

Foundations of Data Science 2020

Boqi Chen

December 12, 2020

Contents

1	Linear Regression	1
1.1	Linear Models	1
1.1.1	Linear Models	1
1.1.2	Learning Linear Models	2
1.1.3	Least Squares Objective Function	2
1.2	Maximum Likelihood	6
1.2.1	Maximum Likelihood Principle	6
1.2.2	Probabilistic Formulation of Linear Regression via Maximum Likelihood	7
1.3	Basis Expansion	10
1.3.1	Polynomial Basis Expansion	10
1.4	Basis Expansion with Kernels	12
1.4.1	Polynomial Kernel	12
1.4.2	Radial Basis Function Kernel	13
1.4.3	Liner Model with RBF Kernel	14
1.5	Bias-Variance Trade-off	17
1.5.1	Bias and Variance	17
1.5.2	Learning Curves	18
1.5.3	Overfitting	19
1.6	Regularization	21
1.6.1	Ridge Regression	21
1.6.2	Lasso Regression	23
1.7	Validation	26
1.7.1	Validation Error	26
1.7.2	K-Fold Cross Validation	27
1.7.3	Kaggle Leaderboard Mechanism	27
1.8	Feature Selection	29
1.8.1	Wrapper Methods	29
1.8.2	Filter Methods	29
2	Optimization	31
2.1	Convex Optimization	31
2.1.1	Convex Set	31

2.1.2	Convex Function	32
2.1.3	Convex Optimization	32
2.1.4	Classes of Convex Optimization Problem	33
2.2	Non-convex Optimization	36
2.2.1	Gradient Descent	37
2.2.2	Stochastic Gradient Descent	38
2.2.3	Sub-Gradient Descent	39
2.2.4	Constrained Convex Optimization	40
2.2.5	Second Order Methods	41
3	Classification	44
3.1	Generative Models for Classification	44
3.1.1	Naïve Bayes Classifier	46
3.1.2	Quadratic Discriminant Analysis	49
3.1.3	Linear Discriminant Analysis	51
3.1.4	Softmax Function	52
3.1.5	Two-Class Linear Decision Boundaries and Sigmoid Function	53
3.2	Logistic Regression	54
3.2.1	Binary-Class Logistic Regression	54
3.2.2	Multi-Class Logistic Regression	58
3.3	Multi-class Classification and Measuring Performance	60
3.3.1	Multi-Class Classification	60
3.3.2	Measuring Performance Binary Classification	60
3.3.3	Measuring Performance Multi-Class Classification	62
3.4	Support Vector Machines	64
3.4.1	Maximum Margin Principle	64
3.4.2	Hinge Loss Optimization	68
3.4.3	Dual SVM Formulation	69
3.4.4	Kernel Methods	73
4	Neural Network	75
4.1	Feed-forward Neural Networks	75
4.1.1	Multi-layer Perceptrons	77
4.1.2	Back-Propagation	80
4.2	Training Neural Networks	86
4.2.1	Initialising Weights and Biases	86
4.2.2	Saturation and Vanishing Gradient	86
4.2.3	Rectified Linear Unit	88
4.2.4	Overfitting and Avoiding Overfitting	89
4.3	Convolutional Neural Networks	92
4.3.1	Convolution	92
4.3.2	Max-Pooling	95
4.3.3	Some Popular Convolutional Neural Networks	95

4.3.4	Back-propagation in Convolutional Neural Networks	97
4.3.5	Number of Parameters in Convolutional Neural Networks	98
5	Unsupervised Learning	99
5.1	Clustering	99
5.1.1	k-Means Formulation of Clustering	99
5.1.2	The k-Means Algorithm	100
5.1.3	Beyond k -Means	101
5.1.4	Transforming Input Formats	101
5.1.5	Hierarchical Clustering	103
5.1.6	Spectral Clustering	106
5.2	Principal Component Analysis	108
5.2.1	PCA: Maximum Variance View	109
5.2.2	PCA: Best Reconstruction View	110
5.2.3	How Many Principal Components to Pick	114

1

Linear Regression

1.1 Linear Models

1.1.1 Linear Models

Suppose the input is a vector $\mathbf{x} \in \mathbb{R}^D$ and the output is $y \in \mathbb{R}$. We have the data $\langle (\mathbf{x}_i, y_i) \rangle_{i=1}^N$. Where D denotes the dimension of the data and N is the size of the dataset. The linear model is given in equation 1.1.

$$y = w_0 + x_1 w_1 + \cdots + x_D w_D + \epsilon \quad (1.1)$$

Where w_0 is the bias or the intercept and ϵ is the noise.

If we have the input vector augmented to 1(i.e. $\mathbf{x} = [1, x_0, \dots, x_D]$), we can simplify the linear model to the dot product of input vector \mathbf{x} and weight \mathbf{w} plus the noise given in equation 1.2.

$$y = \mathbf{x} \cdot \mathbf{w} + \epsilon \quad (1.2)$$

Let's look at a toy example here. Say if we want to predict the commute time into the city center and we have two variables: the distance to the city center and days of the week. We have the data shown in Table 1.1.

Distance (km)	Days	Commute Time (min)
2.7	Fri.	25
4.1	Mon.	33
1.0	Sun.	15
5.2	Tue.	45
2.8	Sat.	22

Table 1.1: Commute Time Data

In order to implement the linear regression, we first have to convert the non-numeric variables $\{Mon, \dots, Sun\}$ to numbers. Because in linear regression model variables have arithmetic meaning(for example $0 <$

1), we cannot simply convert these nominal(unordered) variables to integers $\{0, \dots, 6\}$. However, we can introduce dummy variables for representing these features. The most common method is One-hot Encoding. For example, in this case we can set weekdays to 0 and weekends to 1. Then we have the reorganized data $\langle (\mathbf{x}_i, y_i) \rangle_{i=1}^5$, where $\mathbf{x} \in \mathbb{R}^2, y \in \mathbb{R}$. For every data point $\mathbf{x}_i = [x_{i1}, x_{i2}]$, $x_{i1} \in \mathbb{R}(\text{distance}), x_{i2} \in \{0, 1\}(\text{weekdays or weekends})$. Now the linear regression model for commute time is given in the equation 1.3.

$$y = \mathbf{x} \cdot \mathbf{w} + \epsilon = y = w_0 + x_1 w_1 + x_2 w_2 + \epsilon \quad (1.3)$$

1.1.2 Learning Linear Models

We have the data $\langle (\mathbf{x}_i, y_i) \rangle_{i=1}^N$ where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \mathbb{R}$. Now we have to estimate the model parameter \mathbf{W} where $\mathbf{W} \in \mathbb{R}^D$ to build our linear regression model. This falls into two phases: training phase and testing/deployment phase. The training phase is to learn or estimate the parameter \mathbf{W} from the data, which is shown in figure 1.1.

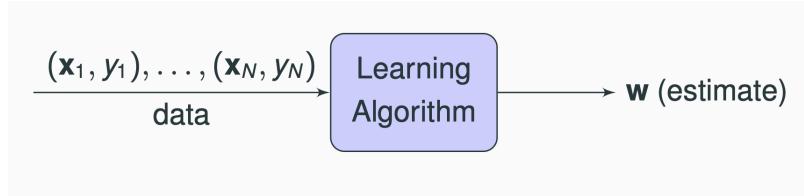


Figure 1.1: Training Phase

The testing phase is to predict the output \hat{y}_{new} given a new input \mathbf{x}_{new} .

$$\hat{y} = \mathbf{x}_{new} \cdot \mathbf{w} \quad (1.4)$$

Obviously, we care about how accurate the prediction is. In another word, we care about how different the prediction \hat{y}_{new} is from the actual observation y_{new} . Thus, we should keep some data aside for testing before deploying the model.

1.1.3 Least Squares Objective Function

We have the dataset $\langle (\mathbf{x}_i, y_i) \rangle_{i=1}^N$ where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \mathbb{R}$.

The prediction is given by:

$$\hat{y}_i = \mathbf{x}_i \cdot \mathbf{w} \quad (1.5)$$

To simplify the illustration, we here only consider the input vector of 1 dimension. For example, we predict the commute time using only the

distance from the city center.

$$\hat{y}_i = \mathbf{x}_i \cdot \mathbf{w} = w_0 + x_1 w_1 \quad (1.6)$$

Where x_1 is the distance from the city center.

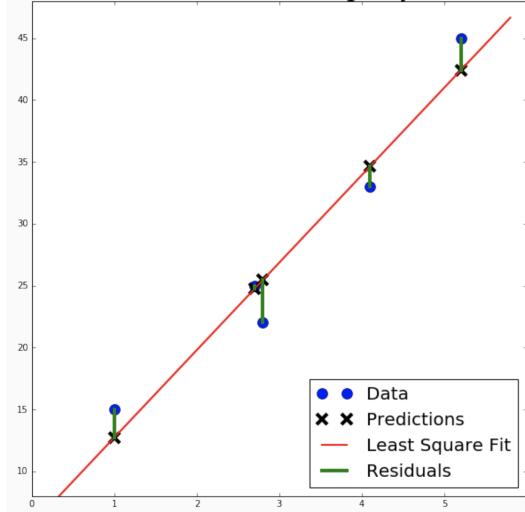


Figure 1.2: Prediction Commute Time with Distance

The performance of the model can be evaluated by the residual between the prediction and the actual observation given an input \mathbf{x}_i .

Obviously, we want to minimize the residual so that the prediction is closest possible to the actual observation. Our goal is to minimize the sum of the residuals over all data in the training dataset.

However, there exists a problem of positive and negative residual canceling each other out. Thus we often utilize the residual sum of squares(RSS) to avoid it. The residual sum of squares here is our Loss Function.

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2 = \frac{1}{2N} \sum_{i=1}^N (\mathbf{x}_i \cdot \mathbf{w} - y_i)^2 \quad (1.7)$$

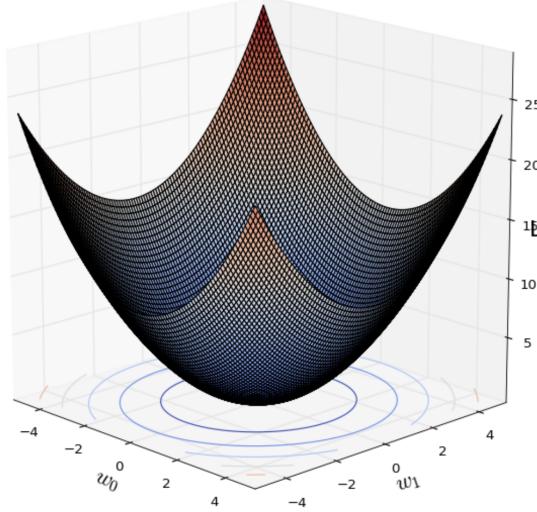


Figure 1.3: Loss Function

Thus, our goal here it to minimize the Loss Function.

We can easily obtain the solution for \mathbf{W} which minimizes the loss function by setting the partial derivative to 0 and solving the resulting system.

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w_0} = 2 \times \frac{1}{2N} \sum_{i=1}^N w_0 + x_i w_1 - y_i = 0 \\ \frac{\partial \mathcal{L}}{\partial w_1} = 2 \times \frac{1}{2N} \sum_{i=1}^N (w_0 + x_i w_1 - y_i) x_i = 0 \end{cases} \quad (1.8)$$

Simplify the equation above we obtain the Normal Equation.

$$\begin{cases} w_0 + w_1 \frac{\sum_i x_i}{N} = \frac{\sum_i y_i}{N} \\ w_0 \frac{\sum_i x_i}{N} + w_1 \frac{\sum_i x_i^2}{N} = \frac{\sum_i x_i y_i}{N} \end{cases} \quad (1.9)$$

Solving for the Normal Equation we obtain:

$$\begin{cases} w_0 = \bar{y} - w_1 \bar{x} \\ w_1 = \frac{\widehat{COV}(x,y)}{\widehat{VAR}(x)} \end{cases} \quad (1.10)$$

Where $\bar{x} = \frac{\sum_i x_i}{N}$, $\bar{y} = \frac{\sum_i y_i}{N}$, $\widehat{VAR}(x) = \frac{\sum_i x_i^2}{N} - \bar{x}^2$ and $\widehat{COV}(x,y) = \frac{\sum_i x_i y_i}{N} - \bar{x} \bar{y}$.

