

Symbols

In general, scalars are printed as normal letters, e. g. \mathbf{p}, \mathbf{q} ; vectors are bolded, e. g. $\boldsymbol{\theta}$; and matrices are denoted by bold upper-case letters, e. g. $\boldsymbol{\Phi}$.

Latin letters

\mathbf{a}	action (RL)
\mathbf{A}	adjacency matrix
\mathcal{A}	action space (RL)
C	inverse regularization parameter (SVMs), cluster
\mathcal{C}	set of classes
\mathcal{C}_k	k-th class
D	number of dimensions
\mathbf{D}	degree matrix
\mathcal{D}	data set, $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}_{i=1}^n$
\mathcal{D}_j	j-th partition of data set \mathcal{D}
\mathcal{E}	set of edges of a graph
\mathcal{F}	feature space
\mathcal{G}	graph, goal state (RL)
$\mathbf{h} / \mathbf{h}_{\boldsymbol{\theta}}$	hypothesis / model
H	kernel function for kernel density estimation
\mathbf{H}	Hessian matrix (matrix of second derivatives)
\mathcal{H}	hypothesis space
i	index over the data examples
\mathbf{I}	identity matrix
j	index over the features
\mathcal{J}	loss function / error function

k	index over classes, number of clusters / neighbors / base classifiers
\mathbf{K}	kernel matrix
\mathcal{K}	kernel function
ℓ	loss
\mathbf{L}	graph Laplacian matrix
\mathcal{L}	likelihood, lagrangian, loss
\mathcal{LL}	log-likelihood
\mathfrak{m}	number of features ($\mathfrak{j} = 1, \dots, \mathfrak{m}$)
M	number of mixture components
\mathcal{M}	margin
\mathfrak{n}	number of examples ($\mathfrak{i} = 1, \dots, \mathfrak{n}$)
N	node in a tree
p	degree of a polynomial
Q	true q value (RL)
\hat{Q}	approximate q value (RL)
r	cost ratio, c_{fp}/c_{fn} , reward (RL)
\mathcal{R}	region
s	state (RL)
s_0	initial state (RL)
s'	next state (RL)
\mathcal{S}	state space (RL)
T	sequence length
v	bin volume
\mathbf{v}	vector
V^π	value function when using policy π (RL)

\mathcal{V}	vocabulary, set of vertices (nodes) of a graph
w	window size
\mathbf{x}	data instance
$\mathbf{x}^{(i)}$	i -th data instance
x_j	j -th feature of \mathbf{x}
$x_j^{(i)}$	j -th feature of the i -th data instance
\mathbf{x}_\perp	orthogonal projection of \mathbf{x}
$\hat{\mathbf{x}}$	data instance with attached 1 entry, $[1\mathbf{x}]$
\mathbf{X}	design matrix / feature matrix / regressor matrix
\mathcal{X}	data input space
y	label
\mathbf{y}	label vector
z	center of radial basis function, activation

Greek letters

α	learning rate, responsibility (mixture models), model weight
β	precision $\equiv \sigma^{-1}$
γ	variance threshold (PCA), discount factor (RL)
δ	error gradient (neural networks), state transition function (RL)
ε	noise, error rate
Θ, θ	trainable parameters of the model
κ	number of classes ($\mathbf{k} = 1, \dots, \kappa$)
λ	regularization parameter, eigenvalue, lagrange multiplier
μ	mean of a distribution
ξ	slack variable (SVMs)

π	prior (mixture models), policy (RL)
π^*	optimal policy (RL)
σ	sigmoid, standard deviation
σ^2	variance
Σ	covariance matrix
τ	trajectory (RL)
$\varphi(\mathbf{x})$	basis function of \mathbf{x}
Φ	design matrix with basis functions

Mathematical symbols (• is a placeholder)

