

Chapter ML:I

I. Introduction

- Examples of Learning Tasks
- Specification of Learning Tasks
- Elements of Machine Learning
- Comparative Syntax Overview
- Classification Approaches

Comparative Syntax Overview

Concept	Stein et al. 2021	Bishop 2006	Hastie et al. 2009	Mitchell 1997
Feature	x, x_i, x_1, \dots, x_p			
Feature vector	\mathbf{x}			

Classification Approaches

Search in hypothesis space

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