

Artificial Intelligence

wu

May 8, 2019

Contents

1	Inference and Reasoning	2
1.1	Propositional logic	2
1.2	Predicate logic	2
1.3	First Order Inductive Learner	2
2	Statistical learning and modeling	2
2.1	Machine Learning: the concept	2
2.1.1	Example and concept	2
2.1.2	supervised learning: important concepts	3
2.2	example: polynomial curve fitting	3
2.3	probability theory review and notation	3
2.4	information theory	6
2.5	model selection	6
2.6	decision theory	6
3	Statistical learning and modeling - Supervised learning	6
3.1	Basic concepts	6
3.2	discriminant functions	8
3.2.1	Two classes	8
3.2.2	K-class	8
3.2.3	Learning the parameters of linear discriminant functions	8
3.3	probalibilistic generative models	10
3.4	probabilistic discriminative models	12
3.5	Boosting	12
3.5.1	AdaBoost	12
4	wef	12
4.1	wfe	12

1 Inference and Reasoning

1.1 Propositional logic

1.2 Predicate logic

1.3 First Order Inductive Learner

knowledge graph: node = entity, edge = relation. triplet (head entity, relation, tail entity)

2 Statistical learning and modeling

2.1 Machine Learning: the concept

2.1.1 Example and concept

Supervised learning problems applications in which the **training data** comprises examples of the input vectors along with their corresponding **target vectors** are known

classification and regression

Unsupervised learning problems the training data consists of a set of input vectors X **without any corresponding target values**

density estimation, clustering, hidden markov models

Reinforcement learning problem finding suitable actions to take in a given situation in order to maximize a reward. Here the learning algorithm is not given examples of optimal outputs, in contrast to supervised learning, but must instead discover them by a process of trial and error. A general feature of reinforcement learning is the trade-off between exploration and exploitation

types of machine learning

- supervised learning
 - classification: the output is categorical or nominal variable
 - regression: the output is real-valued variable
- unsupervised learning
- semi-supervised learning

- reinforcement learning
- deep learning

2.1.2 supervised learning: important concepts

- Data: labeled instances $\langle \mathbf{x}_i, \mathbf{y} \rangle$
- features: attribute-value pairs which characterize each \mathbf{x}
- learning a discrete function: **classification**
- learning a continuous function: **regression**

Classification - A two-step process

- **model construction**
- **model usage**

regression

- Example: price of a used car
 \mathbf{x} : car attributes. $\mathbf{y} = g(\mathbf{x} \mid \boldsymbol{\theta})$: price. g : model. $\boldsymbol{\theta}$ parameter set.

2.2 example: polynomial curve fitting

2.3 probability theory review and notation

rules of probability

- **sum rule** $p(X) = \sum_Y p(X, Y)$
- **product rule** $p(X, Y) = p(Y|X)p(X)$

Bayes' Theorem: $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$. Using sum rule $p(X) = \sum_Y p(X|Y)p(Y)$

probability densities.

$$p(x \in (a, b)) = \int_a^b p(x)dx$$

$$P(z) = \int_{-\infty}^z p(x)dx$$

$$\int_{-\infty}^{\infty} p(x)dx = 1 \quad p(x) \geq 0$$

$$\text{expectation } \mathbb{E}[f] = \begin{cases} \sum_x p(x)f(x) & \text{discrete variables} \\ \int p(x)f(x)dx & \text{continuous variables} \end{cases}. \text{ In either}$$

$$\text{cases, } \mathbb{E}[f] \approx \frac{1}{N} \sum_{n=1}^N f(x_n). \text{ conditional expectation: } \mathbb{E}_x[f|y] = \sum_x p(x|y)f(x).$$

The **variance** of $f(x)$ is

$$\begin{aligned} \text{var}[f] &= \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] \\ &= \mathbb{E}[f(x)^2 - 2f(x)\mathbb{E}[f(x)] + \mathbb{E}[f(x)]^2] \\ &= \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2 \end{aligned}$$

The **covariance** is

$$\begin{aligned} \text{cov}[x, y] &= \mathbb{E}_{x,y}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])] \\ &= \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y] \end{aligned}$$

the variance of the sum of two independent random variables is the sum of variance. Given

