# CS229: Machine Learning

## Junchuan Zhao

## May 23, 2022

**Abstract**

This document includes my notes of CS229: Machine Learning.

# 1 Lecture 2

## 1.1 Linear Regression

Hypothesis: $h(x) = \sum_{j=0}^{n} \theta_j x_j$ ($n$: number of features)

Linear Regression: the hypothesis is the linear combination of the training dataset.

Loss function (goal): $\min_{\theta} \frac{1}{2} \sum_{i=1}^{m} (h_\theta(x^i) - y^i)^2$

Parameters: $\theta$, $m$, $n$, $x$, $y$

## 1.2 Gradient Descent

### 1.2.1 Batch Gradient Descent

Basic Ideas: start with some $\theta$, keep changing $\theta$ to reduce $J(\theta)$, repeat until convergent

$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ ($\alpha$: learning rate)

$\frac{\partial}{\partial \theta_j} J(\theta) = (h_\theta(x) - y) \cdot x_j$

$\theta_j := \theta_j - \alpha \sum_{i=1}^{m} (h_\theta(x)^i - y^i) \cdot x_j^i$ ($\alpha$: learning rate)

Batch Gradient Descent: all the training data as a batch when optimizing the loss function. (-): not good for large scale dataset, very expensive.

### 1.2.2 Stochastic Gradient Descent

---
**Algorithm 1** Stochastic Gradient Descent

---
**Require:** the parameter of $j^{th}$ feature
**Ensure:** the updated parameter of $j^{th}$ feature
 1: **for** $i = 1$ to $m$ **do**
 2:     $\theta_j := \theta_j - \alpha \cdot (h_\theta(x)^i - y^i) \cdot x_j^i$;
 3: **end for**

---

stochastic gradient descent heads over to the global optimum.

### 1.2.3 Total Equations

A Total Equation for Stochastic Gradient Descent: $\nabla_\theta J(\theta) = \mathbf{0}$

If A is square ($A \in \mathbb{R}_{n \times n}$)

   Trace of A: $tr(A) = \sum_i A_{ii}$

Features of Trace:

- $tr(A) = tr(A^T)$
- $f(A) = tr(AB)$, $\nabla_A f(A) = B^T$
- $tr(AB) = tr(BA)$
- $tr(ABC) = tr(CAB)$
- $\nabla_A tr(AA^T C) = CA + C^T A$

Loss Function: $\nabla_\theta J(\theta) = \frac{1}{2}(X\theta - y)^T(X\theta - y)$
$= \frac{1}{2}(\theta^T X^T - y^T)(X\theta - y)$
$= \frac{1}{2}(\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta)$
$= X^T X\theta - X^T y$
Loss Function: $X^T X\theta - X^T y = \mathbf{0} \iff X^T X\theta = X^T y$
$\theta = (X^T X)^{-1} X^T y$

# 2 Lecture 3

## 2.1 Locally Weighted Regression

"Parametric" learning algorithm: fit fixed set of parameters ($\theta_i$) to data.
"Non-parametric learning algorithm": amount of data/parameters you need to keep grows (linearly) with the size of data.

Linear Regression: to evaluate $h$ at certain $x$
Fit $\theta$ to minimize
   $\frac{1}{2}\sum_i (y^i - \theta^T x^i)^2$,
return $\theta^T x$

Linear Regression: to evaluate $h$ at local region of $x$
Fit $\theta$ to minimize
   $\sum_i w^i (y^i - \theta^T x^i)^2$, where $w^i$ is a "weight function".
common choice for $w^i$ is $w^i = e^{(-\frac{(x^i - x)^2}{2\tau^2})}$
If $|x^i - x|$ is small, then $w^i \approx 1$.
If $|x^i - x|$ is large, then $w^i \approx 0$.
$\tau$: bandwidth, control a larger or narrower window.

