

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
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