

CS229 Machine Learning Notes and Formulae

Chang He

July 2019

Contents

1	Supervised Learning and Linear Models	2
1.1	Linear Models	2
1.1.1	Linear Regression	2
1.1.1.1	Ordinary Least Square	2
1.1.1.2	Probabilistic interpretation of Linear Model	2
1.1.2	Linear Models for Classification	3
1.1.2.1	Using Linear Regression for Classification	3
1.1.2.2	Logistic Regression	3
1.2	Newton's Method	4
1.3	Generalized Linear Model	4
1.3.1	Exponential Family Distribution	4
1.3.2	Generalized Linear Models	5
1.3.2.1	Multinomial Distribution and Softmax Regression	6
2	Generative Learning Algorithms	7
2.1	Multivariate Normal Distribution	8
2.1.1	Gaussian Discriminant Analysis	8
2.2	Naive Bayes	10
3	Support Vector Machine	11
3.1	Preliminaries	11
3.1.1	Lagrange Multiplier	11
3.1.2	Distance of a Point to a Hyperplane	12
3.1.3	Lagrange duality	12
3.2	Optimal Margin Classifier	13
3.2.1	Lagrange Dual Form	14
3.2.2	Kernel	14
3.2.3	Soft Margin Classifier	15
3.2.4	Sequential Minimal Optimization	16
3.2.5	Multi-class SVM (Crammer and Singer)	17

1 Supervised Learning and Linear Models

1.1 Linear Models

1.1.1 Linear Regression

Given data set X with n data points and d features, with the associated targets y , minimize the loss function:

$$\frac{1}{2} \sum_{i=1}^n (\theta^T X_{(i)} + b - y_{(i)})^2$$

By changing the definition of X and θ to include the intercept term b , with the 0th feature of X being 1 and 0th term of θ being b , we have:

$$\frac{1}{2} \sum_{i=1}^n (\theta^T X_{(i)} - y_{(i)})^2$$

Vectorizing the loss function, we have the objective function:

$$\min_{\theta} l(\theta) = \min_{\theta} \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

1.1.1.1 Ordinary Least Square To find the potential minimum of this function, we find its gradient first:

$$\begin{aligned} \nabla_{\theta} l(\theta) &= \nabla_{\theta} \frac{1}{2} (X\theta - y)^T (X\theta - y) \\ &= \nabla_{\theta} \frac{1}{2} (\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta + y^T y) \\ &= \frac{1}{2} \nabla_{\theta} \text{tr} (\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta + y^T y) && \text{tr } x = x \text{ if } x \in \mathcal{R} \\ &= \frac{1}{2} \nabla_{\theta} (\text{tr } \theta^T X^T X\theta - \text{tr } \theta^T X^T y - \text{tr } y^T X\theta) && \text{remove the constant} \\ &= \frac{1}{2} \nabla_{\theta} (\text{tr } \theta^T X^T X\theta - \text{tr } y^T X\theta - \text{tr } y^T X\theta) && \text{tr } A^T = \text{tr } A \\ &= \frac{1}{2} (\nabla_{\theta^T} \text{tr } \theta^T X^T X\theta)^T - \nabla_{\theta} \text{tr } y^T X\theta && \nabla_A f(A) = (\nabla_{A^T} f(A))^T \\ &= \frac{1}{2} (\theta^T X^T X + \theta^T X^T X) - y^T X && \nabla_A ABA^T C = CAB + C^T AB^T \\ & && \text{with } A = \theta^T, B = X^T X, C = I \\ &= (\theta^T X^T X + y^T X)X \\ &= (X\theta - y)^T X \\ &= X^T (X\theta - y) \end{aligned}$$

To obtain the critical point, set the gradient to 0 and solve the linear equation:

$$X^T X\theta = X^T y$$

The result is denoted as $\theta = (X^T X)^{-1} X^T y$.

1.1.1.2 Probabilistic interpretation of Linear Model Given data set X with n data points and d features, with the associated targets y , assuming that the data set is modeled by parameter θ and the residual $(y_{(i)} - \theta^T X_{(i)})$ follows a Gaussian distribution, we have:

$$\begin{aligned} y_{(i)} &= \theta^T X_{(i)} + \epsilon_{(i)} && \epsilon_{(i)} \sim \mathcal{N}(0, \sigma^2) \\ P(\epsilon_{(i)}) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon_{(i)})^2}{2\sigma^2}\right) \\ P(y_{(i)}|X_{(i)}; \theta) &= \mathcal{N}(y_{(i)} - \theta^T X_{(i)}, \sigma^2) \\ &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(y_{(i)} - \theta^T X_{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

Assuming $\epsilon_{(i)}$ is Independently and Identically Distributed (IDD), we have:

$$\begin{aligned} P(\mathbf{y}|X; \boldsymbol{\theta}) &= \prod_{i=1}^n P(\mathbf{y}_{(i)}|X_{(i)}; \boldsymbol{\theta}) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mathbf{y}_{(i)} - \boldsymbol{\theta}^T X_{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

Let find the maximum likelihood estimate of $\boldsymbol{\theta}$ by maximizing $P(\mathbf{y}|X; \boldsymbol{\theta})$:

$$\begin{aligned} \max_{\boldsymbol{\theta}} P(\mathbf{y}|X; \boldsymbol{\theta}) &= \max_{\boldsymbol{\theta}} \ln P(\mathbf{y}|X; \boldsymbol{\theta}) && \text{since } \ln \text{ is monotonously increasing} \\ &= \max_{\boldsymbol{\theta}} \sum_{i=1}^n \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mathbf{y}_{(i)} - \boldsymbol{\theta}^T X_{(i)})^2}{2\sigma^2}\right) \right) \\ &= \max_{\boldsymbol{\theta}} \sum_{i=1}^n \left(-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(\mathbf{y}_{(i)} - \boldsymbol{\theta}^T X_{(i)})^2}{2\sigma^2} \right) \\ &= \min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(\frac{(\mathbf{y}_{(i)} - \boldsymbol{\theta}^T X_{(i)})^2}{2\sigma^2} \right) && \text{remove the constant term} \\ &= \min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_{(i)} - \boldsymbol{\theta}^T X_{(i)})^2 \\ &= \min_{\boldsymbol{\theta}} \frac{1}{2} (X\boldsymbol{\theta} - \mathbf{y})^T (X\boldsymbol{\theta} - \mathbf{y}) \\ &= \min_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) \end{aligned}$$

1.1.2 Linear Models for Classification

1.1.2.1 Using Linear Regression for Classification Since the component range of hypothesis function $h_{\boldsymbol{\theta}}(X) = X\boldsymbol{\theta}$ lies far outside of the relevant range for classification $[0, 1]$, linear regression is extremely sensitive to data points with a large norm, so it is generally not recommended to use linear regression for classification.

1.1.2.2 Logistic Regression The hypothesis for logistic regression is:

$$h_{\boldsymbol{\theta}}(X) = S(\boldsymbol{\theta}^T X) = \frac{1}{1 + \exp(-X\boldsymbol{\theta})}$$

, where function S is called sigmoid or logistic function.

Each component of the hypothesis, $h_{\boldsymbol{\theta}}(X_{(i)})$, represents the probability of the associated target $\mathbf{y}_{(i)}$ equal to 1, that is:

$$\begin{aligned} h_{\boldsymbol{\theta}}(X_{(i)}) &= P(\mathbf{y}_{(i)} = 1|X_{(i)}; \boldsymbol{\theta}) \\ 1 - h_{\boldsymbol{\theta}}(X_{(i)}) &= P(\mathbf{y}_{(i)} = 0|X_{(i)}; \boldsymbol{\theta}) \end{aligned}$$

Since $y \in \{0, 1\}$,

$$P(\mathbf{y}_{(i)}|X_{(i)}; \boldsymbol{\theta}) = h_{\boldsymbol{\theta}}(X_{(i)})^y (1 - h_{\boldsymbol{\theta}}(X_{(i)}))^{(1-y)}$$

Using MLE estimation, with IDD assumption, we get:

$$\begin{aligned} &\max_{\boldsymbol{\theta}} \ln P(\mathbf{y}|X; \boldsymbol{\theta}) \\ &= \max_{\boldsymbol{\theta}} \sum_{i=1}^n \ln P(\mathbf{y}_{(i)}|X_{(i)}; \boldsymbol{\theta}) \\ &= \max_{\boldsymbol{\theta}} \sum_{i=1}^n \ln (h_{\boldsymbol{\theta}}(X_{(i)})^y (1 - h_{\boldsymbol{\theta}}(X_{(i)}))^{(1-y)}) \\ &= \max_{\boldsymbol{\theta}} \sum_{i=1}^n (y \ln(h_{\boldsymbol{\theta}}(X_{(i)})) + (1 - y_{(i)}) \ln(1 - h_{\boldsymbol{\theta}}(X_{(i)}))) \end{aligned}$$

Given that $1 - S(x) = \frac{\exp(-x)}{1 + \exp(-x)} = \exp(-x)S(x)$

$$\begin{aligned}
&= \max_{\boldsymbol{\theta}} \sum_{i=1}^n \left(\mathbf{y}_{(i)} \ln(h_{\boldsymbol{\theta}}(X_{(i)})) + (1 - \mathbf{y}_{(i)})(-\boldsymbol{\theta}^T X_{(i)} + \ln(h_{\boldsymbol{\theta}}(X_{(i)}))) \right) \\
&= \max_{\boldsymbol{\theta}} \sum_{i=1}^n \left(\mathbf{y}_{(i)} \ln(h_{\boldsymbol{\theta}}(X_{(i)})) + (-\boldsymbol{\theta}^T X_{(i)} + \ln(h_{\boldsymbol{\theta}}(X_{(i)}))) + (\boldsymbol{\theta}^T X_{(i)} \mathbf{y}_{(i)} - \mathbf{y}_{(i)} \ln(h_{\boldsymbol{\theta}}(X_{(i)}))) \right) \\
&= \max_{\boldsymbol{\theta}} \sum_{i=1}^n \left(\ln(h_{\boldsymbol{\theta}}(X_{(i)})) - (1 - \mathbf{y}_{(i)})\boldsymbol{\theta}^T X_{(i)} \right) \\
&= \min_{\boldsymbol{\theta}} \sum_{i=1}^n \left(\ln(1 + \exp(-\boldsymbol{\theta}^T X_{(i)})) + (1 - \mathbf{y}_{(i)})\boldsymbol{\theta}^T X_{(i)} \right)
\end{aligned}$$