$\mathbb{1}\{\dots\}$	indicator function
$\arg \max_{\mathbf{x}} f(\mathbf{x})$	value of \mathbf{x} which maximizes $f(\mathbf{x})$
$\arg \min_{\mathbf{x}} f(\mathbf{x})$	value of \mathbf{x} which minimizes $f(\mathbf{x})$
$\text{Bern}(\mathbf{x} \mu)$	bernoulli distribution, $\mu^x(1 - \mu)^{1-x}$
$\text{Bin}(\mathbf{m} \mathbf{n}, \mu)$	binomial distribution, $\binom{n}{m}\mu^m(1 - \mu)^{n-m}$
$c(\mathcal{C}_i \mathcal{C}_j)$	cost for predicting class \mathcal{C}_i instead of \mathcal{C}_j
\mathcal{C}_{MAP}	maximum a posteriori class
$d(\dots)$	distance metric
$\text{dom}(f)$	domain of function f
$E(\mathcal{D})$	entropy of a data set, $-\sum_{\mathbf{c} \in \mathcal{C}} \mathbf{p}_{\mathbf{c}} \log_2 \mathbf{p}_{\mathbf{c}}$
$\exp\{\dots\}$	exponential function
$E\{\mathbf{x}\}$	expectation of a random variable \mathbf{x}
$E_{\mathbf{x} \sim \mathbf{p}(\mathbf{x})}\{f(\mathbf{x})\}$	expectation of $f(\mathbf{x})$ where \mathbf{x} is drawn from distribution $\mathbf{p}(\mathbf{x})$
$g(\mathbf{x})$	activation function, sigmoid
$\text{KL}(\mathbf{p} \mathbf{q})$	Kullback-Leibler divergence between \mathbf{p} and \mathbf{q}

$\max_{\mathbf{x}} f(\mathbf{x})$	maximum value of $f(\mathbf{x})$
$\min_{\mathbf{x}} f(\mathbf{x})$	minimum value of $f(\mathbf{x})$
$\mathcal{N}(\mathbf{x} \boldsymbol{\mu}, \sigma)$	normal (Gaussian) distribution
$\mathcal{N}(\boldsymbol{x} \boldsymbol{\mu}, \boldsymbol{\Sigma})$	multivariate normal (Gaussian) distribution
$\hat{p}(\mathbf{x})$	approximate probability of \mathbf{x}
$p(\mathbf{x})$	probability of \mathbf{x} , marginal probability
$p(\mathbf{x}, \mathbf{y})$	probability of \mathbf{x} and \mathbf{y} , joint probability
$p(\mathbf{x} \mathbf{y})$	probability of \mathbf{x} given \mathbf{y}
$R(\alpha_i \boldsymbol{x})$	risk of a decision α_i given \boldsymbol{x}
\mathbb{R}	set of real numbers
\mathbb{R}^+	set of positive real numbers (including 0)
$\text{ReLU}(\mathbf{x})$	ReLU activation function
$\tanh(\mathbf{x})$	tangent hyperbolic activation function
$\text{var}\{ \mathbf{x} \}$	variance of a random variable \mathbf{x}
$\mathbf{x} \sim p(\mathbf{x})$	\mathbf{x} is distributed according to probability distribution $p(\mathbf{x})$
\bullet^T	transpose of \bullet
\bullet^{-1}	inverse of \bullet (of a matrix for example)
$\bullet^\#$	pseudo-inverse of \bullet
$\boldsymbol{a}^T \boldsymbol{b}$	dot product of \boldsymbol{a} and \boldsymbol{b}
$\langle \bullet, \bullet \rangle$	dot product
$\boldsymbol{a} \boldsymbol{b}^T$	outer product
$ \bullet $	number of elements in the set \bullet
$\ \bullet\ $	vector norm
\angle	angle
$\binom{n}{k}$	binomial coefficient, n choose k

$\nabla_{\boldsymbol{\theta}}$	gradient with respect to $\boldsymbol{\theta}$ / nabla operator
∂	partial derivative
$\sum_{i=1}^n \dots$	sum
$\prod_{i=1}^n \dots$	product
\oplus	positive class
\ominus	negative class
$\bar{\mathbf{x}}$	sample set mean of \mathbf{x}
\star	convolution operator

Abbreviations

AUC	area-under-the-curve
CNN	convolutional neural network
fp	false positives
fn	false negatives
GRU	gated recurrent unit
i.i.d.	independent and identically distributed
Lasso	least absolute shrinkage and selection operator
LSTM	long short-term memory
MAE	mean absolute error
MDP	Markov decision process
MLE	maximum Likelihood Estimation
MLP	multi-layer perceptron
MNIST	
OCR	optical character recognition
PSD	positive semi-definite

ROC	receiver operating characteristic
RMSE	root mean square error
RNN	recurrent neural network
s. t.	subject to
SVM	support vector machine
tp	true positives
tn	true negatives
X-Val	cross validation