## 2.2 Why Square Error?

Assume $y^i = \theta^T x^i + \epsilon^i$, where $\epsilon^i \sim N(0, \sigma^2)$ models effects of random noise
$p(y^i | x^i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y^i - \theta^T x^i)^2}{2\sigma^2}}$
$\iff y^i | x^i; \theta \sim N(\theta^T x^i, \sigma^2)$

Difference between Likelihood ($L$) and Probability ($P$): $L$ variess parameters, $P$ varies datapoints. Assume the datapoints are IID,

The "likelihood" of $\theta$: $L(\theta) = p(\boldsymbol{y}|\boldsymbol{x};\theta) = \prod_{i=1}^{m} p(y^i|x^i;\theta)$

$l(\theta) = \log \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi}\sigma} e^{(\cdots)} = \sum_{i=1}^{m} [\log \frac{1}{\sqrt{2\pi}\sigma} + \log e^{(\cdots)}]$

$= m \log \frac{1}{\sqrt{2\pi}\sigma} - \sum_{i=1}^{m} \frac{(y^i - \theta^T x^i)^2}{2\sigma^2}$

MLE: maximum likelihood estimation. Choose $\theta$ to maximize $L(\theta)$

i.e. choose $\theta$ to minimize $\frac{1}{2} \sum_{i=1}^{m} (y^i - \theta^T x^i)^2$, which is actually $J(\theta)$

## 2.3 Classification

Binary Classification: $y \in \{0, 1\}$

### 2.3.1 Logistic Regression

We want $h_\theta(x) \in [0, 1]$, so we define
$h_\theta(x) = g(\theta^T x) = \frac{1}{1+e^{-\theta^T x}}$
Sigmoid Function: $g(z) = \frac{1}{1+e^{-z}}$
An important thing is that: the partial derivative of the sigmoid function is a **concave** function, which means it has the global maximum.
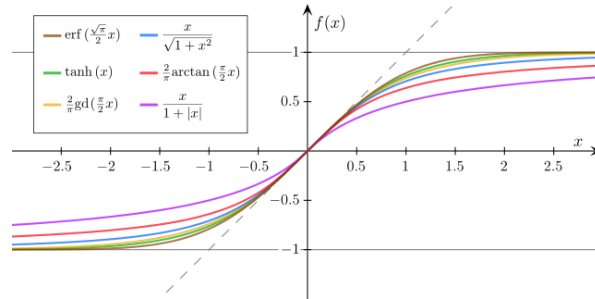


Figure 1: The graph of sigmoid function.

Different from Linear Regression $h_\theta(x) = \theta^T x$, Logistic Regression is $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$.
$p(y = 1|x;\theta) = h_\theta(x)$, $p(y = 0|x;\theta) = 1 - h_\theta(x)$
In sum, $p(y|x;\theta) = h_\theta(x)^y (1 - h_\theta(x))^{1-y}$

$L(\theta) = p(\boldsymbol{y}|\boldsymbol{x};\theta) = \prod_{i=1}^{m} h_\theta(x^i)^{y^i} (1 - h_\theta(x^i))^{1-y^i}$

$l(\theta) = \log(L(\theta)) = \sum_{i=1}^{m} y^i \log h_\theta(x^i) + (1 - y^i) \log(1 - h_\theta(x^i))$

Choose $\theta$ to maximize $l(\theta)$, we use Batch gradient ascent.

Batch Gradient Ascent:

$\theta_j := \theta_j + \alpha \frac{\partial}{\partial \theta_j} l(\theta)$

$\theta_j := \theta_j + \alpha \sum_{i=1}^{m} (y^i - h_\theta(x^i)) x_j^i$, same as the one in linear regression. (A larger category called generalized linear model (GLM))

Note: no **normal equation** for logistic regression solution.