The gradient of which is:

$$\begin{aligned}
&\nabla_{\boldsymbol{\theta}} \sum_{i=1}^n \left(\ln(1 + \exp(-\boldsymbol{\theta}^T X_{(i)})) + (1 - \mathbf{y}_{(i)})\boldsymbol{\theta}^T X_{(i)} \right) \\
&= \sum_{i=1}^n \left(-\frac{1}{1 + \exp(\boldsymbol{\theta}^T X_{(i)})} \circ X_{(i)} + (1 - \mathbf{y}_{(i)})X_{(i)} \right) \\
&= \sum_{i=1}^n \left(1 - \mathbf{y}_{(i)} - \frac{1}{1 + \exp(\boldsymbol{\theta}^T X_{(i)})} \right) \circ X_{(i)} \\
&= \sum_{i=1}^n \left(\frac{1}{1 + \exp(-\boldsymbol{\theta}^T X_{(i)})} - \mathbf{y}_{(i)} \right) \circ X_{(i)} \\
&= (h(X_{(i)}) - \mathbf{y}_{(i)}) \circ X_{(i)}
\end{aligned}$$

1.2 Newton's Method

$$\boldsymbol{\theta} := \boldsymbol{\theta} - (\nabla_{\boldsymbol{\theta}}^2 l(\boldsymbol{\theta}))^{-1} \nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta})$$

Newton's Method enjoys quadratic convergence, meaning that the error of the result will decrease quadratically. However, it is more expensive, both time and memory-wise, to calculate and invert Hessian in each iteration compared to gradient descent. For data with a large number of dimensions, using Newton's Method is problematic also because of proliferation of saddle points, which Newton's Method tends to attract to.

Despite so, for linear models with mean-square loss, Newton's Method converges extremely fast. For a perfect square form loss function, Newton's Method converges in an exactly single iteration.

1.3 Generalized Linear Model

1.3.1 Exponential Family Distribution

Exponential Family Distribution is a class of distributions parameterized by η :

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

Both Bernoulli and Gaussian distribution are special cases of the exponential family distribution. For Bernoulli distribution, let's define:

$$\begin{aligned}
p(y; \phi) &= \phi^y (1 - \phi)^{(1-y)} \\
&= \exp(y \ln \phi + (1 - y) \ln(1 - \phi)) \\
&= \exp(y \ln \frac{\phi}{1 - \phi} + \ln(1 - \phi))
\end{aligned}$$

In this case, we have:

$$b(y) = 1$$

$$\begin{aligned}
\eta &= \frac{\phi}{1 - \phi} \\
T(y) &= y \\
a(\eta) &= -\ln(1 + \exp(\eta)) \\
p(y; \eta) &= \exp(\eta y - \ln(1 + \exp(\eta))) \\
&= \frac{e^{\eta y}}{1 + e^{\eta}} = \text{sigmoid}(\eta) \cdot \exp(\eta(y - 1))
\end{aligned}$$

For Gaussian distribution, we define:

$$\begin{aligned}
p(y; \eta) &= b(y) \exp(\boldsymbol{\eta}^T T(y) - a(\boldsymbol{\eta})) \\
p(y; (\mu, \sigma)) &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right) \\
&= \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(\frac{2y\mu - y^2}{2\sigma^2}\right) \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \\
&= \frac{1}{\sqrt{2\pi}} \sigma^{-1} \exp\left(y \cdot \frac{\mu}{\sigma^2} - y^2 \frac{1}{2\sigma^2} - \frac{\mu^2}{2\sigma^2}\right) \\
&= \frac{1}{\sqrt{2\pi}} \exp\left(y \cdot \frac{\mu}{\sigma^2} - y^2 \frac{1}{2\sigma^2} - \frac{\mu^2}{2\sigma^2} - \ln \sigma\right) \\
b(y) &= \frac{1}{\sqrt{2\pi}} \\
\boldsymbol{\eta} &= \begin{bmatrix} \frac{\mu}{\sigma^2} \\ 1 \\ -\frac{1}{2\sigma^2} \end{bmatrix} \\
T(y) &= \begin{bmatrix} y \\ y^2 \end{bmatrix} \\
a(\boldsymbol{\eta}) &= \frac{\mu^2}{2\sigma^2} + \ln \sigma = -\frac{\boldsymbol{\eta}_{(1)}^2}{4\boldsymbol{\eta}_{(2)}} + \frac{1}{2} \ln \left| \frac{1}{2\boldsymbol{\eta}_{(2)}} \right| \\
p(y; \boldsymbol{\eta}) &= \frac{1}{\sqrt{2\pi}} \exp\left(\boldsymbol{\eta}^T \begin{bmatrix} y \\ y^2 \end{bmatrix} + \frac{\boldsymbol{\eta}_{(1)}^2}{4\boldsymbol{\eta}_{(2)}} - \frac{1}{2} \ln \left| \frac{1}{2\boldsymbol{\eta}_{(2)}} \right| \right)
\end{aligned}$$

1.3.2 Generalized Linear Models

We make following assumptions to define GLMs:

1. $p(\mathbf{y}_{(i)} | X_{(i)}; \boldsymbol{\theta}) = \text{Exp}(\boldsymbol{\eta})$
2. $h(X_{(i)}) = \mathbb{E}[T(\mathbf{y}_{(i)}) | X_{(i)}]$
3. $\boldsymbol{\eta} = \boldsymbol{\theta}^T X_{(i)}$ or $\boldsymbol{\eta}_{(j)} = \boldsymbol{\Theta}_j^T X_{(i)}$

Using the parameters for Bernoulli distribution, we can see that:

$$\begin{aligned}
h(X_{(i)}) &= p(\mathbf{y}_{(i)} = 1 | X_{(i)}) \\
&= \text{sigmoid}(\eta) \\
&= \frac{1}{1 + \exp(\boldsymbol{\theta}^T X_{(i)})}
\end{aligned}$$

and for the parameters of Gaussian distribution, the hypothesis is:

$$\begin{aligned}
h(X_{(i)}) &= \mathbb{E}[T(\mathbf{y}_{(i)}) | X_{(i)}] \\
&= \mu \\
&= \boldsymbol{\Theta}_1^T X_{(i)}
\end{aligned}$$

1.3.2.1 Multinomial Distribution and Softmax Regression Similar to Bernoulli distribution, multinomial distribution is used to represent the discrete probability distribution of a random variable $y \in 0, \dots, k$:

$$\begin{aligned} P(y = i) &= \begin{cases} \phi_i, & i > 0 \\ 1 - \sum_{i=1}^k \phi_i, & i = 0 \end{cases} \\ &= \left(1 - \sum_{i=1}^k \phi_i\right)^{1\{y=0\}} \prod_{i=1}^k \phi_i^{1\{y=i\}} \end{aligned}$$

Given an identity matrix $I \in \mathcal{R}^{k \times k}$, we define $T(y)$ as:

$$\begin{aligned} T(y) &= \begin{cases} I_y, & y > 0 \\ \mathbf{0}, & y = 0 \end{cases} \\ P(y = i) &= \left(1 - \sum_{i=1}^k \phi_i\right)^{(1 - \mathbf{1}^T T(y))} \prod_{i=1}^k \phi_i^{T(y)_{(i)}} \\ &= \exp \left(\left(1 - \mathbf{1}^T T(y)\right) \ln \left(1 - \sum_{i=1}^k \phi_i\right) + \sum_{i=1}^k T(y)_{(i)} \ln \phi_i \right) \\ &= \exp \left(\ln(1 - \mathbf{1}^T \phi) - (\ln(1 - \mathbf{1}^T \phi) \cdot \mathbf{1})^T T(y) + T(y)^T \ln \phi \right) \\ &= \exp \left(T(y)^T \left(\ln \frac{\phi}{1 - \mathbf{1}^T \phi} \right) - (-\ln(1 - \mathbf{1}^T \phi)) \right) \end{aligned}$$

so, if we define ϕ_0 as $\phi_0 = 1 - \mathbf{1}^T \phi$:

$$\begin{aligned} b(y) &= 1 \\ \boldsymbol{\eta} &= \ln \frac{\phi}{\phi_0} \\ a(\boldsymbol{\eta}) &= -\ln(\phi_0) \end{aligned}$$

Using the fact that $\phi_0 + \mathbf{1}^T \phi = 1$, we can derive the response function:

$$\begin{aligned} \exp(\boldsymbol{\eta}) &= \frac{\phi}{\phi_0} \\ \mathbf{1}^T \exp(\boldsymbol{\eta}) &= \frac{\mathbf{1}^T \phi}{\phi_0} = \frac{1 - \phi_0}{\phi_0} \\ 1 &= \phi_0 \mathbf{1}^T \exp(\boldsymbol{\eta}) + \phi_0 \\ \phi_0 &= \frac{1}{1 + \mathbf{1}^T \exp(\boldsymbol{\eta})} \\ \phi &= \frac{\exp(\boldsymbol{\eta})}{1 + \mathbf{1}^T \exp(\boldsymbol{\eta})} \\ a(\boldsymbol{\eta}) &= \ln(1 + \mathbf{1}^T \exp(\boldsymbol{\eta})) \\ P(y; \boldsymbol{\eta}) &= \exp \left(\boldsymbol{\eta}^T T(y) - \ln(1 + \mathbf{1}^T \exp(\boldsymbol{\eta})) \right) \end{aligned}$$