In terms of more general case, we can express everything in matrix notation. We have the data $\langle (\mathbf{x}_i, y_i) \rangle_{i=1}^N$, where $\mathbf{x} \in \mathbb{R}^D$, $y \in \mathbb{R}$. Then we have the data augmented to 1: $\langle (\mathbf{x}_i, y_i) \rangle_{i=1}^N$, where $\mathbf{x} \in \mathbb{R}^{D+1}$, $y \in \mathbb{R}$.

The prediction is given in the following equation.

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} \quad (1.11)$$

Here we have $\hat{\mathbf{y}} \in \mathbb{R}^{N \times 1}$, $\mathbf{X} \in \mathbb{R}^{N \times (D+1)}$ and $\mathbf{w} \in \mathbb{R}^{(D+1) \times 1}$. Then the Loss Function is given by:

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &= \frac{1}{2N} \sum_{i=1}^N (\mathbf{x}_i^T \cdot \mathbf{w} - y_i)^2 = \frac{1}{2N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{2N} (\mathbf{w}^T (\mathbf{X}^T \mathbf{X})\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y}) \\ &= \frac{1}{2N} (\mathbf{w}^T (\mathbf{X}^T \mathbf{X})\mathbf{w} - (\mathbf{y}^T \mathbf{X}\mathbf{w})^T - \mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y}) \\ &= \frac{1}{2N} (\mathbf{w}^T (\mathbf{X}^T \mathbf{X})\mathbf{w} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y})\end{aligned}\tag{1.12}$$

The gradient of the Loss Function $\nabla_{\mathbf{w}} \mathcal{L}$ is given by:

$$\begin{aligned}\nabla_{\mathbf{w}} \mathcal{L} &= \frac{1}{2N} \frac{\partial (\mathbf{w}^T (\mathbf{X}^T \mathbf{X})\mathbf{w} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y})}{\partial \mathbf{w}} \\ &= \frac{1}{2N} \left[\frac{\partial (\mathbf{w}^T (\mathbf{X}^T \mathbf{X})\mathbf{w})}{\partial \mathbf{w}} - 2 \frac{\partial (\mathbf{y}^T \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} + \frac{\partial (\mathbf{y}^T \mathbf{y})}{\partial \mathbf{w}} \right] \\ &= \frac{1}{2N} (\mathbf{X}^T \mathbf{X}\mathbf{w} + (\mathbf{X}^T \mathbf{X})^T \mathbf{w} - 2\mathbf{X}^T \mathbf{y}) \\ &= \frac{1}{N} (\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y})\end{aligned}\tag{1.13}$$

Here we use the rules of Linear Form Expressions(i.e. $\nabla_{\mathbf{w}} (\mathbf{C}^T \mathbf{w}) = \mathbf{C}$), Quadratic Expressions(i.e. $\nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{A}\mathbf{w}) = \mathbf{A}\mathbf{w} + \mathbf{A}^T \mathbf{w}$) and that the transpose of a scalar is itself(i.e. $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$).

By setting the gradient to 0,

$$\nabla_{\mathbf{w}} \mathcal{L} = \frac{1}{N} (\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y}) = 0\tag{1.14}$$

We obtain:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\tag{1.15}$$

Thus the prediction $\hat{\mathbf{y}}$ is given by:

$$\hat{\mathbf{y}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\tag{1.16}$$

$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is called the hat matrix.

In terms of the computational complexity, the computational complexity of $\mathbf{Z} = \mathbf{X}^T \mathbf{X}$, \mathbf{Z}^{-1} , $\mathbf{A} = \mathbf{X}^T \mathbf{y}$ and $\mathbf{w} = \mathbf{Z}^{-1} \mathbf{A}$ is $O(D^2 N)$, $O(D^3)$, $O(N)$ and $O(D^2)$ each respectively . Thus the overall computational complexity is $O(D^2 N + D^3)$. Moreover, if \mathbf{X} is defined by a join of several relations, the number of rows N maybe exponential in the number of relations.

$$N = O(M^{\text{number relations}})\tag{1.17}$$

Thus \mathbf{X} is sparse and can be represented in $O(M)$ space losslessly for acyclic joins which are common in practice. Then \mathbf{w} can be computed in $O(D^2M + D^3)$

There leaves only one question. When do we expect $\mathbf{X}^T \mathbf{X}$ to be invertible?

The matrix $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{(D+1) \times (D+1)}$. $\text{Rank}(\mathbf{X}^T \mathbf{X}) = \text{Rank}(\mathbf{X}) \leq \min\{D+1, N\}$. Thus it is invertible if $\text{Rank}(\mathbf{X}) = D+1$.

Recall our toy example of predicting the computing time. If we use one-hot encoding for a feature like days. $\{x_{Mon}, \dots, x_{Sun}\}$ stand for a 0-1 valued variables in the one-hot encoding. Thus we will always have $x_{Mon} + \dots + x_{Sun} = 1$, which introduces a linear dependance in the columns of \mathbf{X} reducing the rank. In this case, we cannot expect $\mathbf{X}^T \mathbf{X}$ to be invertible. We may have to drop some features to adjust the rank.

1.2 Maximum Likelihood

Recall what we have discussed in the previous section. Generally speaking ,we want to pick a model that is expected to fit our data well enough. So we choose a measure of performance of how well this model fits the data that makes "sense" and can be optimised. Then we can utilize a certain optimisation algorithm to obtain the model parameters. This is the "optimisation" view of machine learning.

In this section we will look at things in a different view - the "probabilistic" view of machine learning. In general, we want to pick a model for our data and explicitly formulate the deviation or uncertainty from the model by probability. We will use the notions from probability theory to define the suitability of various models. Then we can "find" the parameters or make predictions on unseen data utilizing these suitability criteria.

1.2.1 Maximum Likelihood Principle

Given observations x_1, \dots, x_N drawn independently from the same distribution(i.e. independently and identically distributed(i.i.d.)) p which has a parametric form, for example, the parameter vector $\boldsymbol{\theta}$. The likelihood of observing x_1, \dots, x_N equals to the probability of making these observations assuming they are generated according to p (i.e the joint probability distribution of the observations given $\boldsymbol{\theta}$: $p(x_1, \dots, x_N | \boldsymbol{\theta})$).

The *Maximum Likelihood Principle* is to pick the parameter $\boldsymbol{\theta}$ that maximise the likelihood.

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} p(x_1, \dots, x_N | \boldsymbol{\theta}) \quad (1.18)$$

Since the observations are i.i.d.,

$$p(x_1, \dots, x_N | \boldsymbol{\theta}) = \prod_{i=1}^N p(x_i | \boldsymbol{\theta}) \quad (1.19)$$

Let's look at an example here. The probability distribution of a discrete random variable X is given in Table 1.2.

X	0	1	2	3
p(X)	$\frac{2\theta}{3}$	$\frac{\theta}{3}$	$\frac{2(1-\theta)}{3}$	$\frac{1-\theta}{3}$

Table 1.2: Probability Distribution of X

If we have the observations $(3, 0, 2, 1, 3, 2, 1, 0, 2, 1)$. How can we compute the maximum likelihood estimate for θ ?

The likelihood is given by:

$$p(x_1, \dots, x_N | \boldsymbol{\theta}) = \prod_{i=1}^{10} p(x_i | \boldsymbol{\theta}) = \left(\frac{2\theta}{3}\right)^2 \left(\frac{\theta}{3}\right)^3 \left(\frac{2(1-\theta)}{3}\right)^3 \left(\frac{1-\theta}{3}\right)^2 \quad (1.20)$$

To make computation simpler, we take the Log of the likelihood and take the derivative with respect to the parameter θ .

$$\log L(\theta) = 5 \log \theta + 5 \log(1 - \theta) - 10 \log 3 \quad (1.21)$$

In order to maximize the likelihood, we set the derivative to 0 and obtain:

$$\theta = \frac{1}{2} \quad (1.22)$$

1.2.2 Probabilistic Formulation of Linear Regression via Maximum Likelihood

Recall the linear model,

$$y = \mathbf{w}\mathbf{x} + \epsilon \quad (1.23)$$

We can see it as model y given \mathbf{x}, \mathbf{W} as a random variable with mean $\mathbf{W}^T \mathbf{x}$.

$$\mathbb{E}[y | \mathbf{x}, \mathbf{w}] = \mathbf{w}^T \mathbf{x} \quad (1.24)$$

We will be specific in choosing the distribution of y given \mathbf{x}, \mathbf{w} .

If we assume that given \mathbf{x}, \mathbf{w} , y is normal with mean $\mathbf{w}^T \mathbf{x}$ and variance σ^2 , we can obtain,

$$p(y | \mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2) = \mathbf{w}^T \mathbf{x} + \mathcal{N}(0, \sigma^2) \quad (1.25)$$

Alternatively, we may view this model as linear regression with Gaussian noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

Since the input $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ is fixed, we do not need to model a distribution over \mathbf{X} .

Now we have the observed data (\mathbf{X}, \mathbf{y}) made of $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$. If we make an assumption that all data share the same variance σ , the likelihood of the observing data for model parameters \mathbf{w}, σ is given as,

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma) = p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w}, \sigma) = \prod_{i=1}^N p(y_i | \mathbf{x}_i, \mathbf{w}, \sigma) \quad (1.26)$$

According to the model $y_i \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}_i, \sigma^2)$,

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w}, \sigma) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2}} \quad (1.27)$$

Thus,

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w}, \sigma) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2} \quad (1.28)$$

With matrix notations,

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w}, \sigma) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^{\frac{N}{2}} e^{-\frac{1}{2\sigma^2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})} \quad (1.29)$$

Every parameter set (\mathbf{w}, σ^2) defines a probability distribution over observed data. Our goal here is to find parameters \mathbf{w} and σ^2 that maximise the likelihood of observing the data.

To simplify the computation, we first take the negative log-likelihood,

$$NLL(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma) = \frac{N}{2} \log 2\pi\sigma^2 + \frac{1}{2\sigma^2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (1.30)$$

Then our goal is changed to minimise the negative log-likelihood.

We first find \mathbf{W} that minimise the negative log-likelihood. Recall the objective function we use for the least squares estimate in the previous section,

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (1.31)$$

For minimising with respect to \mathbf{w} , the two objectives are the same up to a constant additive and multiplicative factor!

Thus, the solution \mathbf{w}_{ML} to find the maximum likelihood estimator is the same as the least squares estimator:

$$\mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (1.32)$$

Then the maximum likelihood estimate for σ is obtained by setting the derivative of negative log-likelihood with respect to σ^2 , given by:

$$\sigma_{ML}^2 = \frac{1}{N}(\mathbf{X}\mathbf{w}_{ML} - \mathbf{y})^T(\mathbf{X}\mathbf{w}_{ML} - \mathbf{y}) \quad (1.33)$$

Thus given training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, we can obtain the maximum likelihood estimate parameters \mathbf{w}_{ML} and σ_{ML}^2 .

On a new point \mathbf{x}_{new} , we can make prediction,

$$\hat{y} = \mathbf{w}_{ML}\mathbf{x}_{new} \quad (1.34)$$

and give confidence intervals,

$$y_{new} \sim \hat{y} + \mathcal{N}(0, \sigma^2) \quad (1.35)$$

as is shown in figure 1.4.