### 2.3.2 Newton's Method

The basic idea: have some function $f$, we want to fit $\theta$, **s.t.** $f(\theta) = 0 \iff$ want maximizing $l(\theta)$ $\iff$ want $l'(\theta) = 0$
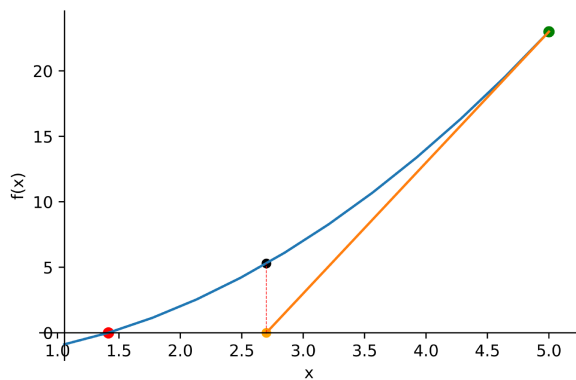


Figure 2: Newton's Method.

$\theta^{t+1} := \theta^t - \frac{f(\theta^t)}{f'(\theta^t)}$, let $f(\theta) = l'(\theta)$

$\theta^{t+1} := \theta^t - \frac{l'(\theta^t)}{l''(\theta^t)}$

"**Quadratic convergent**": $0.01$ error $\longrightarrow 0.0001$ error $\longrightarrow 0.00000001$ error.

When $\theta$ is a vector: $(\theta \in \mathbb{R}^{n+1})$

$\quad \theta^{t+1} := \theta^t + \alpha H^{-1} \nabla_\theta l(\theta)$

where

$$\nabla_\theta l(\theta) = \begin{bmatrix} \frac{\partial l}{\partial \theta_{00}} & \frac{\partial l}{\partial \theta_{01}} & \cdots & \frac{\partial l}{\partial \theta_{0n}} \\ \frac{\partial l}{\partial \theta_{10}} & \frac{\partial l}{\partial \theta_{11}} & \cdots & \frac{\partial l}{\partial \theta_{1n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial l}{\partial \theta_{n0}} & \frac{\partial l}{\partial \theta_{n1}} & \cdots & \frac{\partial l}{\partial \theta_{nn}} \end{bmatrix}$$

$H \in \mathbb{R}^{n+1 \times n+1}$ is the Hessian Matrix, $H_{ij} = \frac{\partial^2 l}{\partial \theta_i \partial \theta_j}$

Note: If the dataset is large, the size of $H$ will be too large to compute.

## 3 Lecture 4

### 3.1 Perceptron

Binary step function:

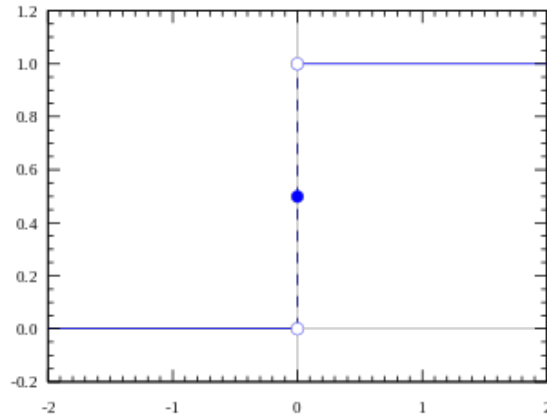$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases} \tag{1}$$

Figure 3: The graph of binary step.

Batch Gradient Ascent: $\theta_j := \theta_j + \alpha(y^i - h_\theta(x^i))x_j^i$

## 3.2  Exponential Family

### 3.2.1  Defination and Examples of Exponential Family

Probability Density Function (PDF): $p(y; \eta) = b(y)\exp[\eta^T T(y) - a(\eta)]$

$y$: data

$\eta$: natural parameter

$T(y)$: sufficient statistic

$b(y)$: base measure

$a(\eta)$: log partition

Example 1: Bernoulli (Binary Data): $p(y; \phi) = \phi^y(1 - \phi)^{(1-y)}$

$\phi =$ probability of the event

$p(y; \phi) = \exp[y \log\left(\frac{\phi}{1-\phi}\right) + \log\left(1 - \phi\right)]$