Now we have parameterized it using $\boldsymbol{\eta}$, we can find its expected value, which leads to the hypothesis for the corresponding GLM:

$$\begin{aligned} h(X_{(i)}) &= \mathbb{E}[T(\mathbf{y}_{(i)}) | X_{(i)}; \boldsymbol{\Theta}] \\ &= \sum_{\mathbf{y}_{(i)}=0}^k \frac{T(\mathbf{y}_{(i)}) \exp(\boldsymbol{\eta}^T T(\mathbf{y}_{(i)}))}{1 + \mathbf{1}^T \exp(\boldsymbol{\eta})} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{1 + \mathbf{1}^T \exp(\boldsymbol{\eta})} \sum_{\mathbf{y}_{(i)}=1}^k T(\mathbf{y}_{(i)}) \exp(\boldsymbol{\eta}^T T(\mathbf{y}_{(i)})) \\
&= \frac{1}{1 + \mathbf{1}^T \exp(\boldsymbol{\eta})} \sum_{\mathbf{y}_{(i)}=1}^k T(\mathbf{y}_{(i)}) \exp\left(\boldsymbol{\eta}(\mathbf{y}_{(i)})\right) \\
&= \frac{\exp(\boldsymbol{\eta})}{1 + \mathbf{1}^T \exp(\boldsymbol{\eta})} = \boldsymbol{\phi} \\
&= \frac{\exp(\boldsymbol{\Theta}^T X_{(i)})}{1 + \mathbf{1}^T \exp(\boldsymbol{\Theta}^T X_{(i)})} \\
&= \begin{bmatrix} \frac{\exp(\boldsymbol{\Theta}_1^T X_{(i)})}{1 + \sum_{j=1}^k \exp(\boldsymbol{\Theta}_j^T X_{(i)})} \\ \frac{\exp(\boldsymbol{\Theta}_2^T X_{(i)})}{1 + \sum_{j=1}^k \exp(\boldsymbol{\Theta}_j^T X_{(i)})} \\ \vdots \\ \frac{\exp(\boldsymbol{\Theta}_k^T X_{(i)})}{1 + \sum_{j=1}^k \exp(\boldsymbol{\Theta}_j^T X_{(i)})} \end{bmatrix}
\end{aligned}$$

This model is called **Softmax Regression**, a generalization of logistic regression for multinomial distribution. Note that its parameter $\boldsymbol{\Theta}$ is a $n \times k$ matrix.

To fit the parameters, we need to find the maximum likelihood estimate:

$$\begin{aligned}
&\max_{\boldsymbol{\Theta}} \prod_{i=1}^m P(\mathbf{y}_{(i)} | X_{(i)}; \boldsymbol{\Theta}) \\
&= \max_{\boldsymbol{\Theta}} \sum_{i=1}^m ((\boldsymbol{\Theta}^T X_{(i)})^T T(\mathbf{y}_{(i)}) - \ln(1 + \mathbf{1}^T \exp(\boldsymbol{\Theta}^T X_{(i)}))) = \max_{\boldsymbol{\Theta}} -l(\boldsymbol{\Theta}) \\
&\nabla_{\boldsymbol{\Theta}} l(\boldsymbol{\Theta}) = - \sum_{i=1}^m \nabla_{\boldsymbol{\Theta}} \left(\text{tr} \left(X_{(i)}^T \boldsymbol{\Theta} T(\mathbf{y}_{(i)}) \right) - \ln(1 + \mathbf{1}^T \exp(\boldsymbol{\Theta}^T X_{(i)})) \right) \\
&= - \sum_{i=1}^m \left(X_{(i)} T(\mathbf{y}_{(i)})^T - \frac{X_{(i)} \exp(\boldsymbol{\Theta}^T X_{(i)})^T}{1 + \mathbf{1}^T \exp(\boldsymbol{\Theta}^T X_{(i)})} \right) \\
&= \sum_{i=1}^m X_{(i)} \left(\frac{\exp(\boldsymbol{\Theta}^T X_{(i)})}{1 + \mathbf{1}^T \exp(\boldsymbol{\Theta}^T X_{(i)})} - T(\mathbf{y}_{(i)}) \right)^T
\end{aligned}$$

Similar to logistic regression, there is no general closed form solution for softmax regression. Use either gradient descent or Newton's Method to maximize the joint probability.

2 Generative Learning Algorithms

All models/algorithms mentioned before in the last chapter are **discriminative learning algorithms**, which try to maximize the joint probability of $P(\mathbf{y}|X; \boldsymbol{\eta})$. This chapter will introduce algorithms are **generative learning algorithms** that models $P(X|\mathbf{y})$ and $P(\mathbf{y})$ (called **prior**), then use **Bayes rule**:

$$P(\mathbf{y}|x) = \frac{P(x|\mathbf{y})P(\mathbf{y})}{P(x)}$$

to derive the the distribution of \mathbf{y} given X . Here, $P(x)$ can be calculated by applying law of total probability:

$$P(x) = \int_{\mathbf{y}} y P(x|\mathbf{y})$$

2.1 Multivariate Normal Distribution

A multivariate distribution a vector of random variables outputting values in an n -dimensional space. To describe variance of high-dimensional data, we define **covariance** as the average outer product of its deviation (from its average) with itself, that is:

$$\text{Cov}(\mathbf{Z}) = \mathbb{E}[(\mathbf{Z} - \mathbb{E}[\mathbf{Z}])(\mathbf{Z} - \mathbb{E}[\mathbf{Z}])^T]$$

We can derive another form of the covariance:

$$\begin{aligned} \text{Cov}(\mathbf{Z}) &= \mathbb{E}[\mathbf{Z}\mathbf{Z}^T - \mathbb{E}[\mathbf{Z}]\mathbf{Z}^T - \mathbf{Z}\mathbb{E}[\mathbf{Z}]^T + \mathbb{E}[\mathbf{Z}]\mathbb{E}[\mathbf{Z}]^T] \\ &= \mathbb{E}[\mathbf{Z}\mathbf{Z}^T] - \mathbb{E}[\mathbf{Z}]\mathbb{E}[\mathbf{Z}]^T - \mathbb{E}[\mathbf{Z}]\mathbb{E}[\mathbf{Z}]^T + \mathbb{E}[\mathbf{Z}]\mathbb{E}[\mathbf{Z}]^T \\ &= \mathbb{E}[\mathbf{Z}\mathbf{Z}^T] - \mathbb{E}[\mathbf{Z}]\mathbb{E}[\mathbf{Z}]^T \end{aligned}$$

Multivariate normal distribution is defined as following:

$$\begin{aligned} P(\mathbf{y}; \boldsymbol{\mu}, \Sigma) &= \mathcal{N}(\boldsymbol{\mu}, \Sigma) \\ &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{y} - \boldsymbol{\mu})\right) \end{aligned}$$

Multivariate normal distribution is an analogy of the single-variable counterpart in higher dimensional space. The probability density of \mathbf{y} is highest at mean vector $\boldsymbol{\mu}$, decreases along a bell curve as \mathbf{y} moves away from the mean vector. We see the probability density of a value \mathbf{y} that is d units away from $\boldsymbol{\mu}$ along the unit vector $\hat{\mathbf{u}}$ is:

$$\begin{aligned} P(\mathbf{y} = d\hat{\mathbf{u}} + \boldsymbol{\mu}; \boldsymbol{\mu}, \Sigma) &= \mathcal{N}(\boldsymbol{\mu}, \Sigma) \\ &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(d\hat{\mathbf{u}})^T \Sigma^{-1}(d\hat{\mathbf{u}})\right) \\ &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{d^2}{2}\hat{\mathbf{u}}^T \Sigma^{-1}\hat{\mathbf{u}}\right) \end{aligned}$$

Since Σ is symmetric and postive semi-definite by definition, there is $\Sigma = Q\Lambda Q^T$, where the column vectors of Q are the eigenvectors of Σ and Λ is a diagonal matrix formed by the respective eigenvalues, all of which are greater or equal to 0:

$$\begin{aligned} &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{d^2}{2}\hat{\mathbf{u}}^T Q\Lambda^{-1}Q^T\hat{\mathbf{u}}\right) \\ &= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{d^2}{2}\hat{\mathbf{u}}^T Q\Lambda^{-1/2}\Lambda^{-1/2}Q^T\hat{\mathbf{u}}\right) \end{aligned}$$

Note that since $\Lambda = \text{Diag}(\boldsymbol{\lambda})$, $\Lambda^{-1} = \text{Diag}(1/\boldsymbol{\lambda})$ and $\Lambda^{-1/2} = \text{Diag}(\boldsymbol{\lambda}^{-1/2})$, so $\Lambda^{-1/2}$ is also a symmetric, positive semi-definite matrix:

=

2.1.1 Gaussian Discriminant Analysis

$$\begin{aligned} y &\sim \text{Bernoulli}(\phi) \\ \mathbf{x}|y=0 &\sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma) \\ \mathbf{x}|y=1 &\sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma) \end{aligned}$$