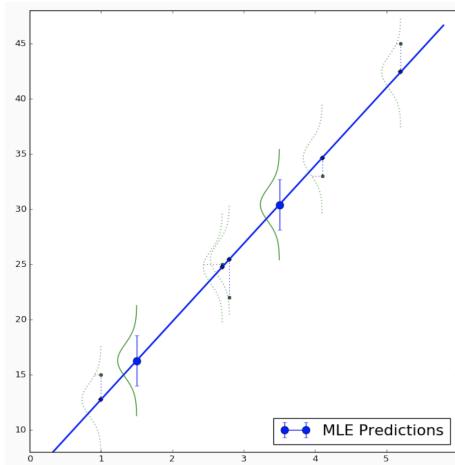


Figure 1.4: Maximum Likelihood Estimator

If the data has outliers, we can model the noise using a distribution that has heavier tails,

$$\epsilon \sim \mathcal{L}(0, b) \quad (1.36)$$

where $L(\mu, b)$ is Laplace Distribution whose density function is given by,

$$p(x) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}} \quad (1.37)$$

In this case, the likelihood of observing the data in terms of model parameters \mathbf{w} and b is given by,

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{w}, b) = \prod_{i=1}^N \frac{1}{2b} e^{-\frac{|y_i - \mathbf{w}^T \mathbf{x}_i|}{b}} = \frac{1}{(2b)^N} e^{-\frac{1}{b} \sum_{i=1}^N (|y_i - \mathbf{w}^T \mathbf{x}_i|)} \quad (1.38)$$

The negative log-likelihood is given by,

$$NLL(\mathbf{y} | \mathbf{X}, \mathbf{w}, b) = \frac{1}{b} \sum_{i=1}^N (|y_i - \mathbf{w}^T \mathbf{x}_i| + N \log(2b)) \quad (1.39)$$

Maximum likelihood estimate can be obtained by minimising the sum of the absolute values of the residuals, which is more robust to outliers in fitting a linear model. However since absolute value function is not differentiable everywhere, there is no closed-form solution for \mathbf{w}_{ML} and b_{ML} . It is also a non-linear objective function which is very hard to optimise. The good news is that it can be transformed into a linear objective subject to linear constraints which is a convex optimisation problem and has a unique solution. The details are discussed in Murphy's book at page 226.

1.3 Basis Expansion

1.3.1 Polynomial Basis Expansion

Consider the following one-dimensional data points generated by a quadratic function (red curve) with some noise shown in Figure 1.5.

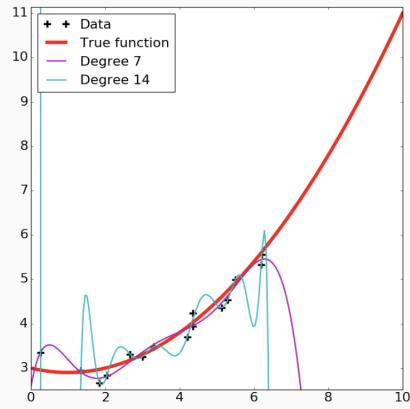


Figure 1.5: Polynomial Expansion

The linear model fitting the data points can be given by the dot product of the weight vector $\mathbf{w} = [w_0, w_1]^T$ and the feature vector $\psi(x) = [1, x]^T$

which was mapped from x by the mapping function $\psi(x)$.

$$\hat{y}_i = \mathbf{w}\psi(x_i) = w_0 + w_1x_i + \epsilon \quad (1.40)$$

Since the data points are generated by a quadratic function, the linear model with a polynomial degree of 1 is not able to fit the data very well.

Thus, we can utilize a mapping function with a higher polynomial degree of d to do a d -degree polynomial basis expansion on the input feature x . The feature vector $\psi(x)$ now is given by,

$$\psi(x) = [1, x, x^2, \dots, x^d] \quad (1.41)$$

Where $d \in \mathbb{N}$.

Thus, the linear model with a polynomial degree of d is given by,

$$\hat{y}_i = \mathbf{w}\psi(x) = [w_0, \dots, w_d][1, x, x^2, \dots, x^d]^T + \epsilon = \sum_{i=0}^d w_i x_i^d + \epsilon \quad (1.42)$$

The purple curve and the green curve in Figure 1.5 shows the linear model with a polynomial degree of 7 and 14 each respectively. As is shown in the figure, the linear models with higher polynomial degrees fit the data points better yet introduce a problem of overfitting. Using more data can avoid overfitting.

The polynomial expansion can also be performed in higher dimensions. In another word, we can also perform polynomial expansion on multi-dimensional input data. In this case, we are still fitting lineal model, but using more features. Using degree d polynomials in k dimensions results in $O(k^d)$ features.

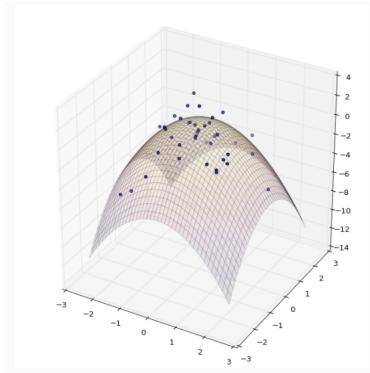


Figure 1.6: Quadratic Model $\psi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1x_2]$

1.4 Basis Expansion with Kernels

Recall the linear regression model with polynomial basis expansion. We want to minimize the loss function given by,

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2N} (\phi(\mathbf{X})\mathbf{w} - \mathbf{y})^T (\phi(\mathbf{X})\mathbf{w} - \mathbf{y}) \quad (1.43)$$

Where $\phi(\mathbf{X}) \in \mathbb{R}^{N \times O(k^d)}$ is the k degree polynomial expansion of \mathbf{X} .

Setting $\nabla_{\mathbf{w}} \mathcal{L}$ to 0 we can obtain the weight minimizing the loss function given by,

$$\mathbf{w} = (\phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^T \mathbf{y} \quad (1.44)$$

If we consider the basis expansion of just one input $\phi(\mathbf{x}) \in \mathbb{R}^{1 \times O(k^d)}$. The computational complexity of $\phi(\mathbf{x})^T \phi(\mathbf{x})$ is already $O(k^d)$, let alone we have N inputs $\phi(\mathbf{X}) \in \mathbb{R}^{N \times O(k^d)}$ and have to compute the inverse matrix $(\phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1}$. In a word, this is insanely computationally-expensive.

However, we can use kernel trick to compute $\phi(\mathbf{X})^T \phi(\mathbf{X})$ to reduce the computational complexity.

For a certain expansion ϕ , a kernel computes the dot product,

$$\kappa(\mathbf{x}', \mathbf{x}) : \chi \times \chi \rightarrow \mathbb{R} = \phi(\mathbf{x}') \phi(\mathbf{x}) \quad (1.45)$$

1.4.1 Polynomial Kernel

Polynomial kernel is given as:

$$\kappa_{poly}(\mathbf{x}', \mathbf{x}) = (\mathbf{x}' \mathbf{x} + \theta)^d \quad (1.46)$$

Assume we want to find expansion function for kernel $\kappa_{poly}(\mathbf{x}', \mathbf{x}) = (\mathbf{x}' \mathbf{x})^d$, where $\mathbf{x}', \mathbf{x} \in \mathbb{R}$.

Recall the Multinomial theorem,

$$(y_1 + \cdots + y_k)^d = \sum_{n_i \geq 0, \sum_i n_i = d} C(d; n_1 + \cdots + n_k) \prod_{j=1}^k y_j^{n_j} \quad (1.47)$$

Where $C(d; n_1 + \cdots + n_k)$ is the multinomial coefficient,

$$C(d; n_1 + \cdots + n_k) = \frac{d!}{n_1! \cdots n_k!} \quad (1.48)$$

C can be viewed as the number of ways to distribute d balls into k bins, where the j^{th} bin holds $n_j \geq 0$ balls.

Now assume $y_i = x_i x'_i$, then Equation 1.45 is given as,

$$(x', \mathbf{x})^d = \sum_{n_i \geq 0, \sum_i n_i = d} \sqrt{C(d; n_1 + \dots + n_k)} \prod_{j=1}^k x_j^{n_j} \sqrt{C(d; n_1 + \dots + n_k)} \prod_{j=1}^k x'_j^{n_j} \quad (1.49)$$

Since the dimension of $\phi_{poly}(\mathbf{x})$ and $\phi_{poly}(\mathbf{x}')$ is $O(k^d)$, the complexity of computing $\kappa_{poly}(\mathbf{x}', \mathbf{x})$ using expansion functions is $O(k^d)$. Alternatively, if we know that the dot product of expansion functions $\phi(\mathbf{x})\phi(\mathbf{x}')$ equals to $(\mathbf{x}', \mathbf{x})^d$, the computational complexity can be reduced to $O(k \log d)$ using kernel $\kappa_{poly}(\mathbf{x}', \mathbf{x})$, which is more computationally-efficient. This is called the kernel trick.

1.4.2 Radial Basis Function Kernel

A Radial Basis Function Kernel is defined as,

$$\kappa_{RBF}(\mathbf{x}, \mathbf{x}') = \frac{1}{e^{\gamma \|\mathbf{x} - \mathbf{x}'\|^2}} \in (0, 1) \quad (1.50)$$

Where the hyperparameter γ is the width and \mathbf{x}' the center of the distribution.

RBF acts as a similarity measure between \mathbf{x} and \mathbf{x}' .

When \mathbf{x} and \mathbf{x}' are far apart(i.e. large squared Euclidian distance) or γ is very large,

$$\kappa_{RBF}(\mathbf{x}, \mathbf{x}') \rightarrow 0, \gamma \|\mathbf{x} - \mathbf{x}'\|^2 \rightarrow \infty \quad (1.51)$$

When \mathbf{x} and \mathbf{x}' are close(i.e. small squared Euclidian distance) or γ is very small,

$$\kappa_{RBF}(\mathbf{x}, \mathbf{x}') \rightarrow 1, \gamma \|\mathbf{x} - \mathbf{x}'\|^2 \rightarrow 0 \quad (1.52)$$

The question now is to find ϕ_{RBF} such that,

$$\kappa_{RBF}(\mathbf{x}', \mathbf{x}) = \frac{1}{e^{\gamma \|\mathbf{x} - \mathbf{x}'\|^2}} = \phi_{RBF}(\mathbf{x})\phi_{RBF}(\mathbf{x}') \quad (1.53)$$

The Taylor expansion of $e^{\mathbf{x}\mathbf{x}'}$ is given by,

$$e^{\mathbf{x}\mathbf{x}'} = \sum_{k=0}^{\infty} \frac{1}{k!} (\mathbf{x}\mathbf{x}')^k = \phi_{poly}(\mathbf{x})\phi_{poly}(\mathbf{x}') \quad (1.54)$$

Without loss of generality, we assume $\gamma = 1$, then,

$$e^{\gamma \|\mathbf{x} - \mathbf{x}'\|^2} = e^{-(\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2 - 2\mathbf{x}\mathbf{x}')} = \sum_{k=0}^{\infty} \left(\sqrt{\frac{1}{k!}} (e^{-\|\mathbf{x}\|^2}) \phi_{poly}(\mathbf{x}) \left(\sqrt{\frac{1}{k!}} (e^{-\|\mathbf{x}'\|^2}) \phi_{poly}(\mathbf{x}') \right) \right)$$

(1.55)

Thus, we can see that $\phi_{RBF} : \mathbb{R}^d \rightarrow \mathbb{R}^\infty$ projects vectors into an infinite dimensional space. Therefore, it is not feasible to compute $\kappa_{RBF}(\mathbf{x}', \mathbf{x})$ using ϕ_{RBF} .