$b(y) = 1$

$T(y) = y$

$\eta = \log\frac{\phi}{1-\phi} \Rightarrow \phi = \frac{1}{1+e^{-\eta}}$

$a(\eta) = -\log\left(1 - \phi\right) = -\log\left(1 - \frac{1}{1+e^{-\eta}}\right) = \log\left(1 + e^\eta\right)$

Example 2: Gaussian (assume variance $= 1$): $p(y; \mu) = \frac{1}{\sqrt{2\pi}}e^{(-\frac{(y-\mu)^2}{2})} = \frac{1}{\sqrt{2\pi}}e^{-\frac{y^2}{2}}\exp(\mu y - \frac{\mu^2}{2})$

$b(y) = \frac{1}{\sqrt{2\pi}}e^{-\frac{y^2}{2}}$

$T(y) = y$

$\eta = \mu$

$a(\eta) = \frac{\mu^2}{2} = \frac{\eta^2}{2}$

### 3.2.2  Properties of Exponential Family

- MLE w.r.t $\eta$ is concave, NLL is convex
- $E[y; \eta] = \frac{\partial}{\partial\eta}a(\eta)$
- $Var[y; \eta] = \frac{\partial^2}{\partial\eta^2}a(\eta)$

Real - Gaussian; Binary - Bernoulli; Count - Poisson; $\mathbb{R}^+$ - Gamma, Exponential; Distribution - Beta, Dirichlet (Bayesian).

## 3.3 Generalized Linear Model (GLM)

### 3.3.1 GLM Assumptions

Assumptions/Design Choices
(1) $y|x; \theta \sim$ Exponential Family $(\eta)$
(2) $\eta = \theta^T x$, $\theta \in \mathbb{R}^n$ and $x \in \mathbb{R}^n$
(3) Test Time: output $h_\theta(x) = E[y|x; \theta]$
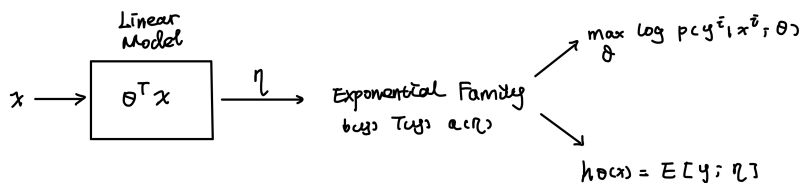
### 3.3.2 GLM Pipeline



Figure 4: The pipeline for GLM.

### 3.3.3 GLM Training

Learning Update Rule: $\theta_j := \theta_j + \alpha(y^i - h_\theta(x^i))x_j^i$

Terminology:
$\eta$ - natural parameter
$g(\eta) = \mu = E(y; \eta)$ - canonical response function
$\eta = g^{-1}(\mu)$ - canonical link function

Three parameterizations:
Model Param: $\theta$ (**Learn Parameters**)
Natural Param: $\eta$ (Design Choice: $\theta^T x$)
Canonical Param: $\phi$ - Bernoulli, $\mu$ $\sigma$ - Gaussian, $\lambda$ - Poisson (Canonical Response: $g(\cdot)$)

### 3.3.4 Review Logistic Regression

$h_\theta(x) = E[y|x; \theta] = \phi$ (the mean of the Bernoulli distribution)
$h_\theta(x) = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}}$

## 3.4 Softmax Regression

Task: multi-class Classification
$k$ - number of classes, $x^i \in \mathbb{R}^n$, Label $y^i \in \{0, 1\}^k$ (one-hot vector)
Each class has its own parameters: $\theta_{class} \in \mathbb{R}^n$
$p(y = i|x; \theta) = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}}$
$l(\theta) = \sum_{i=1}^m \log \prod_{l=1}^k (\frac{e^{\theta_l^T x^i}}{\sum_{j=1}^k e^{\theta_j^T x^i}})^{\{y^i=l\}}$

# 4 Lecture 5

## 4.1 Discriminative v.s. Generative Learning Algorithms

Discriminative: Learn $p(y|x)$ which is the mapping from x to y.