(Negative) Joint Log Probability:

$$\begin{aligned} l(\phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) &= -\ln \prod_{i=1}^m p(X_{(i)}|\mathbf{y}_{(i)}; \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) p(\mathbf{y}_{(i)}; \phi) \\ &= -\ln \prod_{i=1}^m p(X_{(i)}|\mathbf{y}_{(i)} = 1; \boldsymbol{\mu}_1, \Sigma)^{\mathbf{y}_{(i)}} p(X_{(i)}|\mathbf{y}_{(i)} = 0; \boldsymbol{\mu}_0, \Sigma)^{(1 - \mathbf{y}_{(i)})} p(\mathbf{y}_{(i)}; \phi) \\ &= -\sum_{i=1}^m \left(-\frac{n}{2} \mathbf{y}_{(i)} \ln 2\pi - \frac{1}{2} \mathbf{y}_{(i)} \ln |\Sigma| - \mathbf{y}_{(i)} \left(\frac{1}{2} (X_{(i)} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (X_{(i)} - \boldsymbol{\mu}_1) \right) \right) \end{aligned}$$

$$\begin{aligned}
& - \sum_{i=1}^m \left(-\frac{n}{2}(1 - \mathbf{y}_{(i)}) \ln 2\pi - \frac{1}{2}(1 - \mathbf{y}_{(i)}) \ln |\Sigma| - (1 - \mathbf{y}_{(i)}) \left(\frac{1}{2}(X_{(i)} - \boldsymbol{\mu}_0)^T \Sigma^{-1} (X_{(i)} - \boldsymbol{\mu}_0) \right) \right) \\
& - \sum_{i=1}^m \left(\mathbf{y}_{(i)} \ln \phi + (1 - \mathbf{y}_{(i)}) \ln(1 - \phi) \right) \\
& = \sum_{i=1}^m \left(\frac{n}{2} \ln 2\pi + \frac{1}{2} \ln |\Sigma| + \mathbf{y}_{(i)} \left(\frac{1}{2}(X_{(i)} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (X_{(i)} - \boldsymbol{\mu}_1) \right) \right. \\
& \quad \left. + (1 - \mathbf{y}_{(i)}) \left(\frac{1}{2}(X_{(i)} - \boldsymbol{\mu}_0)^T \Sigma^{-1} (X_{(i)} - \boldsymbol{\mu}_0) \right) \right. \\
& \quad \left. - \mathbf{y}_{(i)} \ln \phi - (1 - \mathbf{y}_{(i)}) \ln(1 - \phi) \right) \\
& = \frac{mn}{2} \ln 2\pi + \frac{m}{2} \ln |\Sigma| + \frac{1}{2} \sum_{i=1}^m \mathbf{y}_{(i)} (X - \mathbf{1}\boldsymbol{\mu}_1^T)_{(i)}^T \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_1^T)_{(i)} \\
& \quad + \frac{1}{2} \sum_{i=1}^m (1 - \mathbf{y})_{(i)} (X - \mathbf{1}\boldsymbol{\mu}_0^T)_{(i)}^T \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_0^T)_{(i)} \\
& \quad - \mathbf{1}^T (\mathbf{y} \ln \phi + (1 - \mathbf{y}) \ln(1 - \phi)) \\
& = \frac{mn}{2} \ln 2\pi + \frac{m}{2} \ln |\Sigma| + \frac{1}{2} \sum_{i=1}^m ((X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \text{Diag}(\mathbf{y}))_i^T (\Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_1^T)^T)_i \\
& \quad + \frac{1}{2} \sum_{i=1}^m ((X - \mathbf{1}\boldsymbol{\mu}_0^T)^T \text{Diag}(1 - \mathbf{y}))_i^T (\Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_0^T)^T)_i \\
& \quad - \mathbf{1}^T (\mathbf{y} \ln \phi + (1 - \mathbf{y}) \ln(1 - \phi)) \\
& = \frac{mn}{2} \ln 2\pi + \frac{m}{2} \ln |\Sigma| + \frac{1}{2} \text{tr} \text{Diag}(\mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_1^T) \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \\
& \quad + \frac{1}{2} \text{tr} \text{Diag}(1 - \mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_0^T) \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_0^T)^T \\
& \quad - \mathbf{1}^T (\mathbf{y} \ln \phi + (1 - \mathbf{y}) \ln(1 - \phi))
\end{aligned}$$

Gradient with respect to parameters:

$$\begin{aligned}
dl(\phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) &= -\mathbf{1}^T \left(\frac{\mathbf{y}}{\phi} - \frac{1 - \mathbf{y}}{1 - \phi} \right) d\phi \\
\mathbf{1}^T \mathbf{y} (1 - \phi) &= \mathbf{1}^T (1 - \mathbf{y}) \phi \\
\mathbf{1}^T (2\mathbf{y} - 1) \phi &= \mathbf{1}^T \mathbf{y} \\
\mathbf{1}^T \mathbf{y} - \mathbf{1}^T \mathbf{y} \phi &= \mathbf{1}^T \mathbf{1} \phi - \mathbf{1}^T \mathbf{y} \phi \\
\phi &= \frac{\mathbf{1}^T \mathbf{y}}{\mathbf{1}^T \mathbf{1}} = \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{\mathbf{y}_{(i)} = 1\}
\end{aligned}$$

Gradient with respect to $\boldsymbol{\mu}_1$

$$\begin{aligned}
dl(\phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) &= \frac{1}{2} d \text{tr} \text{Diag}(\mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_1^T) \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \\
&= \frac{1}{2} \text{tr} d \left(\text{Diag}(\mathbf{y}) \mathbf{1}\boldsymbol{\mu}_1^T \Sigma^{-1} \mathbf{1}^T - \text{Diag}(\mathbf{y}) X \Sigma^{-1} \mathbf{1}^T - \text{Diag}(\mathbf{y}) \mathbf{1}\boldsymbol{\mu}_1^T \Sigma^{-1} X^T \right) \\
&= \frac{1}{2} d \left(\text{tr} \mathbf{1}^T \text{Diag}(\mathbf{y}) \mathbf{1}\boldsymbol{\mu}_1^T \Sigma^{-1} \mathbf{1}^T - \text{tr} \mathbf{1}^T \text{Diag}(\mathbf{y}) X \Sigma^{-1} \mathbf{1}^T - \text{tr} \boldsymbol{\mu}_1^T \Sigma^{-1} X^T \text{Diag}(\mathbf{y}) \mathbf{1} \right) \\
&= \frac{1}{2} \left(2 \text{tr} \mathbf{1}^T \text{Diag}(\mathbf{y}) \mathbf{1}\boldsymbol{\mu}_1^T \Sigma^{-1} d\boldsymbol{\mu}_1 - \text{tr} \mathbf{1}^T \text{Diag}(\mathbf{y}) X \Sigma^{-1} d\boldsymbol{\mu}_1 - \text{tr} \mathbf{1}^T \text{Diag}(\mathbf{y}) X \Sigma^{-1} d\boldsymbol{\mu}_1 \right) \\
&= \text{tr} \mathbf{1}^T \text{Diag}(\mathbf{y}) \mathbf{1}\boldsymbol{\mu}_1^T \Sigma^{-1} d\boldsymbol{\mu}_1 - \text{tr} \mathbf{1}^T \text{Diag}(\mathbf{y}) X \Sigma^{-1} d\boldsymbol{\mu}_1 \\
&= \text{tr} \mathbf{y}^T (\mathbf{1}\boldsymbol{\mu}_1^T - X) \Sigma^{-1} d\boldsymbol{\mu}_1 \\
\nabla_{\boldsymbol{\mu}_1} l(\boldsymbol{\mu}_1) &= \Sigma^{-1} (\mathbf{1}\boldsymbol{\mu}_1^T - X)^T \mathbf{y}
\end{aligned}$$

Since Σ is not singular, (or we assume so), $(\mathbf{1}\boldsymbol{\mu}_1^T - X)^T \mathbf{y}$ must be 0 for the product to be 0:

$$(\mathbf{1}\boldsymbol{\mu}_1^T - X)^T \mathbf{y} = 0$$

$$\begin{aligned}\boldsymbol{\mu}_1 \mathbf{1}^T \mathbf{y} &= X^T \mathbf{y} \\ \boldsymbol{\mu}_1 &= \frac{X^T \mathbf{y}}{\mathbf{1}^T \mathbf{y}} = \frac{\sum_{i=1}^m \mathbf{1}\{\mathbf{y}_{(i)} = 1\} X_{(i)}}{\sum_{i=1}^m \mathbf{1}\{\mathbf{y}_{(i)} = 1\}}\end{aligned}$$

Via a similar process, the gradient with respect to $\boldsymbol{\mu}_0$ is:

$$\begin{aligned}dl(\phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) &= \text{tr} (1 - \mathbf{y})^T (\mathbf{1}\boldsymbol{\mu}_0^T - X) \Sigma^{-1} d\boldsymbol{\mu}_0 \\ \nabla_{\boldsymbol{\mu}_0} l(\boldsymbol{\mu}_0) &= \Sigma^{-1} (\mathbf{1}\boldsymbol{\mu}_0^T - X)^T (1 - \mathbf{y}) \\ 0 &= (\mathbf{1}\boldsymbol{\mu}_0^T - X)^T (1 - \mathbf{y}) \\ \boldsymbol{\mu}_0 &= \frac{X^T (1 - \mathbf{y})}{\mathbf{1}^T (1 - \mathbf{y})} = \frac{\sum_{i=1}^m \mathbf{1}\{\mathbf{y}_{(i)} = 0\} X_{(i)}}{\sum_{i=1}^m \mathbf{1}\{\mathbf{y}_{(i)} = 0\}}\end{aligned}$$