1.4.3 Liner Model with RBF Kernel

Recall RBF kernel:

$$\kappa_{RBF}(\mathbf{x}', \mathbf{x}) == \frac{1}{e^{\gamma\|\mathbf{x}-\mathbf{x}'\|^2}} \quad (1.56)$$

If we choose different centers μ_1, \dots, μ_M , we can construct a feature map $\psi(\mathbf{x}) = [1, \kappa_{RBF}(\mu_1, \mathbf{x}), \dots, \kappa_{RBF}(\mu_M, \mathbf{x})]^T$. Then the linear model is given as,

$$\hat{y} = \mathbf{w}\psi(\mathbf{x}) + \epsilon \quad (1.57)$$

Figure 1.7 shows liner model with RBF kernel when input data is one-dimensional.

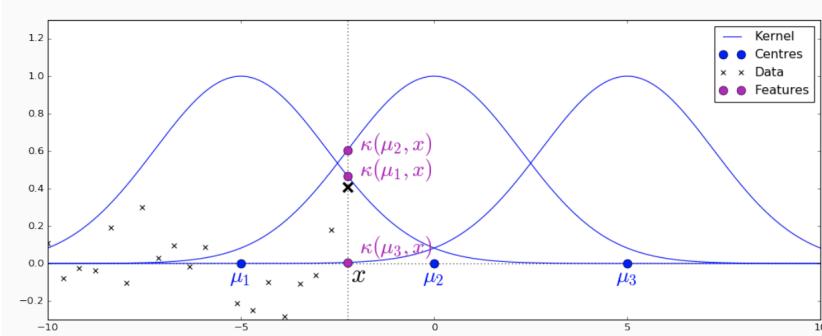


Figure 1.7: Liner Model with RBF Kernel

For hyperparameter μ , a reasonable choice for centers are data points themselves.

For hyperparameter γ , overfitting or underfitting may occur if width parameter is not chosen properly. More specifically, overfitting occurs when the kernel is too narrow(i.e. γ is very large) and vice versa.

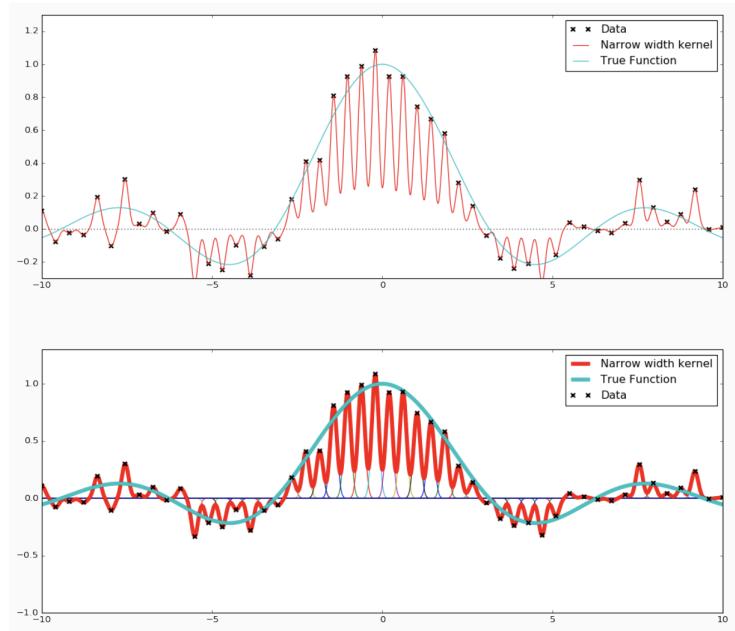
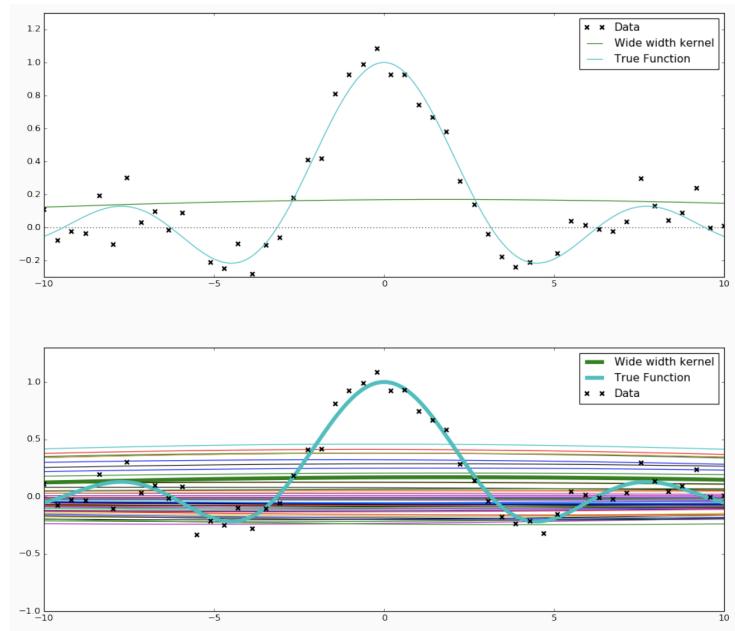
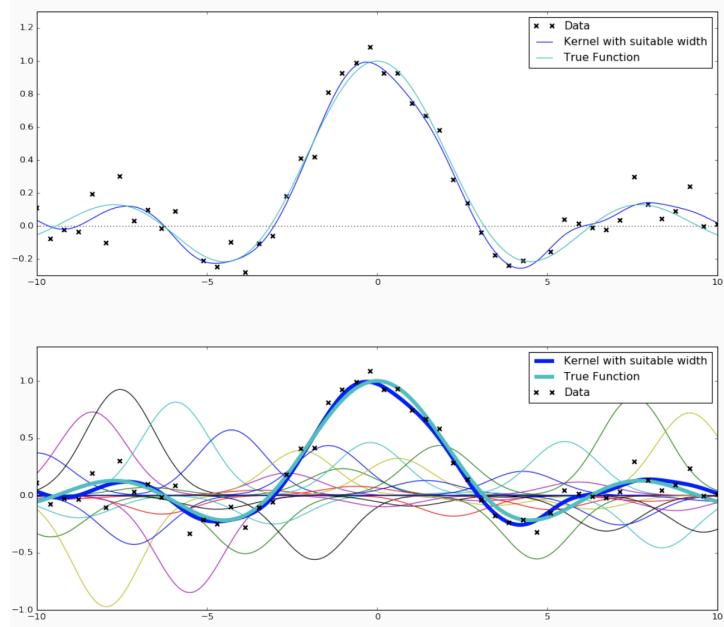
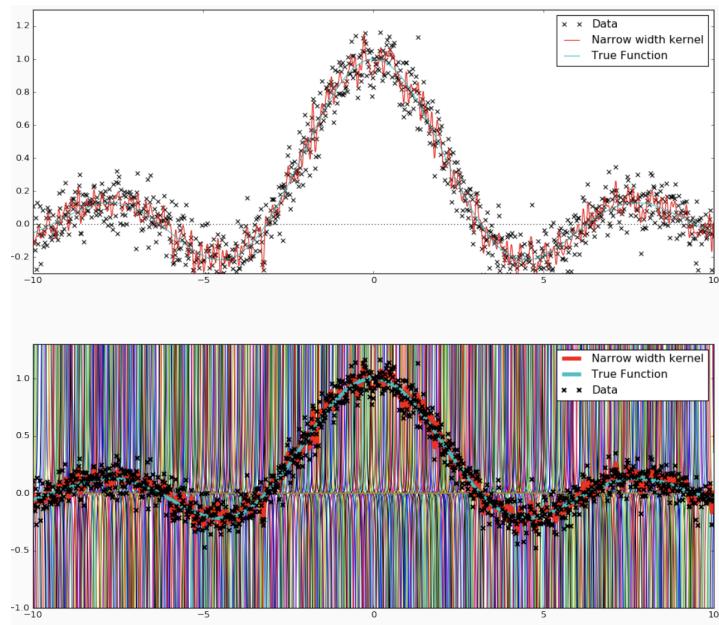
Figure 1.8: Narrow Kernel with Large γ Figure 1.9: Wide Kernel with Small γ

Figure 1.10 shows the model when kernerl is well-chose.

Figure 1.10: Wide Kernel with Well-chosen γ

Having more data can help reduce the overfitting if the kernel is too narrow.

Figure 1.11: Narrow Kernel with Large γ , Big Data

However, extra data does not help reduce underfitting when the kernel

is too wide.

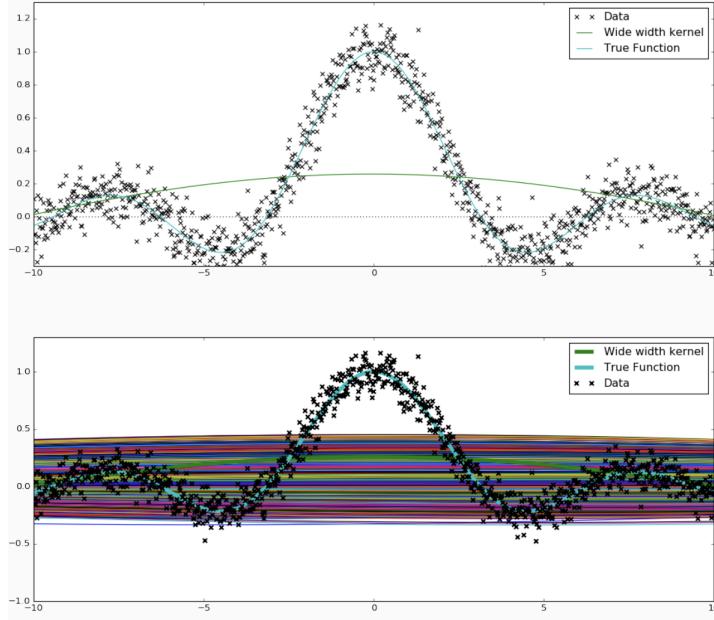


Figure 1.12: Wide Kernel with Small γ , Big Data

1.5 Bias-Variance Trade-off

1.5.1 Bias and Variance

The fundamental goal of machine learning is to generalize beyond the examples in the training data. Thus, test data is highly important.

Also, data alone is not enough. Every learner must embody knowledge/assumptions beyond the data. No learner can beat random guessing over all possible functions to be learned(Wolpert's "No Free Lunch"). Our hope is that the functions to learn in the real world are not drawn uniformly from the set of all mathematically possible function. The reasonable assumptions behind machine learning's success is that similar examples have similar classes and limited dependence or limited complexity.

Bias is the tendency to consistently learn the same wrong thing. For example, linear learner has high bias. When the frontier between two classes is not a hyperplane the learner is unable to induce it. However, decision trees do not have this problem because they can represent any Boolean function. Having high bias means underfitting.

Variance is the tendency to learn random things irrespective of the real signal. For example, decision trees suffer from high variance because when learned on different training sets generated by the same phenomenon they

are often very different, when in fact they should be the same. Having high variance means overfitting.

Thus, more powerful learners are not necessarily better than less powerful ones.

Generalization error is the combination of bias and variance,

$$\text{Generalization Error} = \text{Bias} + \text{Variance} \quad (1.58)$$

Figure 1.13 shows the example of high bias and high variance.

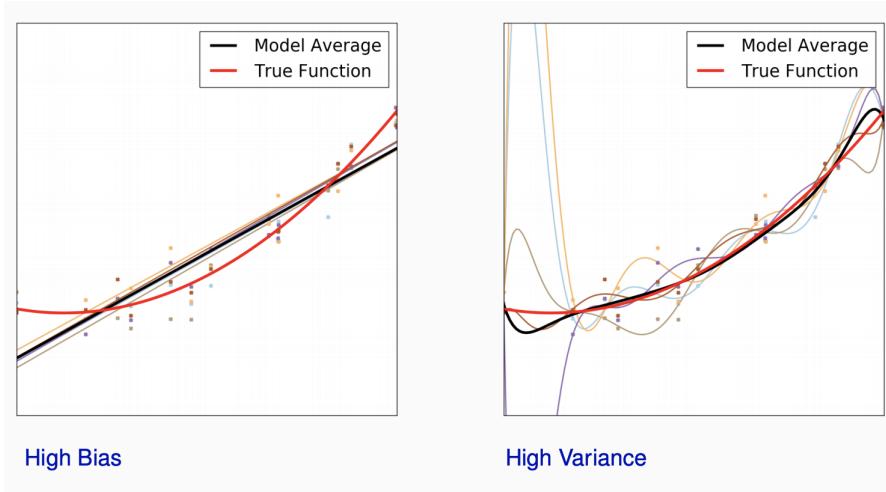


Figure 1.13: Experiment: 1-D Points Generated by a Quadratic Function

In fact, the term bias and variance are precisely defined in statistical notions. The description is in Section 5.4 in the Book *Deep Learning* by Goodfellow, Ian and et al.