Generative: Learn $p(x|y)$, $p(y)$ (class prior)

According to Bayes rule: $p(y|x) = \frac{p(x|y)p(y)}{p(x)}$, $p(x) = \sum_{i=0}^{k} p(x|y=i)p(y=i)$

## 4.2 Gaussian Discriminant Analysis (GDA)

### 4.2.1 Assumptions and Fundamental Knowledge

Suppose $x \in \mathbb{R}^n$ (convention: drop $x_0 = 1$)
Assume $p(x|y)$ is Gaussian

Basic Knowledge for Multivariate Gaussian:
$z \sim N(\boldsymbol{\mu}, \Sigma)$, $(z \in \mathbb{R}^n)$
$p(z) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} exp(-\frac{1}{2}(x-\boldsymbol{\mu})^T \Sigma^{-1}(x-\boldsymbol{\mu}))$ $\quad E[z] = \boldsymbol{\mu}$
$Cov(z) = E[(z-\mu)(z-\mu)^T] = E_{zz^T} - E[z]E[z]^T$

### 4.2.2 GDA model

Parameters: $\boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma$ (use same covariance matrix)
$p(x|y=0) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} exp(-\frac{1}{2}(x-\boldsymbol{\mu_0})^T \Sigma^{-1}(x-\boldsymbol{\mu_0}))$
$p(x|y=1) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} exp(-\frac{1}{2}(x-\boldsymbol{\mu_1})^T \Sigma^{-1}(x-\boldsymbol{\mu_1}))$
$p(y) = \phi^y (1-\phi)^y$ $(p(y=1) = \phi)$

Training set: $\{(x^i, y^i)\}_{i=1}^m$
Joint Likelihood: $L(\phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma) = \prod_{i=1}^{m} p(x^i, y^i; \phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}) = \prod_{i=1}^{m} p(x^i|y^i)p(y^i)$
$\quad$ Discriminative Learning (**Conditional Likelihood**): $\prod_{i=1}^{m} p(y^i|x^i; \theta)$

Maximize Likelihood Estimation (MLE):
$\max_{\phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma} l(\phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma)$
Results of MLE:
$\phi = \frac{\sum_{i=1}^{m} y^i}{m} = \frac{\sum_{i=1}^{m} \mathbb{1}\{y^i=1\}}{m}$ $(\mathbb{1} : \mathbb{1}\{\text{true}\} = 1, \mathbb{1}\{\text{false}\} = 0)$
$\boldsymbol{\mu_0} = \frac{\sum_{i=1}^{m} \mathbb{1}\{y^i=0\}x^i}{\sum_{i=1}^{m} \mathbb{1}\{y^i=0\}} = \frac{\text{sum of feature vectors for samples with y=0}}{\text{number of samples with y=0}}$
$\boldsymbol{\mu_0} = \frac{\sum_{i=1}^{m} \mathbb{1}\{y^i=1\}x^i}{\sum_{i=1}^{m} \mathbb{1}\{y^i=1\}}$
$\Sigma = \frac{1}{m} \sum_{i=1}^{m} (x^i - \mu_{y^i})(x^i - \mu_{y^i})^T$

Prediction:
$\arg\max_{y} p(y|x) = \arg\max_{y} \frac{p(x|y)p(y)}{p(x)}$ $(p(x)$ is constant$)$
$= \arg\max p(x|y)p(y)$

7

### 4.2.3 Compare GDA to Logistic Regression

For fix parameters: $\phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma$, plot $p(y|x; \phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma)$ as a function of $x$.

By Bayes Rule: $p(y=1|x; \phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma) = \frac{p(x|y=1; \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma) p(y=1; \phi)}{p(x; \phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma)}$.

In the end, $p(y=1|x; \phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma) = \frac{1}{1+exp(-\theta^T x)}$, where $\theta$ is a function of $\phi, \boldsymbol{\mu_0}, \boldsymbol{\mu_1}, \Sigma$.