X	probability
x_1	p_1
\dots	\dots
x_n	p_n

Y	probability
y_1	q_1
\dots	\dots
y_m	q_m

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y)$$

In case of two vectors of random variables \mathbf{x} and \mathbf{y} , the covariance is a matrix

$$\begin{aligned} \text{cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}_{\mathbf{x}, \mathbf{y}}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T])] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^T] \end{aligned}$$

Bayesian probabilities: $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$. For a data set $\mathcal{D} = \{t_1, \dots, t_n\}$ and assumption w , $p(w|\mathcal{D}) = \frac{p(\mathcal{D}|w)p(w)}{p(\mathcal{D})}$. $p(w)$ is **prior probability**, $p(\mathcal{D}|w)$ is **likelihood** (the probability \mathcal{D} happens). Hence

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

Gaussian distribution.

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$

μ is called **mean**, σ^2 is called **variance**, σ **standard deviation**, $\beta = 1/\sigma^2$ **precision**

$$\begin{aligned} \mathbb{E}[x] &= \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x dx = \mu \\ \mathbb{E}[x^2] &= \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2 \\ \text{var}[x] &= \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2 \end{aligned}$$

For D -dimensional vector \mathbf{x} of continuous variables

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$$

To determine values for the unknown parameters given μ and σ^2 by maximizing the likelihood function. Use log.

$$\begin{aligned} P(\mathbf{X}|\mu, \sigma^2) &= \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2) \\ \Rightarrow \ln P(\mathbf{X}|\mu, \sigma^2) &= -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) \end{aligned}$$

Hence $\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$, $\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$ by partial derivative.

Maximum likelihood estimator for mean is unbiased, that is, $\mathbb{E}(\mu_{ML}) = \mu$. Maximum likelihood estimator for variance is biased. $\mathbb{E}(\sigma_{ML}^2) = \mathbb{E}(x^2) - \mathbb{E}(\mu_{ML}^2) = \frac{N-1}{N} \sigma_x^2$

2.4 information theory

entropy: measuring uncertainty of a random variable X . $H(X) = H(p) = - \sum_{x \in \Omega} p(x) \log p(x)$ where Ω is all possible values and define $0 \log 0 = 0$, $\log = \log_2$

$H(X) = \sum_{x \in \Omega} p(x) \log_2 \frac{1}{p(x)} = E(\log_2 \frac{1}{p(x)})$. And "information of x " = "#bits to code x " = $-\log p(x)$

Kullback-Leibler divergence: comparing two distributions

2.5 model selection

The technique of S -fold cross-validation, illustrated here for the case of $S = 4$, involves taking the available data and partitioning it into S groups (in the simplest case these are of equal size). Then $S - 1$ of the groups are used to train a set of models that are then evaluated on the remaining group. This procedure is then repeated for all S possible choices for the held-out group, indicated here by the red blocks, and the performance scores from the S runs are then averaged.



cross-validation

split training data into **training set** and **validation set**. Train different models on training set and choose model with minimum error on validation set.

2.6 decision theory

Suppose we have an input vector \mathbf{x} together with a corresponding vector \mathbf{t} of target variables and our goal is to predict \mathbf{t} given new value for \mathbf{x} . The joint probability distribution $p(\mathbf{x}, \mathbf{t})$ provides a complete summary of the uncertainty with these variables

3 Statistical learning and modeling - Supervised learning

3.1 Basic concepts

- **Linearly separable**

- decision regions:
input space is divided into several regions
- decision boundaries:

- * under linear models, it's a linear function
- * (D-1)-dimensional hyper-plane within the D-dimensional input space
- **representation of class labels**
 - Two classes $K = 2$
 - K classes
 - * 1-of-K coding scheme $\mathbf{t} = (0, 0, 1, 0, 0)^T$
 - Predict discrete class labels
 - * linear model prediction $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ w: weight vector, w_0 bias/threshold
 - * nonlinear function $f(\cdot) : \mathcal{R} \rightarrow (0, 1)$
 - * generalized linear models $y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$ f: activation function
 - * decision surface $y(\mathbf{x}) = \text{constant} \rightarrow \mathbf{w}^T \mathbf{x} + w_0 = \text{constant}$
- **Three classification approaches**
 - discriminant function
 - * least squares approach
 - * fisher's linear discriminant
 - * the perceptron algorithm of rosenblatt
 - use discriminant functions directly and don't compute probabilities