1.5.2 Learning Curves

Suppose we have trained a model and used it to make predictions. In reality, these predictions are usually very poor. The question is, how can we know whether we have high bias(underfitting) or high variance(overfitting) or neither. This is very important because it helps us to decide in which way we can improve our model. Should we add more features(e.g. higher degree polynomials, narrower kernels, etc.) to make the model more complex? Or should we simplify the model(e.g. lower degree polynomials, wider kernels, etc.) to reduce the number of parameters? Or if neither of the above works, should we try and obtain more data (which comes with higher computational and monetary cost)?

We can split the data into a training set and testing set. Then we train on increasing sizes of the data and plot the training error and test error as

a function of training data size. The curve obtained is called the learning curve.

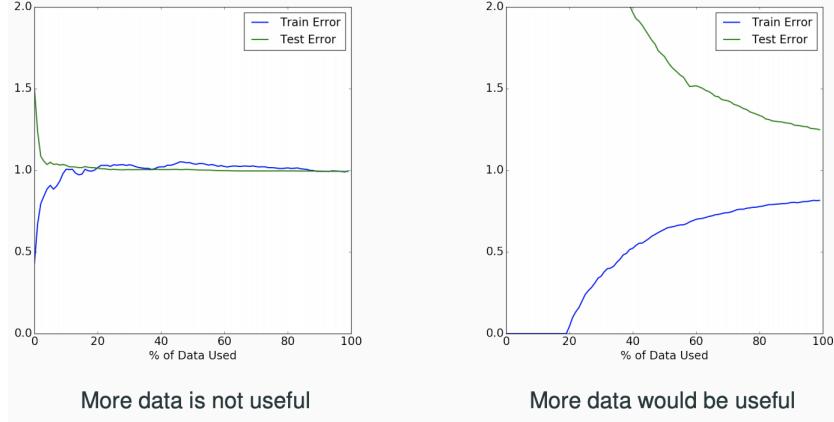


Figure 1.14: Learning Curve

1.5.3 Overfitting

When dealing with high-dimensional data(which may be caused by basis expansion) even for a linear model we have many parameters. For example with $D = 100$ input variables and using degree 10 polynomial expansion we have approximately 10^{20} parameters!

Suppose we have the data $\mathbf{X} \in \mathbb{R}^{100 \times 100}$ with every entry of \mathbf{X} is drawn from $\mathcal{N}(0, 1)$. Let $y_i = x_{i,1} + \mathcal{N}(0, \sigma^2)$ for $\sigma = 0.2$. Thus, the label y_i only depends on the first feature.

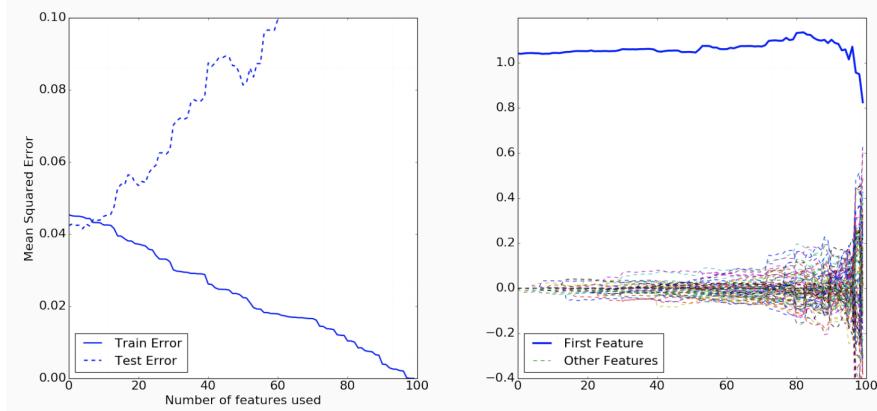


Figure 1.15: Overfitting

As is shown in the Figure 1.15, the more features are used during the training, the larger the testing error will be. In fact, we can see from the

left figure that when the number of features used is relatively small almost all the weight is on the first feature while when the the number of features used is relatively large the first feature becomes less relevant and the other features starts to fit the noise.

There are several ways to combat overfitting. For example, we can add regularization term to the evaluation function. Thus, it will penalise models with more structure and favor smaller models with less room to overfit. Also, we can do cross-validation or statistical significance test like chi-square before adding new structure to decide whether the prediction is different without this structure.

1.6 Regularization

Suppose we have data $\mathbf{X} \in \mathbb{R}^{N \times D}$, where $D \gg N$. Our idea to avoid overfitting is to add a penalty term for weights.

1.6.1 Ridge Regression

Recall the Least Square Estimate Objective,

$$\mathcal{L}(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (1.59)$$

Ridge regression Objective is given by,

$$\mathcal{L}(\mathbf{w})_{Ridge} = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \sum_{i=1}^D w_i^2 \quad (1.60)$$

We add a penalty term $\lambda \sum_{i=1}^D w_i^2$ for all weights except w_0 (since the effect of varying w_0 is output translation and not on model complexity) to control the model complexity.

However, since the translating and scaling inputs contribute to the model complexity, we should try and avoid dependence on scaling and translation of variable by standardize the input.

Feature standardization is a method used to normalize the range of independent variables or features of data. It is also known as data normalization and generally performed during the data preprocessing step.

Feature standardization makes the values of each feature in the data have zero-mean and unit-variance(i.e values drawn from $\mathcal{N}(0, 1)$).

The general method of calculation is to determine the distribution mean and standard deviation for each feature given by,

$$x' = \frac{x - \bar{x}}{\sigma} \quad (1.61)$$

Where \bar{x} is the mean and σ is the standard deviation.

Suppose the data $\mathbf{X} \in \mathbb{R}^{N \times D}$ with inputs standardized and output centered, the loss function of Ridge regression is given by,

$$\mathcal{L}(\mathbf{w})_{Ridge} = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \quad (1.62)$$

The gradient of the loss function is given by,

$$\nabla_{\mathbf{w}} \mathcal{L}_{Ridge} = 2((\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D) \mathbf{w} - \mathbf{X}^T \mathbf{y}) \quad (1.63)$$

Set the gradient to $\mathbf{0}$ and solve for \mathbf{X} , we obtain,

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (1.64)$$

Thus,

$$\mathbf{w}_{Ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} \quad (1.65)$$

Alternatively, we can consider the Ridge regression as a constrained optimization problem,

$$\arg \max_{\mathbf{w}} \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \text{ s.t. } \mathbf{w}^T \mathbf{w} \leq R$$

Then the former loss function $\mathcal{L}(\mathbf{w})_{Ridge}$ is the Lagrangian formulation of this optimization problem. Thus, the solution to Ridge regression can be viewed as a function of λ and R , as is shown in Figure 1.16.

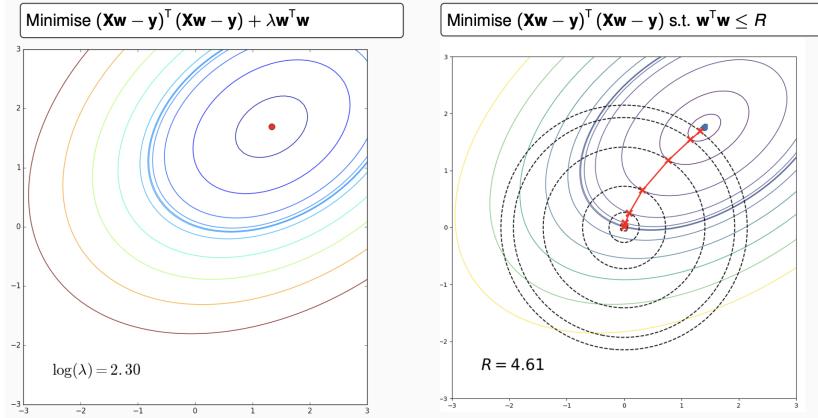


Figure 1.16: Ridge Regression v.s. Constrained Optimization

From Figure 1.16 we can see that decreasing λ equals to loose the constraint(i.e. increase R). If the constraint R is well-chosen, we will derive the same solution \mathbf{w} for both problems.

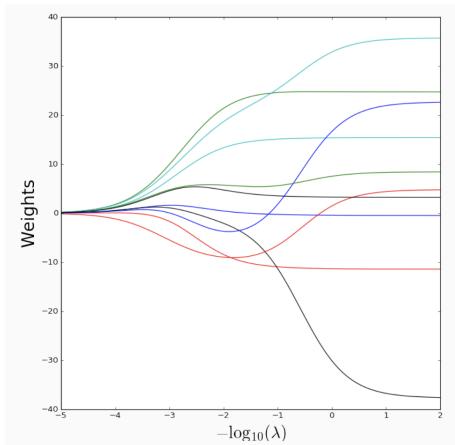


Figure 1.17: Magnitude of Weights v.s. $\log \lambda$

1.6.2 Lasso Regression

Recall the objective of Ridge regression,

$$\mathcal{L}(\mathbf{w})_{Ridge} = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda\mathbf{w}^T\mathbf{w} \quad (1.66)$$

In Ridge regression, in addition to the residual sum of the squares we penalise the sum of squares of the weights. This is called l_2 regularization or weight-decay. By adding it, we will encourage fitting signal rather than just noise by penalising weights.

Similarly, we can do l_1 regularization and obtain the Lasso objective given by,

$$\mathcal{L}(\mathbf{w})_{Lasso} = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda\|\mathbf{w}\|_1 \quad (1.67)$$

Thus,

$$\mathcal{L}(\mathbf{w})_{Lasso} = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \sum_{i=1}^D |w_i| \quad (1.68)$$

As with Ridge regression, there is a penalty on the weights. However, the absolute value function does not allow for a simple closed-form solution.

The Lasso regression can also be seen as a optimization problem,

$$\arg \max_{\mathbf{w}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) \text{ s.t. } \sum_{i=1}^D |w_i| \leq R$$

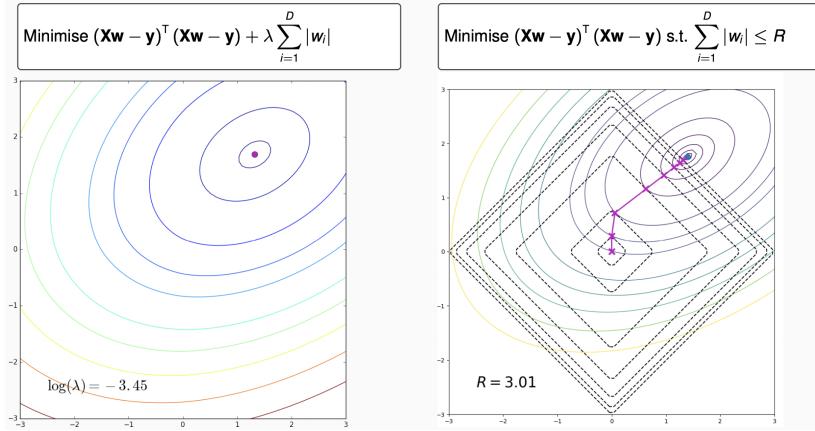
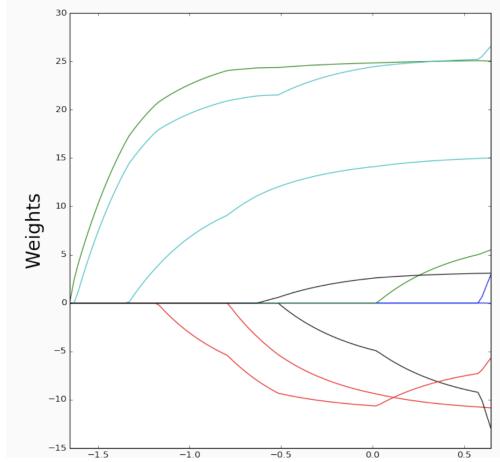


Figure 1.18: Lasso Regression v.s. Constrained Optimization

Figure 1.19: Magnitude of Weights v.s. $\log \lambda$

As is shown in Figure 1.19, in Lasso regression, weights are often exactly 0. Thus, Lasso gives a sparse model.

