate of $\vec{\theta}$
$\bar{\vec{\theta}}$	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 $N \times D$
$\bar{\vec{x}}$	Empirical mean $\bar{\vec{x}} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i$
$\tilde{\vec{x}}$	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