Given discriminant functions $f_1(\mathbf{x}), \dots, f_K(\mathbf{x})$. Classify \mathbf{x} as class \mathcal{C}_k iff $f_k(\mathbf{x}) > f_j(\mathbf{x}), \forall j \neq k$

 - * **least-squares approach**: making the model predictions as close as possible to a set of target values
 - * **fisher's linear discriminant**: maximum class separation in the output space
 - * **the perceptron algorithm of rosenblatt**
 - generative approach
 - * model the class-conditional densities and the class priors
 - * compute posterior probabilities through Bayes's theorem

$$\underbrace{p(\mathcal{C}_k | \mathbf{x})}_{\text{posterior for class}} = \frac{\overbrace{p(\mathbf{x} | \mathcal{C}_k)}^{\text{class conditional density}} \overbrace{p(\mathcal{C}_k)}^{\text{class prior}}}{p(\mathbf{x})} = \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j) p(\mathcal{C}_j)}$$

3.2 discriminant functions

3.2.1 Two classes

- Linear discriminant function $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - Decision surface $\Omega : y(\mathbf{x}) = 0$
 - the normal distant from the origin to the decision surface $\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$
 - if x_A, x_B lie on the decision surface $y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0$, then $\mathbf{w}^T(\mathbf{x}_A - \mathbf{x}_B) = 0$. hence \mathbf{w} is orthogonal to every vector lying within $\cdot \frac{\mathbf{w}}{\|\mathbf{w}\|}$ is the normal vector of
 - $\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$ hence $r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$. $y(\mathbf{x}_\perp) = 0 \rightarrow \mathbf{w}^T \mathbf{x} = -w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|}$
 - $\tilde{\mathbf{w}} = (w_0, \mathbf{w}), \tilde{\mathbf{x}} = (x_0, \mathbf{x}), y(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$

3.2.2 K-class

- One-versus-the-rest classifier K - 1 classifiers each of which solves a two-class problem
- One-versus-one classifier $K(K-1)/2$ binary discriminant functions
- single K-class discriminant comprising K linear functions $y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k_0}$
 - assigning a point \mathbf{x} to class \mathcal{C}_k if $y_k(\mathbf{x}) > y_j(\mathbf{x})$ for all $j \neq k$
 - decision boundary between class $\mathcal{C}_k, \mathcal{C}_j$ is given $y_k(\mathbf{x}) = y_j(\mathbf{x}) \rightarrow (\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k_0} - w_{j_0}) = 0$
 - \mathcal{R}_k is singly connected convex
 - $\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B$ where $0 \leq \lambda \leq 1$, $y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B)$ and hence $\hat{\mathbf{x}}$ also lies inside \mathcal{R}_k

3.2.3 Learning the parameters of linear discriminant functions

1. Linear basis function models **linear regression**: $y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D = \mathbf{w}^T \mathbf{x}$.

For nonlinear functions ϕ_j , $y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$

where $\phi_j(\mathbf{x})$ are **basis functions**

2. parameter optimization via maximum likelihood

Assume target variable t is given by a deterministic function $y(\mathbf{x}, \mathbf{w})$ with additive Gaussian noise so that $t = y(\mathbf{x}, \mathbf{w}) + \epsilon$ where ϵ is a zero mean Gaussian random variable with precision β , hence we can write

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

and $\mathbb{E}(t|\mathbf{x}) = \int t p(t|\mathbf{x}) dt = y(\mathbf{x}, \mathbf{w})$

For data set $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, $\mathbf{t} = (t_1, \dots, t_n)^T$, $p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$$

$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\}^2 = \frac{1}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2$ is sum-of-squares error function

solve \mathbf{w} by maximum likelihood.

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \right\} \phi(\mathbf{x}_n)^T$$

$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left(\sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right)$$

Hence we get

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Φ is **design matrix**.

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \dots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

For bias parameter w_0 . $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right\}^2$.