The assumption of GDA $\Rightarrow p(y=i|x) = \frac{1}{1+e^{-\theta^T x}}$

However, $p(y=i|x) = \frac{1}{1+e^{-\theta^T x}} \nRightarrow$ GDA assumptions.

| GDA | Logistic Regression |
|---|---|
| stronger assumptions | weaker assumptions |
| perform better when assumptions are correct | fine with all the scenarios (distributions) |
| small dataset | large dataset |
| computationally efficient | fits the era of the big data |

Table 1: The comparison between GDA and Logistic Regression.

## 4.3 Naive Bayes

For a spam email classification problem:

1. create a word dictionary (size $m$).

If represent $x$ by a binary feature vector, $x \in \{0,1\}^n$. ($x_i = \mathbb{1}\{$word i appears in email$\}$)

   This cause $2^n$ possible values of x, which is not possible for large dictionary.

2. Assume $x_i$**'s are conditionally independent given y**, which means

   $p(x_1, ..., x_{10000}|y) = p(x_1|y)p(x_2|x_1, y)p(x_3|x_1, x_2, y)...p(x_{10000}|..., y)$

   $\overset{\text{assume}}{=\joinrel=\joinrel=} p(x_1|y)p(x_2|y)...p(x_{10000}|y) = \prod\limits_{i=1}^{n} p(x_i|y)$

Parameters:

$\phi_{j|y=1} = p(x_j = 1|y = 1)$

$\phi_{j|y=0} = p(x_j = 1|y = 0)$

$\phi_y = p(y = 1)$

Joint Likelihood: $L(\phi_i, \phi_{j|y}) = \prod\limits_{i=1}^{m} p(x^i, y^i; \phi_j, \phi_{j|y})$

MLE results: $\phi_y = \frac{\sum\limits_{i=1}^{m} \mathbb{1}\{y^i=1\}}{m}$

$\phi_{j|y=1} = \frac{\sum\limits_{i=1}^{m} \mathbb{1}\{x_j^i=1, y^i=1\}}{\sum\limits_{i=1}^{m} \mathbb{1}\{y^i=1\}}$

$\phi_{j|y=0} = \frac{\sum\limits_{i=1}^{m} \mathbb{1}\{x_j^i=1, y^i=0\}}{\sum\limits_{i=1}^{m} \mathbb{1}\{y^i=0\}}$

At prediction time:

$p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x|y=1)p(y=1)+p(x|y=0)p(y=0)}$

Problem: when a word does not exist $- \frac{0}{0}$

# 5 Lecture 6

## 5.1 Laplace Smoothing

$x \in \{1, ..., k\}$, estimate $p(x = j) = \frac{\sum\limits_{j=1}^{m} \mathbb{1}\{x^i = j\} + 1}{m + + k}$

$\phi_{j|y=0} = \frac{\sum\limits_{i=1}^{m} \mathbb{1}\{x_j^i = 1, y^i = 0\} + 1}{\sum\limits_{i=1}^{m} \mathbb{1}\{y^i = 0\} + 2}$

## 5.2 Multinomial Event Model

$x \in \begin{bmatrix} 600 \\ 800 \\ 1600 \\ 6200 \end{bmatrix}, \; x \in \mathbb{R}^n$

$x_j \in \{1, 2, ..., 10000\}$, $n = $ length of email, $m = $ number of emails

Parameters: $\phi_y = p(y = 1), \phi_{k|y=0} = p(x_j = k|y = 0), \phi_{k|y=1} = p(x_j = k|y = 1)$

MLE results: $\phi_{k|y=0} = \dfrac{\sum\limits_{i=1}^{m} (\mathbb{1}\{y^i = 0\} \sum\limits_{j=1}^{n_i} \mathbb{1}\{x_j^i = k\})}{\sum\limits_{i=1}^{m} \mathbb{1}\{y^i = 0\} n_i}$

$\underline{\underline{\text{laplace smoothing}}} = \dfrac{\sum\limits_{i=1}^{m} (\mathbb{1}\{y^i = 0\} \sum\limits_{j=1}^{n_i} \mathbb{1}\{x_j^i = k\}) + 1}{\sum\limits_{i=1}^{m} \mathbb{1}\{y^i = 0\} n_i + 10000}$

The advantages for GDA and Naive Bayes: quick to train, non-iterative.