Gradient with respect to Σ :

$$\begin{aligned}dl(\phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) &= \frac{m}{2} \text{tr} \Sigma^{-1} d\Sigma - \frac{1}{2} \text{tr} \text{Diag}(\mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_1^T) \Sigma^{-1} (d\Sigma) \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \\ &\quad - \frac{1}{2} \text{tr} \text{Diag}(1 - \mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_0^T) \Sigma^{-1} (d\Sigma) \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_0^T)^T \\ &= \frac{m}{2} \text{tr} \Sigma^{-1} \Sigma \Sigma^{-1} d\Sigma - \frac{1}{2} \text{tr} \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \text{Diag}(\mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_1^T) \Sigma^{-1} d\Sigma \\ &\quad - \frac{1}{2} \text{tr} \Sigma^{-1} (X - \mathbf{1}\boldsymbol{\mu}_0^T)^T \text{Diag}(1 - \mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_0^T) \Sigma^{-1} d\Sigma \\ &= \text{tr} \Sigma^{-1} \left(m\Sigma - (X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \text{Diag}(\mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_1^T) - (X - \mathbf{1}\boldsymbol{\mu}_0^T)^T \text{Diag}(1 - \mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_0^T) \right) \Sigma^{-1} d\Sigma \\ \nabla_{\Sigma} l(\Sigma) &= \Sigma^{-1} \left(m\Sigma - (X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \text{Diag}(\mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_1^T) - (X - \mathbf{1}\boldsymbol{\mu}_0^T)^T \text{Diag}(1 - \mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_0^T) \right) \Sigma^{-1} \\ \Sigma &= \frac{1}{m} \left((X - \mathbf{1}\boldsymbol{\mu}_1^T)^T \text{Diag}(\mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_1^T) + (X - \mathbf{1}\boldsymbol{\mu}_0^T)^T \text{Diag}(1 - \mathbf{y}) (X - \mathbf{1}\boldsymbol{\mu}_0^T) \right) \\ &= \frac{1}{m} (X - (\mathbf{y}\boldsymbol{\mu}_1^T + (1 - \mathbf{y})\boldsymbol{\mu}_0^T)) (X - (\mathbf{y}\boldsymbol{\mu}_1^T + (1 - \mathbf{y})\boldsymbol{\mu}_0^T))^T\end{aligned}$$

Gaussian discriminant analysis is a special case of logistic regression. It assumes a stronger assumption that the data is normally-distributed around the respective center of mass of the classes. When this assumption is right, the GDA is **asymptotically efficient**, which means there is no better algorithms that GDA in terms of accuracy given a set amount of training data.

However, when this assumption is not true, logistic regression finds its better fit, making it more robust when the distribution of the data is not well known.

2.2 Naive Bayes

$$\begin{aligned}y &\sim \text{Bernoulli}(\phi) \\ \mathbf{x}_{(j)} | y = 0 &\sim \text{Bernoulli}(\phi_{j|y=0}) \\ \mathbf{x}_{(j)} | y = 1 &\sim \text{Bernoulli}(\phi_{j|y=1})\end{aligned}$$

Naive Bayes Assumption: $\mathbf{x}_{(j)}$'s are conditionally independent given y , i.e.

$$p(\mathbf{x}|y) = \prod_{j=1}^n p(\mathbf{x}_{(j)}|y)$$

Although Naive Bayes Assumption is almost never strictly true, the algorithm resulting from usually works well on many problems. Now let's maximize the joint probability of the model and gradients for our parameters:

$$\begin{aligned}\min l(\phi, \phi_0^T) &= [\phi_{1|y=0} \dots \phi_{n|y=0}], \phi_1^T = [\phi_{1|y=1} \dots \phi_{n|y=1}] \\ &= -\ln \prod_{i=1}^m p(X_{(i)} | \mathbf{y}_{(i)}) p(\mathbf{y}_{(i)}) \\ &= -\ln \prod_{i=1}^m \prod_{j=1}^n p(X_{ij} | \mathbf{y}_{(i)}) \phi^{\mathbf{y}_{(i)}} (1 - \phi)^{1 - \mathbf{y}_{(i)}}\end{aligned}$$

$$\begin{aligned}
&= -\ln \prod_{i=1}^m \prod_{j=1}^n p(X_{ij} | \mathbf{y}_{(i)} = 1)^{\mathbf{y}_{(i)}} p(X_{ij} | \mathbf{y}_{(i)} = 0)^{1-\mathbf{y}_{(i)}} \phi^{\mathbf{y}_{(i)}} (1-\phi)^{1-\mathbf{y}_{(i)}} \\
&= -\ln \prod_{i=1}^m \prod_{j=1}^n \phi_{j|y=1}^{\mathbf{y}_{(i)} X_{ij}} (1-\phi_{j|y=1})^{\mathbf{y}_{(i)}(1-X_{ij})} \phi_{j|y=0}^{(1-\mathbf{y}_{(i)})X_{ij}} (1-\phi_{j|y=0})^{(1-\mathbf{y}_{(i)})(1-X_{ij})} \phi^{\mathbf{y}_{(i)}} (1-\phi)^{1-\mathbf{y}_{(i)}} \\
&= -\sum_{i=1}^m \sum_{j=1}^n \left(\mathbf{y}_{(i)} X_{ij} \ln \phi_{j|y=1} + \mathbf{y}_{(i)} (1-X_{ij}) \ln(1-\phi_{j|y=1}) + (1-\mathbf{y}_{(i)}) X_{ij} \ln \phi_{j|y=0} + (1-\mathbf{y}_{(i)}) (1-X_{ij}) \ln(1-\phi_{j|y=0}) \right) \\
&\quad - n \sum_{i=1}^m \left(\mathbf{y}_{(i)} \ln \phi + (1-\mathbf{y}_{(i)}) \ln(1-\phi) \right) \\
&= -\sum_{i=1}^m \left(\mathbf{y}_{(i)} X_{(i)}^T \ln \phi_1 + \mathbf{y}_{(i)} (1-X_{(i)})^T \ln(1-\phi_1) + (1-\mathbf{y}_{(i)}) X_{(i)}^T \ln \phi_0 + (1-\mathbf{y}_{(i)}) (1-X_{(i)})^T \ln(1-\phi_0) \right) \\
&\quad - n(\mathbf{1}^T \mathbf{y} \ln \phi + \mathbf{1}^T (1-\mathbf{y}) \ln(1-\phi)) \\
&= -\mathbf{y}^T (X \ln \phi_1 - (1-X) \ln(1-\phi_1)) - (1-\mathbf{y})^T (X \ln \phi_0 + (1-X) \ln(1-\phi_0)) - n(\mathbf{1}^T \mathbf{y} \ln \phi + \mathbf{1}^T (1-\mathbf{y}) \ln(1-\phi))
\end{aligned}$$

Gradient:

$$\begin{aligned}
&\text{d}l(\phi, \phi_0, \phi_1) \\
&= \mathbf{y}^T \left(X \text{Diag} \left(\frac{1}{\phi_1} \right) \text{d}\phi_1 - (1-X) \text{Diag} \left(\frac{1}{1-\phi_1} \right) \text{d}\phi_1 \right) + (1-\mathbf{y})^T \left(X \text{Diag} \left(\frac{1}{\phi_0} \right) \text{d}\phi_0 - (1-X) \text{Diag} \left(\frac{1}{1-\phi_0} \right) \text{d}\phi_0 \right) \\
&\quad + n \left(\mathbf{1}^T \mathbf{y} \frac{\text{d}\phi}{\phi} - \mathbf{1}^T (1-\mathbf{y}) \frac{\text{d}\phi}{1-\phi} \right)
\end{aligned}$$

Solution for ϕ_1 :

$$\begin{aligned}
&\text{Diag} \left(\frac{1}{\phi_1} \right) X^T \mathbf{y} = \text{Diag} \left(\frac{1}{1-\phi_1} \right) (1-X)^T \mathbf{y} \\
&\text{Diag} (1-\phi_1) X^T \mathbf{y} = \text{Diag} (\phi_1) (1-X)^T \mathbf{y} \\
&X^T \mathbf{y} = \text{Diag} (\phi_1) ((1-X)^T \mathbf{y} + X^T \mathbf{y}) \\
&X^T \mathbf{y} = \phi_1 \circ (1_X^T \mathbf{y}) = \phi_1 (1^T \mathbf{y}) \\
&\phi_1 = \frac{X^T \mathbf{y}}{1^T \mathbf{y}}
\end{aligned}$$

Similarly:

$$\phi_2 = \frac{X^T (1-\mathbf{y})}{1^T (1-\mathbf{y})}$$

For ϕ , we have:

$$\begin{aligned}
&\frac{1^T \mathbf{y}}{\phi} = \frac{1^T (1-\mathbf{y})}{1-\phi} \\
&1^T \mathbf{y} - 1^T \mathbf{y} \phi = 1^T (1-\mathbf{y}) \phi \\
&1^T \mathbf{y} = 1^T \mathbf{y} \phi = m \phi \\
&\phi = \frac{1^T \mathbf{y}}{m}
\end{aligned}$$

3 Support Vector Machine

3.1 Preliminaries

3.1.1 Lagrange Multiplier

Lagrange Multiplier is a method of finding critical points of a scalar-valued function subject to equality constraints, namely

$$\max_{\mathbf{x}} f(\mathbf{x}) \quad \text{given } g(\mathbf{x}) = 0$$

or:

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{given } g(\mathbf{x}) = 0$$

Such critical points always lies on points where the contours of the two functions are parallel, and since the contour of function is always perpendicular to the gradient at any point, that means we need to find point \mathbf{x}' such that:

$$\begin{aligned}\nabla_{\mathbf{x}} f(\mathbf{x}') &= \lambda \nabla_{\mathbf{x}} g(\mathbf{x}') \\ g(\mathbf{x}') &= 0\end{aligned}$$

This can be cleverly written in term of the gradient of **Lagrangian** being 0:

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \lambda) &= f(\mathbf{x}) - \lambda g(\mathbf{x}) \\ \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}', \lambda) &= \begin{bmatrix} \nabla_{\mathbf{x}} f(\mathbf{x}') - \lambda \nabla_{\mathbf{x}} g(\mathbf{x}') \\ g(\mathbf{x}') \end{bmatrix} = 0\end{aligned}$$

The solutions can be maxima, minima, or saddles, but if $g(\mathbf{x}) = 0$ contains a closed and bounded region, there has to be at least 2 solutions, corresponding to its minimum and maximum.