Recall our experiment. Suppose we have the data $\mathbf{X} \in \mathbb{R}^{100 \times 100}$ with every entry of \mathbf{X} is drawn from $\mathcal{N}(0, 1)$. Let $y_i = x_{i,1} + \mathcal{N}(0, \sigma^2)$ for $\sigma = 0.2$. We have discussed the result of simple linear regression with no regularization before. Now let's look at some results of Ridge regression and Lasso regression.

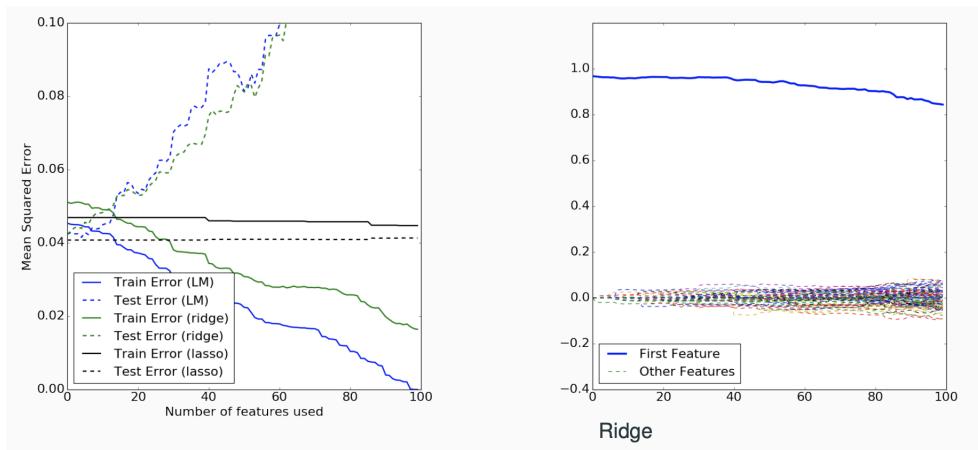


Figure 1.20: Ridge Regression Experiment

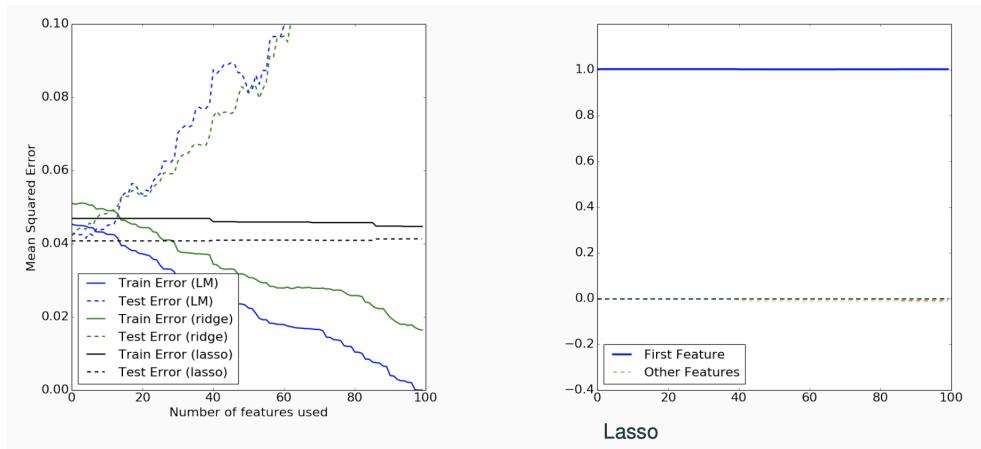


Figure 1.21: Lasso Regression Experiment

1.7 Validation

So far, we have discussed how to estimate the weight parameters \mathbf{w} . For Ridge and Lasso regression, we need to choose λ . If we perform basis expansion, we have to pick the width parameter γ and for polynomial expansion we have to choose polynomial degree d . For more complex models there may be even more hyper-parameters to choose.

Then the question is, how can we choose hyper-parameters so that our model works well. This can be done by validation.

1.7.1 Validation Error

If we have the data $\mathbf{X} \in \mathbb{R}^{D \times N}$, we can divide the dataset into three parts, namely training, validation and testing data.

Thus, for training stage, we can train the model using the training set and evaluate on the validation set. In this case, we will not only get a training error but also a validation error. If we have hyper-parameters to choose, for example λ , we can pick the value of λ that minimises the validation error.

Suppose we have the training data $\mathbf{X} \in \mathbb{R}^{2000 \times 2000}$. If we divide the data into training set and validation set (typically 80% for training and 20% for validation) and plot the training error and validation error v.s. λ for Lasso regression, we will get training and validation curves as is shown in Figure 1.22.

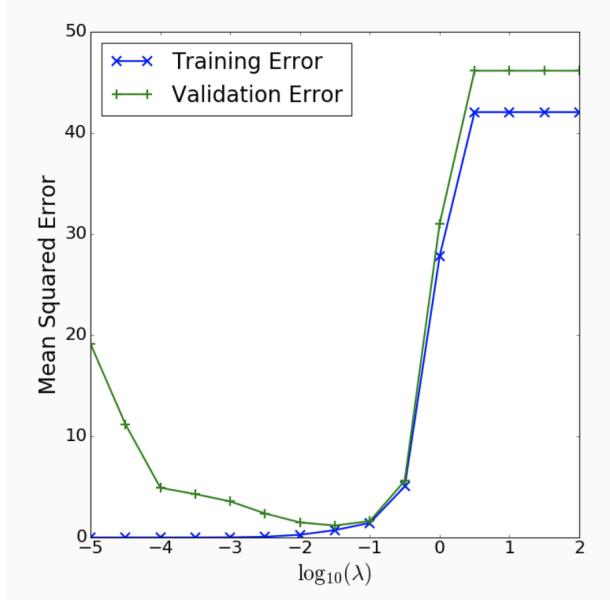


Figure 1.22: Training and Validation Curves

We can see from the U-shaped validation error curve that the left half

is overfitting(i.e. low training error and high validation error) and the right half is underfitting(i.e. high training error and high validation error). Thus, we should pick the hyper-parameter λ at the bottom of the curve.

What if we have many hyper-parameters?

The trivial approach is grid search. We first assume a small domain D_i for each hyper-parameter λ_i where $1 \leq \lambda_i \leq k$. Next we iterate over all possible combinations of hyper-parameter value $D_1 \times \dots \times D_k$. Then we perform cross-validation for each such combination and pick the combination with the lowest validation search.

1.7.2 K-Fold Cross Validation

When the data is scarce, instead of splitting it into training and validation set, we divide the data into k folds(parts).

We use $k - 1$ folds for training and 1 fold as validation. When k equals to the number of data points, it is called Leave and Out Cross Validation(LOOCV).

The validation error for fixed hyper-parameter values is the average of all runs.

Run 1	train	train	train	train	valid
Run 2	train	train	train	valid	train
Run 3	train	train	valid	train	train
Run 4	train	valid	train	train	train
Run 5	valid	train	train	train	train

Figure 1.23: K-Fold Cross Validation

1.7.3 Kaggle Leaderboard Mechanism

Suppose we do all the things right. We train the model on the training set and choose hyper-parameters using k-fold validation. Then when we test on the test set(real world), the error is still somehow unacceptably high. This suggests that we might overfit on the validation set.

Kaggle competitions use a leaderboard mechanism, which is a classic hold-out method. The training set is publicly available. The hold-out set(which is split randomly into 30% hold-out labels and 70% test labels) is also publicly available but without the labels.

The valid submission is the list of predicted labels, one for each point in the testing set. These is also a score function such as misclassification rate to evaluate the submission.

The public ranking is given by the score on the hold-out labels (without the test labels) and the private ranking is given by the score on the test labels. The team can submit repeatedly for public ranking for quite some time before the final private ranking.

Suppose the task is to predict N binary labels, i.e. $\mathbf{y} \in \{0, 1\}^N$. We can achieve a high score on public ranking by so-called Wacky Boosting Algorithm.

We firstly pick $\mathbf{y}_1, \dots, \mathbf{y}_k \in \{0, 1\}^N$ randomly. Next we select $\mathbf{Y} = \{\mathbf{y}_i | \text{score}(\mathbf{y}_i) \leq 0.5\}^N$. Then we simply output $\hat{\mathbf{y}} = \text{majority}(\mathbf{Y})$.

The problem is that there is a statistical dependence the hold-out data and the submission. Thus, submission may incorporate information about the hold-out labels realised through the leaderboard. Due to this feedback loop, the public score is no longer an unbiased estimator of the score. Eventually, the submission will overfit the hold-out set. This is an example of disconnect between theory f static data analysis and practice of interactive data analysis.

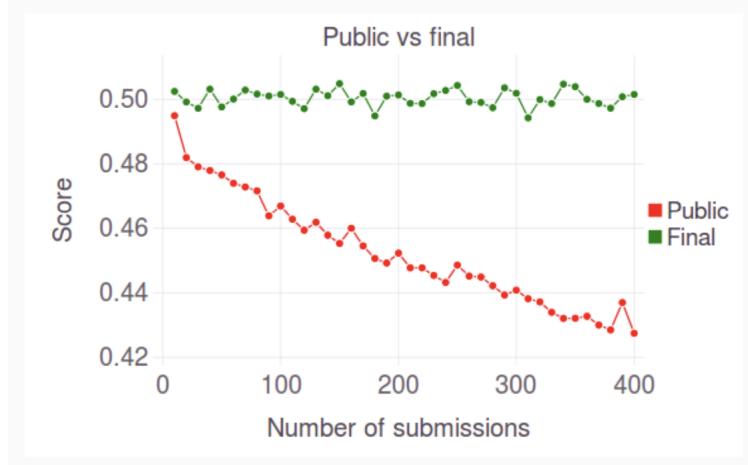


Figure 1.24: Public Score v.s. Final Score

How can we deal this problem? We can limit the rate of re-submissions, control the bit precision numbers or have winners determined on a separate test set.

1.8 Feature Selection

The goal of feature selection is to select features relevant to the model to learn given a set of possible features.

In practice, data usually contains redundant or irrelevant features which can be removed without causing information loss. Also, relevant features maybe redundant in the presence of another relevant feature(i.e. the two features are strongly correlated).

If we remove these redundant features, we can enhance the generalization by reducing overfitting, avoid the curse of dimensionality and shorten the training time. Also, simpler models are easier to interpret by the users.

1.8.1 Wrapper Methods

Wrapper method is to train a new model for each feature subset. The score is given by the count of number of mistakes on hold-out set. It is expensive yet usually provides the best performing feature set.

Let's look at an example of wrapper methods - Forward Step-wise Selection.

For n features, we ask the following sequence of n questions($1 < i \leq n$):

Q_1 : What is the best 1-feature model? Let the chosen feature be f_1 .

\vdots

Q_i : What is the best i -feature model that also has the previous selected features f_1, \dots, f_{i-1} ? Let the newly chosen feature be f_i .

The output is the best seen k -feature model, for $1 < k \leq n$.

