

Naive Bayes – model is a generative model, so it estimates priors and conditional densities.

Linear regression is a supervised learning technique.

Logistic Regression – 2-class logistic regression is a linear classifier. Logistic regression requires the use of gradient ascent. Logistic regression the parameter η is called the learning rate. It control the speed at which changes happen to the w parameters.

Naïve bayes is generative classifier model learns a model of the point probability $P(x,y)$ logistic regression is Discriminative. Classifier models the posterior $p(y | x)$ directly.

Discriminative model – $P(y | x)$ take the form of a logistic sigmoid function – Call a logistic regression. Logistic regression use the logistic function for modeling $P(y | x)$, considering only the case $y \in \{0,1\}$. The logistic function

$$\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t}$$

$P(y = 0 | x) = 1 / (1 + e^{w^T x}) = 1 - \sigma(w^T x)$
 $P(y = 1 | x) = e^{w^T x} / (1 + e^{w^T x}) = \sigma(w^T x)$
 $P(y=0|x) \geq P(y=1|x)$
 $1 \geq e^{w^T x}$
 $0 \geq w^T x \Leftrightarrow$ if $w^T x$ we classify x as 1
This is a linear classifier $w^T x$ (linear function)
 $P(y|x) = \sigma(w^T x)^y (1 - \sigma(w^T x))^{1-y}$

w^1 and w^2 two sets of parameters, whichever giving a larger $P(y|x)$ should be a better parameter.

Call this $L(w)$ the conditional likelihood. The conditional log likelihood

We can use a commonly-used optimization technique, gradient descent/ascent, to find the solution. gradient ascent

The algorithm
Iterate until converge
 $w^{(k+1)} = w^{(k)} + \eta \nabla_{w^{(k)}} L(w)$
 $\eta > 0$ is a constant called the learning rate.

Linear Machine – Part 1 Basics
SVM is a time of linear machine.

Linear Classifier - $w^T x \leq 0$, Class 0, $w^T x > 0$, Class 1 - $g(x) = w^T x$ is called discriminant function. Is a linear classifier. The learning task is to use the training samples to estimate the parameters of the classifier.

Linear discriminant functions give arise to linear decision boundaries.

If we can find at least one vector w such that $g(x) = w^T x$ classifies all samples.

We say the samples are linearly separable. Solving for the Weight Vector Theoretical: Lagrange or Karush-Kuhn-Tucker In practice: eg. Gradient-descent-based search.

Margins for linear classifiers. The normal vector of the decision line/plane is w . $g(x) = w^T x + w_0$, $g(x) = 0$

$g(x) = 0$ be a decision plane.

$g(x)$ gives an algebraic measure of the distance from x to the decision plane.
 $r = g(x) / ||w||$

for a given set of samples S , the margin is the smallest margin over all $x \in S$. This study source was downloaded by 10000082946624 from CourseHero.com on 04-08-2022 11:46:43 GMT-05:00. For a given set, a classifier that gives rise the larger margin will be better.

A nonlinear (quadratic) optimization problem with linear inequality constraints. Reformulate the problem using lagrange multipliers α . Lagrangian Primal or Dual Problem,

Lagrangian Primal

Support Vector Machine (SVM) – Discriminate Classifier formally define by a separating hyperplane.

Support vectors are the data points that lie closest to the decision surface. Goal of SVM, maximize the margin.

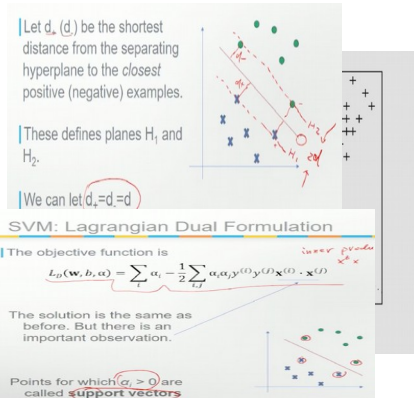
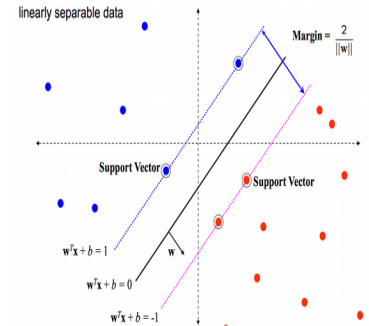
The final w is given by