3.1.2 Distance of a Point to a Hyperplane

Using Lagrange Multiplier, Let's calculate the distance of a point \mathbf{x}_0 in to a hyperplane $\boldsymbol{\theta}^T \mathbf{x} + b = 0$

$$\min_{\mathbf{x}} \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{given } \boldsymbol{\theta}^T \mathbf{x} + b = 0$$

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \lambda) &= (\mathbf{x} - \mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) - \lambda (\boldsymbol{\theta}^T \mathbf{x} + b) \\ \nabla_{\mathbf{x}, \lambda} \mathcal{L}(\mathbf{x}, \lambda) &= \begin{bmatrix} 2(\mathbf{x} - \mathbf{x}_0) - \lambda \boldsymbol{\theta} \\ -(\boldsymbol{\theta}^T \mathbf{x} + b) \end{bmatrix} = 0 \\ \lambda \boldsymbol{\theta} &= 2(\mathbf{x} - \mathbf{x}_0) \\ \lambda &= \frac{2\boldsymbol{\theta}^T (\mathbf{x} - \mathbf{x}_0)}{\boldsymbol{\theta}^T \boldsymbol{\theta}} = \frac{-2(b + \boldsymbol{\theta}^T \mathbf{x}_0)}{\boldsymbol{\theta}^T \boldsymbol{\theta}} & \boldsymbol{\theta}^T \mathbf{x} = -b \\ \frac{\lambda \boldsymbol{\theta}}{2} &= (\mathbf{x} - \mathbf{x}_0) \\ \frac{\lambda^2 \boldsymbol{\theta}^T \boldsymbol{\theta}}{4} &= (\mathbf{x} - \mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) \\ (\mathbf{x} - \mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) &= \frac{(\boldsymbol{\theta}^T \mathbf{x}_0 + b)^2}{\boldsymbol{\theta}^T \boldsymbol{\theta}} \\ \|\mathbf{x} - \mathbf{x}_0\| &= \frac{|\boldsymbol{\theta}^T \mathbf{x}_0 + b|}{\|\boldsymbol{\theta}\|}\end{aligned}$$

3.1.3 Lagrange duality

Given an optimization problem of following form:

$$\begin{aligned}\min_{\mathbf{w}} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & \mathbf{g}(\mathbf{w}) \leq 0 \\ & \mathbf{h}(\mathbf{w}) = 0\end{aligned}$$

define its **generalized Lagrangian** to be:

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{w}) + \boldsymbol{\alpha}^T \mathbf{g}(\mathbf{w}) + \boldsymbol{\beta}^T \mathbf{h}(\mathbf{w})$$

Note that if we try to maximize the Lagrangian with the condition that $\boldsymbol{\alpha} \geq 0$, it is only bounded if the constraints are satisfied, otherwise, that it is bounded, it maximize to $f(\mathbf{w})$, in other words, we have:

$$\theta_{\mathcal{P}}(\mathbf{w}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}: \boldsymbol{\alpha} \geq 0} \mathcal{L}(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \begin{cases} f(\mathbf{w}) & \mathbf{g}(\mathbf{w}) \leq 0 \wedge \mathbf{h}(\mathbf{w}) = 0 \\ \infty & \text{otherwise} \end{cases}$$

Minimizing this problem is exactly same as minimizing $f(\mathbf{w})$ conforming to the given constraints, and this is what we call **primal problem**:

$$\min_{\mathbf{w}} \theta_{\mathcal{P}}(\mathbf{w}) = \min_{\mathbf{w}} \max_{\alpha, \beta: \mathbf{a} \geq 0} \mathcal{L}(\mathbf{w}, \alpha, \beta)$$

If we reverse the order of min and max, we have the **dual problem**:

$$\max_{\alpha, \beta: \mathbf{a} \geq 0} \theta_{\mathcal{D}}(\mathbf{w}) = \max_{\alpha, \beta: \mathbf{a} \geq 0} \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \alpha, \beta)$$

Let p^* be the solution to the primal problem and d^* be the solution to the dual problem:

$$\begin{aligned} d^* &= \max_{\alpha, \beta: \mathbf{a} \geq 0} \min_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \alpha, \beta) \leq \min_{\mathbf{w}} \max_{\alpha, \beta: \mathbf{a} \geq 0} \mathcal{L}(\mathbf{w}, \alpha, \beta) = p^* \\ \max_{\alpha, \beta: \mathbf{a} \geq 0} \min_{\mathbf{w}} \left(f(\mathbf{w}) + \alpha^T \mathbf{g}(\mathbf{w}) + \beta^T \mathbf{h}(\mathbf{w}) \right) \end{aligned}$$

This property is called weak duality, which holds for any continuous function \mathcal{L} . However when **Slater's condition** are met:

- f and $\mathbf{g}_{(i)}$ are convex for all i
- $\mathbf{h}_{(i)}$ is affine for all i
- there exists \mathbf{w} such that $\mathbf{g}_{(i)}(\mathbf{w}) < 0$ for all i

then, **strong duality** holds, i.e:

$$d^* = p^*$$

Another condition, that are both sufficient and necessary when strong duality hold is the **Karush-Kuhn-Tucker condtion**, formulated as following, given a feasible solution \mathbf{w}^* :

$$\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}^*, \alpha, \beta) &= 0 \\ \alpha \circ \mathbf{g}(\mathbf{w}^*) &= 0 \\ \alpha &\geq 0 \end{aligned}$$

3.2 Optimal Margin Classifier

Given a **linearly separable** dataset $(X_{(i)}, \mathbf{y}_{(i)}) \in (\mathbb{R}^n, \{-1, 1\})$, let $\boldsymbol{\theta}$ and b define a hyperplane that properly separates this dataset. The hypothesis for this classifier is:

$$h_{\boldsymbol{\theta}}(X_{(i)}) = \text{sign}(\boldsymbol{\theta}^T X_{(i)} + b)$$

We define the **geometric margin** of the i th data point to be:

$$\gamma_i = \frac{\mathbf{y}_{(i)}(\boldsymbol{\theta}^T X_{(i)} + b)}{\|\boldsymbol{\theta}\|}$$

Note that margin for a particular data point is the distance of that data point to the separating hyperplane when the hypothesis matches its label, and the negative of that distance when the hypothesis mismatches. We can use this to denotes the confidence we have towards the hypothesis for this particular data point. For the entire dataset, we use the minimal margin in the dataset:

$$\gamma = \min_{i=1, \dots, m} \frac{\mathbf{y}_{(i)}(\boldsymbol{\theta}^T X_{(i)} + b)}{\|\boldsymbol{\theta}\|} = \frac{1}{\|\boldsymbol{\theta}\|} \min_{i=1, \dots, m} \mathbf{y}_{(i)}(\boldsymbol{\theta}^T X_{(i)} + b) = \frac{\hat{\gamma}}{\|\boldsymbol{\theta}\|}$$

To find the separating hyperplane, we maximize the margin for our dataset:

$$\begin{aligned} \max_{\hat{\gamma}, \boldsymbol{\theta}} & \frac{\hat{\gamma}}{\|\boldsymbol{\theta}\|} \\ \text{s.t. } & \mathbf{y}_{(i)}(\boldsymbol{\theta}^T X_{(i)} + b) \geq \hat{\gamma} \end{aligned}$$

Since our model is invariant to scaling, let's scale it down by $\hat{\gamma}$:

$$\begin{aligned} \max_{\hat{\gamma}, \boldsymbol{\theta}} & \frac{1}{\|\boldsymbol{\theta}\|} \\ \text{s.t. } & \mathbf{y}_{(i)}(\boldsymbol{\theta}^T X_{(i)} + b) \geq 1 \end{aligned}$$

Then we change maximize $\|\boldsymbol{\theta}\|$ to minimize $\|\boldsymbol{\theta}\|^2$ to allow our model solved by quadratic programming:

$$\begin{aligned} \min_{\hat{\gamma}, \boldsymbol{\theta}} & \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} \\ \text{s.t. } & X\boldsymbol{\theta} + b \leq -\mathbf{y} \end{aligned}$$

3.2.1 Lagrange Dual Form

Let's write our primal in a more amenable form:

$$\begin{aligned} \min \quad & \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} \\ \text{s.t.} \quad & 1 - (X\boldsymbol{\theta} + b) \circ \mathbf{y} \leq 0 \end{aligned}$$

From the constraints of this problem, we can see that only $\mathbf{g}(\boldsymbol{\theta})_{(i)} = 0$ only when $X_{(i)}$ is a support vector, so to satisfy $\mathbf{a} \circ \mathbf{g}(\boldsymbol{\theta}) = 0$, $\boldsymbol{\alpha}_{(i)}$ must be zero unless $X_{(i)}$ is a support vector.

