

Chapter ML:I

I. Introduction

- ❑ Examples of Learning Tasks
- ❑ Specification of Learning Tasks
- ❑ Elements of Machine Learning
- ❑ Notation Overview
- ❑ Classification Approaches Overview

Notation Overview

Data, Sets, and Distributions

Symbol	Semantics
x, x_i, x_1, \dots, x_p	Feature
$\mathbf{x} = (x_1, \dots, x_p)^T \in \mathbf{R}^p$	Feature vector
$\mathbf{x} = (1, x_1, \dots, x_p)^T \in \mathbf{R}^{p+1}$, i.e., $x_0 = 1$	Extended feature vector
\mathbf{X}	Feature space, Cartesian product of the domains of the p dimensions of a feature vector \mathbf{x} .
$X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$	Multiset of feature vectors
X	Random variable (randomness regarding feature x of an object o)
\mathbf{X}	Multivariate random variable, random vector (randomness regarding feature vector \mathbf{x} of an object o)

Notation Overview

Indexing

Running	Sequence	Semantics of maximum
\square_s	$\in \{\square_1, \dots, \square_d\}$	Number of layers in a multilayer perceptron
\square_i	$\in \{\square_1, \dots, \square_k\}$	Number of classes Number of folds during cross validation
\square_l	$\in \{\square_1, \dots, \square_m\}$	Number of elements in a domain of a feature Number of hyperparameter values during model selection
\square_i	$\in \{\square_1, \dots, \square_n\}$	Number of elements in a data set D
\square_j	$\in \{\square_1, \dots, \square_p\}$	Dimension of a feature space or a feature vector

Notation Overview

Functions

Function definition	Function name	Occurrence
$I_{\neq}(a, b) = \begin{cases} 0 & a = b \\ 1 & a \neq b \end{cases}$	Indicator function	Part II: Machine Learning Basics Part III: Linear Models
$f(x) = \dots$	function	Part :

Notation Overview

Algorithms

Signature	Algorithm name	Occurrence
$\text{LMS}(D, \eta)$	Least Mean Squares	Part I: Introduction Examples of Learning Tasks
$\text{ALG}(\dots)$	algorithm	Part : ...

Classification Approaches Overview

Search in hypothesis space

Taxonomy	Model function	Classification rule	Optimization principle	Optimization objective (loss/cost function [+ regularization])	Optimization approach (algorithm)	
Classification approaches	Discriminative approaches	$y(\mathbf{x}) = \text{heaviside}(\mathbf{w}^T \mathbf{x})$ Linear function: $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ Logistic function: $y(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$ SVM w/o kernel (aka linear kernel) $y(\mathbf{x}) = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}))$ $y(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}}$ SVM with nonlinear kernel $y(\mathbf{x}) = \sigma(W^o(\sigma^1(W^h \mathbf{x})))$ Multilayer percep.: $y(\mathbf{x}) = \bigwedge_{i=1}^p x_i = v_i$ Nominal feat. $\bigvee_{j=1, \dots, \text{leaves} } \bigwedge_{i=1, \dots, \text{depth}(l_i)} x_{ij} = v_{ij}$ Arbitrary features: DNF ($\bigvee_i \wedge_j$) on domain predicates	$\mathbf{w}^T \mathbf{x} \begin{cases} \geq 0 \\ < 0 \end{cases}$ $\mathbf{w}^T = (w_0, \dots, w_p)$ $x_0 = 1$ $\mathbf{w}^T \mathbf{x} - b \begin{cases} \geq 1 \\ \leq -1 \end{cases}$ $\mathbf{w}^T = (w_1, \dots, w_p)$ $\mathbf{w}^T \phi(\mathbf{x}) \begin{cases} \geq 0 \\ < 0 \end{cases}$ $\mathbf{w}^T = (w_0, \dots, w_{ \mathbf{w} })$ $\phi_0(\mathbf{x}) = 1$ $\mathbf{w}^T \phi(\mathbf{x}) - b \begin{cases} \geq 1 \\ \leq -1 \end{cases}$ $\mathbf{w} = \sum_{i=1}^n \alpha_i c \phi(\mathbf{x}_i)$ $\text{argmax}_{c \in C} \{ y_c(\mathbf{x}) \}$ Test if \mathbf{x} is a model for α (= fulfills α). α is a formula in DNF.	Exploit misclassified examples individually: Hebbian learning Linear regression Logistic regression Empirical risk minimization Linear regression (nonlinear in predictors) Log. regression (nonlinear in predictors) Empirical risk minimization Regression Maximize version space Decision tree: (greedy) feature-wise splitting of example set	No misclassified example Squared loss (residual sum of squares, RSS) + L_1 or L_2 norm on $\mathbf{w} _{1, \dots, p}$ Logistic loss (derived via ML) Regularized hinge loss Squared loss + L_1 or L_2 norm on $\mathbf{w} _{1, \dots, \mathbf{w} }$ Logistic loss Regularized hinge loss Squared loss (residual sum of squares, RSS) No misclassified example 0/1 Loss (= number of misclassified examples) + Tree height, external path length	Perceptron training algorithm Gradient descent: – batch – incremental – stochastic Newton-Raphson, BFGS Quadratic prog., sub-grad. descent Gradient descent: – batch – incremental – stochastic Newton-Raphson, BFGS Quadratic prog., sub-grad. descent Backpropagation algorithm Candidate elimination algorithm Algorithms: ID3, C4.5, C5.0, CART (exhaustive) search in space of domain splittings
	Generative approaches	Statistical approaches	$X \sim N(\mu, \sigma^2)$ (or other family) Bayes rule for combined conditional events	Maximum a-posteriori hypothesis	Goodness of fit, e.g. according to chi-squared, Kolmogorov-Smirnov	