The total number of models to train and test is $O(n^2)$ because at step i , it trains and tests $n - i + 1$ models. For linear regression, it is the same complexity as building one model.

1.8.2 Filter Methods

Filter methods use proxy measure as score instead of the actual test score. It exposes relationship between different features. Typical scores are mutual information and Pearson correlation coefficient. Filter methods are fast to compute yet results feature set not tunes to a specific type of models.

The mutual information for two random variables X and Y is given by,

$$I(X, Y) = \sum_x \sum_y p(X = x, Y = y) \log \frac{p(X = x, Y = y)}{p(X = x)p(Y = y)} \quad (1.69)$$

The probabilities can be empirically obtained from the training set. For example, $p(X = x)$ is the fraction of the number of samples with $X = x$ over the number of all samples. For variables with continuous domain, we first have to discretise their domains.

The feature selection approach via mutual information is to first compute the mutual information for each feature and the label(target) and then only keep the features that provide information about the output. In this approach, we get the ranking of features instead of finding the best features. The cutting-off point is obtained by using cross-validation.

Covariance for two random variables X and Y measures how the random variables change jointly,

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}(X))(Y - \mathbb{E}(Y)))] \quad (1.70)$$

The feature selection can also be done via (Pearson) correlation coefficient.

The (Pearson) correlation coefficient normalize the covariance of variables X and Y to give a value between -1 and 1 .

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{cov}(X, X)\text{cov}(Y, Y)}} \quad (1.71)$$

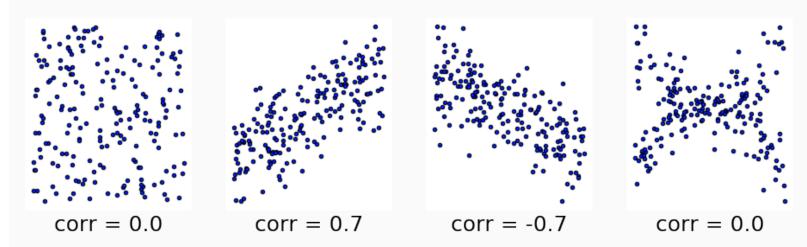


Figure 1.25: Correlation

Note that independent variables are uncorrelated, but uncorrelated variables are not necessarily independent.

2

Optimization

So far we have discussed machine learning problems with closed-form solutions such as minimization of least squares and ridge regression objectives.

However, this is not always the case. Most interesting learning problems do not admit closed-form solutions.

To solve the problems beyond the closed-form solutions, two approaches can be utilized. We can frame the objective of machine learning problem as a mathematical problem and when objectives can be formulated as convex optimization problem we can use existing black-box solver for such problems.

Also, we can use gradient-based optimization methods. These methods are not black-box and thus optimization hyper-parameters affect the performance.

In fact, most machine learning problems can be cast as optimization problems.

2.1 Convex Optimization

2.1.1 Convex Set

Definition 2.1 A set $C \in \mathbb{R}^D$ is convex if for any given $\mathbf{x}, \mathbf{y} \in C$ and $\lambda \in [0, 1]$, it holds that $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in C$.

Proposition 2.2 Given convex sets C_1, \dots, C_n , the set $\cap_{i=1}^n C_i$ is convex.

Proposition 2.3 For any L-norm $\|\cdot\|$, the set $B = \{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\| \leq 1\}$ is convex.

Proof Take $\mathbf{x}, \mathbf{y} = \{\mathbf{x} \in \mathbb{R}^D : \|\mathbf{x}\| \leq 1\}$, to show $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in B$

$$\|\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}\| \leq \lambda\|\mathbf{x}\| + (1 - \lambda)\|\mathbf{y}\| \leq 1 \quad (2.1)$$

□

Proposition 2.4 Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, the polyhedron $P = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{Ax} \leq \mathbf{b}\}$ is convex.

Proof Take $\mathbf{x}, \mathbf{y} \in P$, to show $\mathbf{A}(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) \in P$

$$\mathbf{A}(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) = \lambda\mathbf{Ax} + (1 - \lambda)\mathbf{Ay} \leq \lambda\mathbf{b} + (1 - \lambda)\mathbf{b} \leq \mathbf{b} \quad (2.2)$$

□

Proposition 2.5 The set of positive semi-definite(PSD) matrices is convex.

Proof Take $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{D \times D}$ are PSD, to show $\lambda\mathbf{A} + (1 - \lambda)\mathbf{B}$ is PSD.

$$\mathbf{x}^T(\lambda\mathbf{A} + (1 - \lambda)\mathbf{B})\mathbf{x} = \lambda\mathbf{x}^T\mathbf{Ax} + (1 - \lambda)\mathbf{x}^T\mathbf{Bx} \geq 0 \quad (2.3)$$

□

2.1.2 Convex Function

Definition 2.6 A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ defined on a convex domain is convex if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$ where f is defined and $0 \leq \lambda \leq 1$,

$$f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \quad (2.4)$$

Proposition 2.7 Affine functions: $f(\mathbf{x}) = \mathbf{b}^T \cdot \mathbf{x} + c$ are convex functions.

Proposition 2.8 Quadratic functions: $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \cdot \mathbf{A} \cdot \mathbf{x} + \mathbf{b}^T \cdot \mathbf{x} + c$ are convex functions.

Proposition 2.9 Nonnegative weighted sums of convex functions: Given convex functions f_1, \dots, f_n and $w_1, \dots, w_n \in \mathbb{R} \leq 0$, the following is a convex function

$$f(\mathbf{x}) = \sum_{i=1}^k w_i \cdot f_i(\mathbf{x}) \quad (2.5)$$

Proposition 2.10 Norms: $\|\cdot\|_p$ are convex functions except for $p=0$

2.1.3 Convex Optimization

Given convex functions $f(\mathbf{x}), g_1(\mathbf{x}), \dots, g_m(\mathbf{x})$ and affine functions $h_1(\mathbf{x}), \dots, h_n(\mathbf{x})$, a convex optimization problem is of the form:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{s.t.} && g_i(\mathbf{x}) \leq 0 \\ & && h_j(\mathbf{x}) = 0 \end{aligned} \quad (2.6)$$

The goal is to find an optimal value of a convex problem:

$$v^* = \min\{f(\mathbf{x}) : g_i(\mathbf{x}) \leq 0, i \in \{1, \dots, m\}, h_j(\mathbf{x}), j \in \{0, \dots, n\}\} \quad (2.7)$$

Whenever $f(\mathbf{x}^*) = v^*$ then \mathbf{x}^* is a (not necessarily unique) optimal point.

Note that for infeasible instances(feasible: fulfils all constraints g_i and h_i), $v^* \stackrel{\text{def}}{=} +\infty$; for unbounded instances(unbounded: the set of feasible instances has no infimum), $v^* \stackrel{\text{def}}{=} -\infty$.

\mathbf{x} is locally optimal if:

- \mathbf{x} is feasible and,
- There is $B > 0$ s.t. $f(\mathbf{x}) \leq f(\mathbf{y})$ for all feasible \mathbf{y} with $\|\mathbf{x} - \mathbf{y}\|_2 \leq B$.

\mathbf{x} is globally optimal if:

- \mathbf{x} is feasible and,
- $f(\mathbf{x}) \leq f(\mathbf{y})$ for all feasible \mathbf{y} .

Theorem 2.11 *For a convex optimization problem, all locally optimal points are globally optimal.*

2.1.4 Classes of Convex Optimization Problem

Linear Programming

Linear programming is to look for solutions $\mathbf{x} \in \mathbb{R}^n$ to the following optimization problem:

$$\begin{aligned} & \text{minimize } \mathbf{c}^T \mathbf{x} + d \\ & \text{s.t. } \mathbf{A}\mathbf{x} \leq \mathbf{e} \\ & \quad \mathbf{B}\mathbf{x} = \mathbf{f} \end{aligned} \quad (2.8)$$

The interesting thing of linear programming is that the optimal is usually achieved at the vertices of the constraint polytope.

We are not going to the algorithms of how to solve linear programming problems. However, efficient algorithms exist, both in theory and practice, for tens of thousands of variables.

Let's look at an example of linear model with absolute loss. Suppose we have data (\mathbf{X}, \mathbf{y}) and we want to minimise the following objective:

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^N |\mathbf{x}_i^T \mathbf{w} - y_i| \quad (2.9)$$

We would like to transform this optimisation problem into a linear program. We introduce one ζ for each data point. The linear program in the $D + N$ variables $w_1, \dots, w_D, \zeta_i, \dots, \zeta_N$ is given as,

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^N \zeta_i \\ \text{s.t. } & \mathbf{x}_i^T \mathbf{w} - y_i \leq \zeta_i, \quad i = 1, \dots, N \\ & y_i - \mathbf{x}_i^T \mathbf{w} \leq \zeta_i, \quad i = 1, \dots, N \end{aligned} \tag{2.10}$$

Since a solution of this program always exists, we let $\mathbf{w} \in \mathbb{R}^D$ be any vector and $\zeta_i = |\mathbf{x}_i^T \mathbf{w} - y_i|$. Then, $w_1, \dots, w_D, \zeta_1, \dots, \zeta_N$ represent a feasible solution yet not necessarily optimal to the linear program.

Let \mathbf{w}^*, ζ^* be the optimal solution. Clearly, $\zeta^* \geq |\mathbf{x}_i^T \mathbf{w} - y_i|$ due to our two constraints. This also suggests that $\zeta^* \geq 0$.

Since our linear program minimize the object $\sum_{i=1}^N \zeta_i^*$, we take the ζ_i^* to be the smallest possible, i.e. $\zeta_i^* = |\mathbf{x}_i^T \mathbf{w}^* - y_i|$. Thus, \mathbf{w}^* must be the optimal solution for \mathcal{L} . That is, the solution of this linear programming problem gives \mathbf{w} that minimizes the objective \mathcal{L} .

Quadratically Constrained Quadratic Programming

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \mathbf{x}^T \mathbf{B} \mathbf{x} + \mathbf{c}^T \mathbf{x} + d \\ \text{s.t. } & \frac{1}{2} \mathbf{x}^T \mathbf{Q}_i \mathbf{x} + \mathbf{r}_i^T \mathbf{x} + s_i \leq 0 \\ & \mathbf{A} \mathbf{x} = \mathbf{B} \end{aligned} \tag{2.11}$$

Recall the negative log-likelihood objective we discussed in the previous chapter,

$$NLL(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma) = \frac{N}{2} \log 2\pi\sigma^2 + \frac{1}{2\sigma^2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \tag{2.12}$$

We can form it into a convex quadratic optimisation problem with no constraints given by,

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2\sigma^2} \mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - \frac{1}{\sigma^2} (\mathbf{X}^T \mathbf{y})^T \mathbf{w} + \\ & \quad \left(\frac{1}{2\sigma^2} \mathbf{y}^T \mathbf{y} + \frac{N}{2} \log 2\pi\sigma^2 \right) \end{aligned} \tag{2.13}$$

What about minimising the Lasso Object,

$$\mathcal{L}_{Lasso} = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} + \lambda \sum_{i=1}^D |w_i| \tag{2.14}$$

In this case, we do not have a closed-form solution because of the absolute value. We cannot frame it into a linear programming problem either because we have the quadratic loss.

However, what we can do is to rephrase it into a quadratic optimisation problem with linear constraint(as what we did in transforming the absolute loss into linear programming previously).

Semi-definite Programming

$$\begin{aligned} & \text{minimize } \operatorname{tr}(\mathbf{c}^T \mathbf{x}) \\ & \text{s.t. } \operatorname{tr}(\mathbf{A}_i \mathbf{x}) = \mathbf{B}_i \end{aligned} \tag{2.15}$$

Where \mathbf{x} is positive semi-definite and $\operatorname{tr}(\mathbf{A})$ is the trace of the matrix \mathbf{A} .