The Lagrangian is:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha}) &= \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + \boldsymbol{\alpha}^T (1 - (X\boldsymbol{\theta} + b) \circ \mathbf{y}) \\ \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha}) &= \boldsymbol{\theta} - X^T \text{Diag}(\mathbf{y}) \boldsymbol{\alpha} = 0 \\ \boldsymbol{\theta} &= X^T \text{Diag}(\mathbf{y}) \boldsymbol{\alpha} \\ \nabla_b \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha}) &= \mathbf{y}^T \boldsymbol{\alpha} = 0 \end{aligned}$$

Plugging the definition of $\boldsymbol{\theta}$ to the Lagrangian, we have:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha}) &= \frac{1}{2} \boldsymbol{\alpha}^T \text{Diag}(\mathbf{y}) X X^T \text{Diag}(\mathbf{y}) \boldsymbol{\alpha} + \boldsymbol{\alpha}^T \mathbf{1} - \boldsymbol{\alpha}^T \text{Diag}(\mathbf{y}) X X^T \text{Diag}(\mathbf{y}) \boldsymbol{\alpha} - (\boldsymbol{\alpha}^T \mathbf{y}) b \\ &= \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} (\boldsymbol{\alpha} \circ \mathbf{y})^T X X^T (\boldsymbol{\alpha} \circ \mathbf{y}) \end{aligned}$$

To find $\boldsymbol{\alpha}^*$, maximize the dual problem:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \left(\boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} (\boldsymbol{\alpha} \circ \mathbf{y})^T X X^T (\boldsymbol{\alpha} \circ \mathbf{y}) \right) \\ \text{s.t.} \quad & \boldsymbol{\alpha} \geq 0 \\ & \boldsymbol{\alpha}^T \mathbf{y} = 0 \end{aligned}$$

then calculate $\boldsymbol{\theta}^*$ with:

$$\begin{aligned} \boldsymbol{\theta}^* &= X^T (\mathbf{y} \circ \boldsymbol{\alpha}) \\ \mathbf{b}^* &= - \frac{\max_{i: y_{(i)} = -1} \boldsymbol{\theta}^{*T} X_{(i)} + \min_{i: y_{(i)} = 1} \boldsymbol{\theta}^{*T} X_{(i)}}{2} \end{aligned}$$

To exploit the fact that $\boldsymbol{\alpha}$ is mostly zero, let X_s denotes a matrix of support vectors (in rows), \mathbf{y}_s and \mathbf{a}_s denote corresponding label and $\boldsymbol{\alpha}_{(i)}$, we have:

$$\begin{aligned} X\boldsymbol{\theta} &= X X^T (\mathbf{y} \circ \boldsymbol{\alpha}) \\ &= X X_s^T (\mathbf{y}_s \circ \boldsymbol{\alpha}_s) \end{aligned}$$

where $X X_s^T$ is called **kernel** of the SVM; as we will see later, we can change the kernel to classify linearly inseparable data.

3.2.2 Kernel

Kernels like $K(\mathbf{x}, \mathbf{z})$ are functions that map its arguments to higher dimensions then calculating their inner products. For example, define kernel K as:

$$\begin{aligned} K(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z} + c)^2 \\ &= ((x_1 z_1)^2 + (x_2 z_2)^2 + (x_3 z_3)^2 + 2x_1 z_1 x_2 z_2 + 2x_1 z_1 x_3 z_3 + 2x_2 z_2 x_3 z_3 + 2x_1 z_1 c + 2x_2 z_2 c + 2x_3 z_3 c + c^2) \\ &= \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{z}) \\ \boldsymbol{\phi}(\mathbf{x}) &= \begin{bmatrix} x_1^2 & x_2^2 & x_3^2 & \sqrt{2}x_1 x_2 & \sqrt{2}x_1 x_3 & \sqrt{2}x_2 x_3 & \sqrt{2c}x_1 & \sqrt{2c}x_2 & \sqrt{2c}x_3 & c \end{bmatrix}^T \end{aligned}$$

Kernels with higher degrees like $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + c)^d$ projects the features to an $\frac{(n+d)!}{n!d!}$ feature space, while keep $O(n)$ computational complexity. There is also **Gaussian kernel**, which:

$$\begin{aligned} K(\mathbf{x}, \mathbf{z}) &= \exp\left(-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}\right) \\ &= 1 - \frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2} + \frac{\|\mathbf{x} - \mathbf{z}\|^4}{4\sigma^4} - \frac{\|\mathbf{x} - \mathbf{z}\|^6}{8\sigma^6} + \dots \end{aligned}$$

maps the features onto an infinite-dimensional feature space, as shown above.

Each kernel $K(\mathbf{x}, \mathbf{z})$ is also associated with a **kernel matrix** K , defined as such:

$$\begin{aligned} K_{ij} &= K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{given a list of vectors } \mathbf{x}_1 \dots \mathbf{x}_m &\in \mathbb{R}^n \end{aligned}$$

We can show that K is a symmetric, positive semi-definite matrix:

$$K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = K(\mathbf{x}_j, \mathbf{x}_i) = K_{ji}$$

To prove its definiteness, let's define $\mathbf{y}_1 \dots \mathbf{y}_m \in \mathbb{R}^m$ where \mathbf{y}_j contains the j th elements of $\phi(\mathbf{x}_1) \dots \phi(\mathbf{x}_m)$, for any vector $\mathbf{z} \in \mathbb{R}^m$:

$$\mathbf{z}^T K \mathbf{z} = \mathbf{z}^T \left(\sum_{i=1}^m \mathbf{y}_i \mathbf{y}_i^T \right) \mathbf{z} = \sum_{i=1}^m \mathbf{z}^T \mathbf{y}_i \mathbf{y}_i^T \mathbf{z} = \sum_{i=1}^m (\mathbf{y}_i^T \mathbf{z})^T (\mathbf{y}_i^T \mathbf{z}) \geq 0$$

Thus, a kernel matrix is always positive-semidefinite. Moreover, given a symmetric, positive semi-definite matrix, its corresponding kernel function is valid:

$$\begin{aligned} \mathbf{z}^T K \mathbf{z} &= \sum_{j=1}^m \left(\mathbf{z}_{(j)} \sum_{i=1}^m K(\mathbf{x}_i, \mathbf{x}_j) \mathbf{z}_{(i)} \right) \\ &= \sum_{j=1}^m \sum_{i=1}^m \mathbf{z}_{(j)} K(\mathbf{x}_i, \mathbf{x}_j) \mathbf{z}_{(i)} \\ &= \end{aligned}$$

3.2.3 Soft Margin Classifier

Not all data are linearly separable. Instead of requiring all data points to have $\mathbf{y} \circ (X\boldsymbol{\theta} + b) \geq 1$, we can relax this constraint by subtracting $\xi_{(i)}$ from each term while making such relaxations incur a penalty on the objective function by adding the sum of ξ :

$$\begin{aligned} \min_{\boldsymbol{\theta}, \xi} \quad & \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + C(\mathbf{1}^T \xi) \\ \text{s.t.} \quad & \mathbf{y} \circ (X\boldsymbol{\theta} + b) \geq 1 - \xi \\ & \xi \geq 0 \end{aligned}$$

Here C controls the weight of penalty slack variable incurs. When C is large, the classifier will behave more like an optimal margin classifier, ensuring most of the data points have a margin greater the "minimum" margin of the dataset. When C is small, the classifier will optimize towards a larger margin, despite many of the data pointing probably having a margin less than the "minimum" margin.

Let's rewrite the constraints of the problem to that of standard quadratic programming:

$$\begin{aligned} \min_{\boldsymbol{\theta}, \xi} \quad & \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + C(\mathbf{1}^T \xi) \\ \text{s.t.} \quad & 1 - \xi - \mathbf{y} \circ (X\boldsymbol{\theta} + b) \leq 0 \\ & -\xi \leq 0 \end{aligned}$$

The dual program of this quadratic program can be derived, again, by differentiating its Lagrangian:

$$\mathcal{L}(\boldsymbol{\theta}, \xi, \boldsymbol{\alpha}) = \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + C(\mathbf{1}^T \xi) + \boldsymbol{\alpha}^T (1 - \xi - \text{Diag}(\mathbf{y}) X \boldsymbol{\theta} - b \mathbf{y}) - \gamma^T \xi$$

where α and γ are our Lagrange multipliers, since we have two constraints; now calculate the gradients and set them to 0:

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\theta, \xi, \alpha) &= \theta - X^T \text{Diag}(\mathbf{y}) \alpha = 0 \\ \nabla_b \mathcal{L}(\theta, \xi, \alpha) &= -\alpha^T \mathbf{y} = 0 \\ \nabla_{\xi} \mathcal{L}(\theta, \xi, \alpha) &= C \mathbf{1} - \alpha - \gamma = 0\end{aligned}$$

Since $b\alpha^T \mathbf{y} = 0$ and $C(\mathbf{1}^T \xi) - \alpha^T \xi - \gamma^T \xi = 0$, we get the same solution for dual program objective:

$$\max_{\alpha} \alpha - \frac{1}{2}(\alpha \circ \mathbf{y})^T X X^T (\alpha \circ \mathbf{y})$$

and since $\gamma = C - \alpha \geq 0$, we have constraints:

$$\begin{aligned}\text{subject to } \alpha^T \mathbf{y} &= 0 \\ 0 &\leq \alpha \leq C\end{aligned}$$

The KKT dual-complementarity condition for the dual problem is:

$$\begin{aligned}\alpha \circ (1 - \xi - \mathbf{y}(X\theta + b)) &= 0 \\ \gamma \circ \xi &= 0\end{aligned}$$

We can make find cases to test if a data points conforms to these condition:

$$\begin{aligned}\alpha_{(i)} = 0 &\implies 1 - \xi_{(i)} - (\theta^T X_{(i)} + b) \leq 0 \wedge \gamma_{(i)} = C > 0 \\ &\implies 1 - \xi_{(i)} - (\theta^T X_{(i)} + b) \leq 0 \wedge \xi_{(i)} = 0 \\ &\implies (\theta^T X_{(i)} + b) \geq 1 \\ 0 < \alpha_{(i)} < C &\implies 1 - \xi_{(i)} - (\theta^T X_{(i)} + b) = 0 \wedge \gamma_{(i)} = C - \alpha_{(i)} > 0 \\ &\implies 1 - \xi_{(i)} - (\theta^T X_{(i)} + b) = 0 \wedge \xi_{(i)} = 0 \\ &\implies \theta^T X_{(i)} + b = 1 \\ \alpha_{(i)} = C &\implies 1 - \xi_{(i)} - (\theta^T X_{(i)} + b) = 0 \wedge \gamma_{(i)} = C - C = 0 \\ &\implies 1 - \xi_{(i)} - (\theta^T X_{(i)} + b) = 0 \wedge \xi_{(i)} \geq 0 \\ &\implies \theta^T X_{(i)} + b \leq 1\end{aligned}$$