Hence $w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j$, $\bar{t} = \frac{1}{N} \sum_{n=1}^N t_n$, $\bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n)$.

$$frac{N}{2}\beta = E_D(\mathbf{w}). \quad \frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \left\{ t_n - \mathbf{w}_{ML}^T \phi(\mathbf{x}_n) \right\}^2$$

3. Least-squares approach

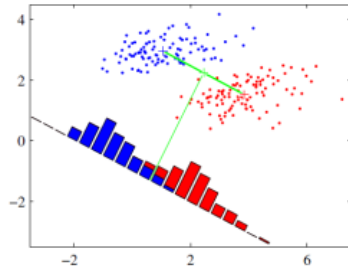
- Problem

- Each class \mathcal{C}_k is described by its own linear model $y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$
- group together: $y(\mathbf{x}) = \widetilde{\mathbf{W}}^T \tilde{\mathbf{x}}, \quad \tilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T, \quad \tilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$

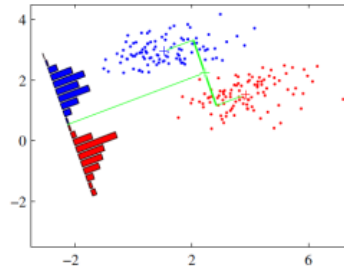
- Learning

- minimizing SSE function sum-of-squares $SSE = \sum_{i=1}^n (y_i - f(x_i))^2$
- $$E_D(\widetilde{\mathbf{W}}) = 1/2 \text{Tr}\{(\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T})^T(\widetilde{\mathbf{X}}\widetilde{\mathbf{W}} - \mathbf{T})\}$$
- $$\widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T}$$

4. fisher's linear discriminant



$y \geq -w_0$ as class \mathcal{C}_1



belonging to \mathcal{C}_1 if $y(\mathbf{x}) \geq y_0$

from the view of dimensionality reduction $y \geq -w_0$ as class \mathcal{C}_1

$$m_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} x_n, m_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} x_n \xrightarrow{y=\mathbf{w}^T \mathbf{x}} m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)$$

5. the perceptron algorithm of rosenblatt

3.3 probabilistic generative models

A probabilistic view of classification from simple assumptions about the distribution of the data

$$\begin{aligned}
p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\
&= \frac{1}{1 + \exp(-a)} = \sigma(a)
\end{aligned}$$

where

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

and $\sigma(a)$ is the **logistic sigmoid** function defined by

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

and $\sigma(-a) = 1 - \sigma(a)$, its inverse is **logit** function

$$a = \ln\left(\frac{\sigma}{1 - \sigma}\right)$$

For case of $K > 2$ classes, we have the following **multi-class generalization**

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}, a_k = \ln [p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)]$$

The **normalized exponential** is known as the **softmax function** as it represents a *smoothed version of the max function*

$$\text{if } a_k \ll a_j, \forall j \neq k, \text{ then } p(\mathcal{C}_k|\mathbf{x}) \approx 1, p(\mathcal{C}_j|\mathbf{x}) \approx 0$$

For **continuous inputs**, assume

$$p(\mathbf{x}|\mathcal{C}_k) = \frac{1}{(2\pi)^{D/2}|\mathbf{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

1. 2 classes

$$\begin{aligned}
p(\mathcal{C}_1|\mathbf{x}) &= \sigma(\mathbf{w}^T \mathbf{x} + w_0) \\
\mathbf{w} &= \mathbf{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\
w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T \mathbf{\Sigma}^{-1}\boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \mathbf{\Sigma}^{-1}\boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}
\end{aligned}$$

2. K classes

$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(\mathcal{C}_k)$$

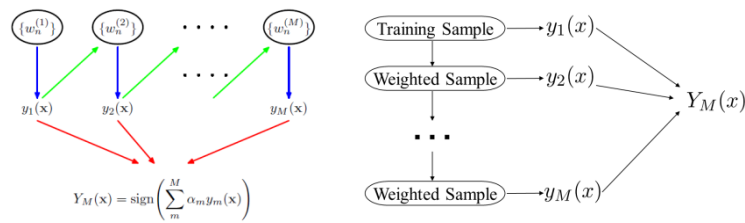
3.4 probabilistic discriminative models

3.5 Boosting

Originally designed for classification problems.

Motivation: a procedure that combines the outputs of many "weak" classifiers to produce a strong/accurate classifier

3.5.1 AdaBoost



4 wef

4.1 wfe

K-means