## 5.3 Support Vector Machines SVM

SVM helps to find non-linear boundaries for classification problems.

### 5.3.1 Optimal Margin Classifier

Functional Margin: how confidentaly and accurately to classify an example.

Logistic Regression: $h_\theta(x) = g(\theta^T x)$, predict 1 if $\theta^T x \geq 0$ and 0 otherwise.

We want if $y^i = 1$, hope that $\theta^T x^i \gg 0$ and 0 when $\theta^T x^i \ll 0$. – correct or confident prediction.

Notation: Labels $y \in \{-1, +1\}$,

$$g(z) = \begin{cases} 1 & z \geq 0 \\ -1 & z < 0 \end{cases} \tag{2}$$

$h_{W,b}(x) = g(W^T x + b)$, $W \in \mathbb{R}^n$ and $b \in \mathbb{R}$. (Compare to Logistic Regression, $b$ is $\theta_0$)

Functional Margin of hyperplane defined by $(W, b)$ w.r.t $(x^i, y^i)$:

$\hat{\gamma}^i = y^i(W^T x^i + b)$, we want $\gamma^i \gg 0$

Funtional Margine w.r.t training set: $\hat{\gamma} = \min\limits_{i} \hat{\gamma}^i$, $i = 1, ..., m$

Rescaling the parameters: $(W, b) \rightarrow (\frac{W}{\|W\|}, \frac{b}{\|W\|})$

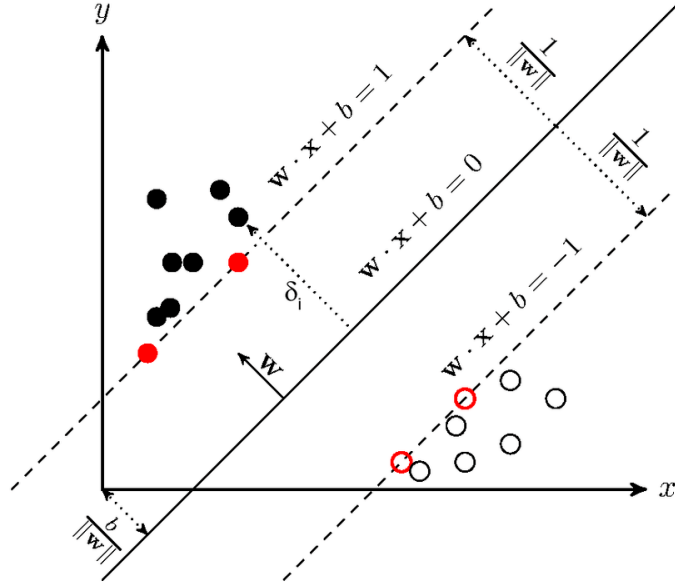Geometric Margin: the distance from the data point to the hyperplane.

Figure 5: Geometric Margin.

Geometric Margin of hyperplane $(W,b)$ w.r.t $(x^i,y^i)$:

$\gamma^i = \frac{y^i(W^T x^i + b)}{\|W\|}$

Geometric Margin w.r.t training set: $\gamma = \min_i \gamma^i$

The relationship between geometric and functional margin: $\gamma^i = \frac{\hat{\gamma}^i}{\|W\|}$

$\hat{\gamma}$ is functional margin; $\gamma$ is geometric margin.

Optimal margin classifier: choose $w$, $b$ to maximize $\gamma$:

$\max_{\gamma,w,b} \gamma$ s.t. $\frac{y^i(W^T x^i + b)}{\|W\|} \geq \gamma$, $(i = 1, ..., m.)$ $\Rightarrow \min_{w,b} \|W\|^2$ s.t. $y^i(W^T x^i + b) \geq 1$