3.2.4 Sequential Minimal Optimization

Sequential Minimal Optimization (SMO) is an algorithm that optimizes the dual program of SVM via coordinate ascent. Coordinate ascent optimizes a problem by optimizing each dimension of the variable sequentially, while fixing other dimensions during this process. It is an efficient algorithm for convex optimization provided the problem can be solved efficiently along each dimension. In the case SVM dual program, we optimize one **pair** of dimension a time to satisfy the constraint $\alpha^T \mathbf{y} = 0$. For example, if we select $\alpha_{(p)}, \alpha_{(q)}$ to optimize, given the constraint $\alpha^T \mathbf{y} = 0$, we have:

$$\begin{aligned}\alpha_{(p)} \mathbf{y}_{(p)} + \alpha_{(q)} \mathbf{y}_{(q)} &= - \sum_{i=1; i \notin \{p, q\}}^m \alpha_{(i)} \mathbf{y}_{(i)} = C_1 \\ \alpha_{(q)} &= \frac{C_1 - \alpha_{(p)} \mathbf{y}_{(p)}}{\mathbf{y}_{(q)}} = \mathbf{y}_{(q)} (C_1 - \alpha_{(p)} \mathbf{y}_{(p)}) \quad \mathbf{y} \in \{-1, 1\}^m\end{aligned}$$

Plugging in this formula for $\alpha_{(q)}$ and treating other variables as constants, our dual program becomes:

$$\max_{\alpha_{(p)}} \alpha_{(p)} + \mathbf{y}_{(q)} (C_1 - \alpha_{(p)} \mathbf{y}_{(p)}) - \sum_{j=1}^m M(p, j) - \sum_{j=1}^m M(q, j) - \sum_{i=1}^m M(i, p) - \sum_{i=1}^m M(i, q) + M(q, p) + M(p, q) + M(p, p) + M(q, q) + C_2$$

where $M(i, j) = \mathbf{y}_{(i)} \mathbf{y}_{(j)} \alpha_{(i)} \alpha_{(j)} K(X_{(i)}, X_{(j)})$. Since $M(i, j) = M(j, i)$, this is equivalent to following:

$$\max_{\alpha_{(p)}} \alpha_{(p)} + \mathbf{y}_{(q)} (C_1 - \alpha_{(p)} \mathbf{y}_{(p)}) - 2 \sum_{j=1}^m M(p, j) - 2 \sum_{j=1}^m M(q, j) + 2M(p, q) + M(p, p) + M(q, q) + C_2$$

$$\iff \max_{\alpha_{(p)}} \alpha_{(p)} + \mathbf{y}_{(q)}(C_1 - \alpha_{(p)}\mathbf{y}_{(p)}) - M(p, p) - 2M(p, q) - 2 \sum_{j=1; j \notin \{p, q\}}^m M(p, j) - M(q, q) - 2 \sum_{j=1; j \notin \{p, q\}}^m M(q, j) + C_2$$

This is a one dimensional quadratic function, just make its derivative 0 to maximize it:

$$\begin{aligned} & \frac{d}{d\alpha_{(p)}} \left(\alpha_{(p)} + \mathbf{y}_{(q)}(C_1 - \alpha_{(p)}\mathbf{y}_{(p)}) - M(p, p) - 2M(p, q) - 2 \sum_{j=1; j \notin \{p, q\}}^m M(p, j) - M(q, q) - 2 \sum_{j=1; j \notin \{p, q\}}^m M(q, j) + C_2 \right) \\ &= 1 - \mathbf{y}_{(p)}\mathbf{y}_{(q)} - \frac{dM(p, p)}{d\alpha_{(p)}} + 2 \frac{dM(p, q)}{d\alpha_{(p)}} + \frac{dM(q, q)}{d\alpha_{(p)}} + 2 \sum_{j=1; j \notin \{p, q\}}^m \frac{dM(p, j)}{d\alpha_{(p)}} + 2 \sum_{j=1; j \notin \{p, q\}}^m \frac{dM(q, j)}{d\alpha_{(p)}} \end{aligned}$$

where we have:

$$\begin{aligned} \frac{dM(p, p)}{d\alpha_{(p)}} &= 2\mathbf{y}_{(p)}^2 \alpha_{(p)} K(X_{(p)}, X_{(p)}) = 2\alpha_{(p)} K(X_{(p)}, X_{(p)}) \\ \frac{dM(p, q)}{d\alpha_{(p)}} &= \mathbf{y}_{(p)}(\mathbf{y}_{(q)}^2(C_1 - \alpha_{(p)}\mathbf{y}_{(p)}))K(X_{(p)}, X_{(q)}) - (\mathbf{y}_{(q)}\mathbf{y}_{(p)})\mathbf{y}_{(p)}\mathbf{y}_{(q)}\alpha_{(p)}K(X_{(p)}, X_{(q)}) \\ &= \mathbf{y}_{(p)}C_1K(X_{(p)}, X_{(q)}) - \mathbf{y}_{(p)}^2\alpha_{(p)}K(X_{(p)}, X_{(q)}) - \mathbf{y}_{(q)}^2\mathbf{y}_{(p)}^2\alpha_{(p)}K(X_{(p)}, X_{(q)}) \\ &= \mathbf{y}_{(p)}C_1K(X_{(p)}, X_{(q)}) - \alpha_{(p)}K(X_{(p)}, X_{(q)}) - \alpha_{(p)}K(X_{(p)}, X_{(q)}) \\ &= (\mathbf{y}_{(p)}C_1 - 2\alpha_{(p)})K(X_{(p)}, X_{(q)}) \\ \frac{dM(q, q)}{d\alpha_{(p)}} &= -2\mathbf{y}_{(q)}(C_1 - \alpha_{(p)}\mathbf{y}_{(p)})K(X_{(q)}, X_{(q)})\mathbf{y}_{(q)}\mathbf{y}_{(p)} = -2(C_1 - \alpha_{(p)}\mathbf{y}_{(p)})K(X_{(q)}, X_{(q)})\mathbf{y}_{(p)} \\ &= 2\alpha_{(p)}K(X_{(q)}, X_{(q)}) - 2C_1K(X_{(q)}, X_{(q)})\mathbf{y}_{(p)} \\ &= 2K(X_{(q)}, X_{(q)})(\alpha_{(p)} - C_1\mathbf{y}_{(p)}) \\ \frac{dM(p, j)}{d\alpha_{(p)}} &= \mathbf{y}_{(p)}\mathbf{y}_{(j)}\alpha_{(j)}K(X_{(p)}, X_{(j)}) \\ \frac{dM(q, j)}{d\alpha_{(p)}} &= -\mathbf{y}_{(q)}\mathbf{y}_{(j)}\alpha_{(j)}K(X_{(q)}, X_{(j)}) \end{aligned}$$

Plugging these derivatives in, we have:

$$\begin{aligned} &= 1 - \mathbf{y}_{(p)}\mathbf{y}_{(q)} - 2\alpha_{(p)}K(X_{(p)}, X_{(p)}) - 2(\mathbf{y}_{(p)}C_1 - 2\alpha_{(p)})K(X_{(p)}, X_{(q)}) - 2K(X_{(q)}, X_{(q)})(\alpha_{(p)} - C_1\mathbf{y}_{(p)}) \\ &- 2 \sum_{j=1; j \notin \{p, q\}}^m \mathbf{y}_{(p)}\mathbf{y}_{(j)}\alpha_{(j)}K(X_{(p)}, X_{(j)}) - 2 \sum_{j=1; j \notin \{p, q\}}^m \mathbf{y}_{(q)}\mathbf{y}_{(j)}\alpha_{(j)}K(X_{(q)}, X_{(j)}) = 0 \\ &(K(X_{(p)}, X_{(q)}) - 2K(X_{(p)}, X_{(p)}) - K(X_{(q)}, X_{(q)}))\alpha_{(p)} \\ &= \sum_{j=1; j \notin \{p, q\}}^m \mathbf{y}_{(j)}\alpha_{(j)}(\mathbf{y}_{(p)}K(X_{(p)}, X_{(j)}) - \mathbf{y}_{(q)}K(X_{(q)}, X_{(j)})) - \frac{1 - \mathbf{y}_{(p)}\mathbf{y}_{(q)}}{2} + \mathbf{y}_{(p)}C_1K(X_{(p)}, X_{(q)}) - 2K(X_{(q)}, X_{(q)})C_1\mathbf{y}_{(p)} \end{aligned}$$

3.2.5 Multi-class SVM (Crammer and Singer)

Let's start with an linear multi-class classifier for a dataset of $S = \{(X_{(1)}, \mathbf{y}_{(1)}), \dots, (X_{(m)}, \mathbf{y}_{(m)})\}$ where $X \in \mathbb{R}^{m \times n}$ and $\mathbf{y} \in \{1, \dots, k\}^m$:

$$h_M(\mathbf{x}) = \arg \max_{r=1}^k M_r^T \mathbf{x}$$

If we let M be a n by r matrix with M_r being its r th column, the result is just:

$$h_M(\mathbf{x}) = \arg \max M^T \mathbf{x}$$

assuming 1-based indexing and one-hot encoding, which represents \mathbf{y} as a m by k matrix Y with each row $Y_{(i)}$ corresponding to the one-hot encoded vector of $\mathbf{y}_{(i)}$. We can represent the margin $\gamma_{(r)}$ as:

$$\gamma_{(r)} = \min_{i=1}^m \frac{M_r^T X_{(i)}}{1}$$