$$w = \sum_i \alpha_i y^{(i)} x^{(i)}$$

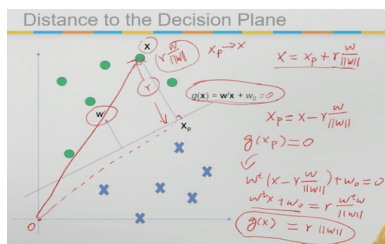
and b is given by

$$y^{(k)} - w^T x^{(k)} \text{ for any } k \text{ such that } \alpha_k > 0$$

The region between two separate planes $H_1 : w^T x + b = 1$ $H_2 : w^T x + b = -1$ as the margin the width is $2 / ||w||$. When the data has noise, there is a problem in drawing a clear hyperplane without misclassifying Distance = $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$ C is a parameter to control how much penalty is assigned to errors.



We will use both notations:
 $g(x) = w^T x$ or $g(x) = w^T x + w_0$
 $w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix}$, $x = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix}$
What is the decision rule for the classifier?
 $g(x) \geq g_j(x), \forall j \neq i$
classify x as class i



Revisit the Lagrange Dual Formulation for SVM
 $L_D(w, b, \alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} \cdot x^{(j)}$
Introduce a kernel function
 $L_D(w, b, \alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y^{(i)} y^{(j)} K(x^{(i)}, x^{(j)})$
Common Kernel Functions
Polynomials of degree d : $K(x^{(i)}, x^{(j)}) = (x^{(i)T} \cdot x^{(j)})^d$
Polynomials of degree up to d : $K(x^{(i)}, x^{(j)}) = ((x^{(i)T} \cdot x^{(j)} + 1)^d)$
Gaussian kernels: $K(x^{(i)}, x^{(j)}) = \exp(-\frac{\|x^{(i)} - x^{(j)}\|^2}{2\sigma^2})$
Sigmoid kernel: $K(x^{(i)}, x^{(j)}) = \tanh(\eta(x^{(i)T} \cdot x^{(j)} + \nu))$

The Kernel Independence of Events
Let (Ω, \mathcal{F}, P) be a probability space, $\forall A, B \in \mathcal{F}$, we say A and B are independent if $P(A \cap B) = P(A)P(B)$.
Mercer's Theorem
If there exists a kernel K such that $K(x, y) = \int \phi(x) \phi(y) d\rho$ where ϕ is a function and ρ is a measure.
Using a kernel K to define a new inner product $\langle \phi(x), \phi(y) \rangle = K(x, y)$ where linear boundaries in the original space can be highly non-linear.
Maximize the margin and also minimize the error

For Large datasets, logistic perform better than naïve bayes. Naïve Bayes Converges to its asymptotic estimates faster than logistic regression. Generative classifier learns a model of the joint probability $p(x,y)$. models the posterior $p(y|x)$ directly.

If the true value of μ is known, then the MLE estimator of σ^2 is $\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$
Is the estimation of σ^2 unbiased?
No
Yes

Basic Machine Learning Paradigms.

Numerical - Values or observations that can be measures

 \mathcal{B} , that satisfies

- $P(A) \geq 0$ for all $A \in \mathcal{B}$
- $P(\Omega) = 1$
- If $A_1, A_2, \dots \in \mathcal{B}$ are pairwise disjoint then
 $P(\cup A_i) = \sum P(A_i)$ (i.e., $A_i \cap A_j = \emptyset, \forall i \neq j$)

Conditional Probability

$$P(B|H) = P(BH) / P(H)$$

and call $P(B|H)$ the conditional probability of B , given H .

- | Inherent ambiguity of many real-world problems

Let (Ω, \mathcal{B}, P) be a probability space, and let $\{H_j\}$ be pairwise disjoint events in \mathcal{B} (i.e., $H_j H_k = \emptyset \forall j \neq k$) and $\bigcup_{j=1, \dots, \infty} H_j = \Omega$. Suppose $P(H_j) > 0, \forall j$, then

$$P(B) = \sum_{j=1, \dots, n} P(H_j)P(B|H_j)$$

-- Such $\{H_j\}$ is called a partition of Ω .