2.2 Non-convex Optimization

Let's first recap some calculus background.

Theorem 2.12 *If a function f is differentiable, the gradient of f at a point is either 0 or perpendicular to the contour of f at that point.*

For $f : \mathbb{R}^n \rightarrow \mathbb{R}$, its gradient $\nabla f(\mathbf{a}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined as the vector,

$$\nabla f(\mathbf{a}) = [\frac{\partial f}{\partial x_1}(\mathbf{a}) \ \dots \ \frac{\partial f}{\partial x_n}(\mathbf{a})]^T \quad (2.16)$$

Each component of the gradient says how fast the function changes with respect to the standard basis,

$$\frac{\partial f}{\partial x_i}(\mathbf{a}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{a} + h\hat{\mathbf{i}}) - f(\mathbf{a})}{h} \quad (2.17)$$

Where $\hat{\mathbf{i}}$ is the unit vector in the direction of x , which captures information in which direction we move.

For changing with respect to the direction of a certain arbitrary vector \mathbf{v} , the directional derivative in the direction of \mathbf{v} is give as,

$$\nabla_{\mathbf{v}} f(\mathbf{a}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{a} + h\hat{\mathbf{v}}) - f(\mathbf{a})}{h} = \nabla f(\mathbf{a}) \cdot \mathbf{v} \quad (2.18)$$

Geometrically speaking, $\nabla_{\mathbf{v}} f(\mathbf{a})$ simply projects $\nabla f(\mathbf{a})$ to \mathbf{v} ,

$$\nabla f(\mathbf{a}) \cdot \mathbf{v} = \|\nabla f(\mathbf{a})\| \|\mathbf{v}\| \cos\theta \quad (2.19)$$

Where $\cos\theta$ is the angle between $\nabla f(\mathbf{a})$ and \mathbf{v} . The maximum value is obtained when $\cos\theta = 1$ or $\theta = 0$, where $\nabla f(\mathbf{a})$ and \mathbf{v} have the same direction.

If we want to capture the curvature of the surface, we will have to consider the matrix of all second-order partial derivatives of f , namely the Hessian matrix \mathbf{H} .

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function taking as input a vector $\mathbf{x} \in \mathbb{R}^n$ and outputting a scalar $f(\mathbf{x}) \in \mathbb{R}$. If all second partial derivatives of f exist and are continuous over the domain of the function, then the Hessian matrix \mathbf{H} of f is a square $n \times n$ matrix defined as,

$$(\mathbf{H}f)_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j} \quad (2.20)$$

The Hessian matrix is a symmetric matrix.

If $(\mathbf{H}f)$ has positive(or negative) eigenvalues, what we find at $\nabla f = 0$ is local minimum(or maximum).

If $(\mathbf{H}f)$ has mixed eigenvalues, what we find at $\nabla f = 0$ is saddle point.

If $(\mathbf{H}f)$ has eigenvalues of 0, it means that there is no inverse and thus the gradient is locally unchanging.

2.2.1 Gradient Descent

Gradient descent is one of the simplest yet very general optimization algorithms. It is an iterative algorithm, producing a new vector \mathbf{w}_{t+1} at each iteration as,

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t \quad (2.21)$$

Where \mathbf{g}_t denotes the gradient $\nabla f(\mathbf{w}_t)$ and $\eta_t > 0$ is the learning rate or step size.

We can easily see that at each iteration, it moves in the direction in the steepest descent. For a convex function, no matter where it starts, it will reach the same global minimum.

For deciding the number of iterations, we usually have three tests for convergence:

- Fixed number of iterations: Terminate if $t > T$
- Small increase: Terminate if $f(\mathbf{w}_t - \eta_t \mathbf{g}_t) - f(\mathbf{w}) < \epsilon$
- Small change: Terminate if $\|\mathbf{w}_{t+1} - \mathbf{w}_t\| \leq \epsilon$

For very large and wide datasets(i.e. very large N and D), the computational complexity of each gradient calculation is linear in N and D while that of inverting matrices in a closed-form solution is relatively much higher.

The choice of a good step-size is also crucial . If the step-size is too large, the algorithm may never converge; if the step-size is too small, the convergence may be too slow. Sometimes, we may even want a time-varying step-size.

We can choose a constant step size,

$$\eta_t = c \quad (2.22)$$

or a decaying step-size,

$$\eta_t = c/t \quad (2.23)$$

Other different rates of decay such as $\frac{1}{\sqrt{t}}$ is also common in practice.

We can also perform the backtracking line search:

- Start with $\eta_t = C/t$ (usually a large value)
- Check for a decrease: Is $f(\mathbf{w}_t - \eta_t \mathbf{g}_t) < f(\mathbf{w})$?
- If decrease condition is not met, multiply η_t by a decaying factor, e.g. 0.5
- Repeat until the decrease condition is met

2.2.2 Stochastic Gradient Descent

Suppose we minimize the objective function over data points $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$,

$$\mathcal{L}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{w}; \mathbf{x}_i, y_i) + \lambda \mathcal{R}(\mathbf{w}) \quad (2.24)$$

The gradient of the objective function is,

$$\nabla_{\mathbf{w}} \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} \ell(\mathbf{w}; \mathbf{x}_i, y_i) + \lambda \nabla_{\mathbf{w}} \mathcal{R}(\mathbf{w}) \quad (2.25)$$

For a random point (\mathbf{x}_i, y_i) , the expectation of the gradient at that point $\mathbf{g}_i = \nabla_{\mathbf{w}} \ell(\mathbf{w}; \mathbf{x}_i, y_i)$ is given by,

$$\mathbb{E}[\mathbf{g}_i] = \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} \ell(\mathbf{w}; \mathbf{x}_i, y_i) \quad (2.26)$$

The expectation of \mathbf{g}_i is exactly the same as the entire gradient except for the regularization term. In another word, by just taking the local gradient and doing our gradient descent approach will eventually leads to the very same outcome. This justifies using \mathbf{g}_i instead of $\nabla_{\mathbf{w}} \mathcal{L}$ in our update program for gradient descent.

The updating rule particular for stochastic gradient descent is the same as for gradient descent. But instead of computing the average gradient over the entire data set, in stochastic gradient descent, we just take a single random data point and compute the gradient of the loss function for that point.

If the objective also has a regularization term, the gradient of the regularization has to be added to the overall gradient. We have to consider the scale of the data to the regularization term. If the objective is the average loss of the data set, then we can just add the regularization at that point. If however, our objective is the sum of the loss, the gradient of the regularization term has to be scaled appropriately when using stochastic gradient descent.

The point of the stochastic gradient descent is that we compute the gradient at one data point instead of all data point. This enables so-called online-learning. In this particular case, we do not see the whole data set. Also, it is computationally cheaper. This is why stochastic gradient descent is rather appealing these days.

However, this is not necessarily always the case. There are recent papers showing doing batch gradient descent for a specific type of data that is very sparse is preferred over stochastic gradient descent. The following experiments compares the stochastic gradient descent with batch gradient descent.

Suppose we have 1000 data points for training and 1000 data points for test. There are two features drawn from the normal distribution $x_1 \in \mathcal{N}(0, 5)$ and $x_2 \in \mathcal{N}(0, 8)$. We take the labels to be centered so that there is no intercept. Thus the least-squares linear regression model is given as $f_w(\mathbf{x}) = x_1 w_1 + x_2 w_2$. The results of the experiment with parameters (w_1, w_2) initialized as $(-2, -3)$ and finally converging to $(1, 1)$ is shown in Figure 2.1.

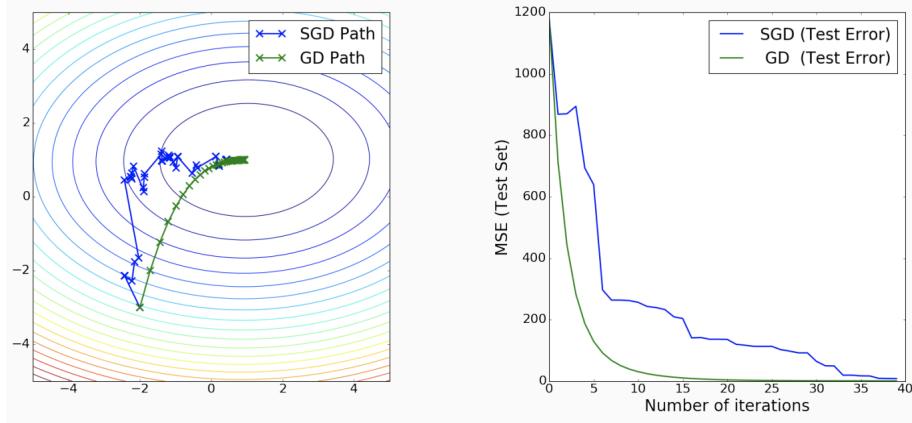


Figure 2.1: Experiment Results: SGD v.s. GD

As is shown in the Figure 2.1, it only takes gradient descent approach 15 to 20 iterations to converge while it takes stochastic gradient descent 40 iterations to reach the same result. However, there is another aspect that is equally relevant we did not take into account. For gradient descent, we need to look at the whole data set. Getting 15 or 20 of such iterations is much more computationally expensive than see 40 data points in 40 iterations of stochastic gradient descent.

In practice, what happens is that we do not necessarily use pure stochastic gradient descent. We can use what is called mini-batch gradient descent which uses batches that are fairly small. It reduces the variance in the gradients and hence it is more stable than stochastic gradient descent.

2.2.3 Sub-Gradient Descent

Recall the linear model trained with least-squares model and ℓ_1 regularization,

$$\mathcal{L}_{lasso}(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \sum_{i=1}^D |w_i| \quad (2.27)$$

When applying gradient descent on the lasso objective, we will encounter the problem that the objective function is not differentiable.

In mathematics, the sub-gradient generalizes the derivative to convex functions which are not necessarily differentiable. Sub-gradient arise in convex analysis, the study of convex functions, often in connection to convex optimization.

For a convex function f , a \mathbf{g} satisfying,

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \mathbf{g}^T(\mathbf{x} - \mathbf{x}_0) \quad (2.28)$$

is called a sub-gradient at \mathbf{x}_0 .

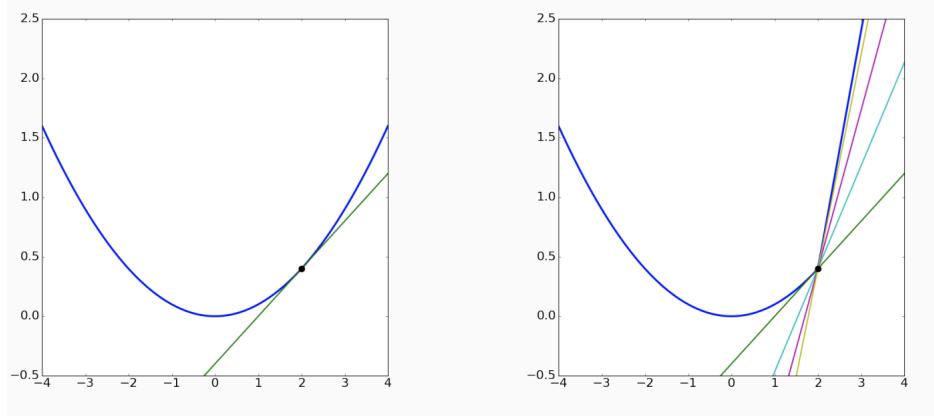


Figure 2.2: Sub-gradient

2.2.4 Constrained Convex Optimization

Recall what we did with gradient descent. We want to minimize $f(\mathbf{x})$ by moving in the negative gradient direction at each step,

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t) \quad (2.29)$$

Here we have no constraint on the parameters. We may choose them as we like so that we can reach the minimizer.

But sometimes in practice, it is important to minimize a function subject to additional constraints $\mathbf{w}_C \in C$. Here C is a set of possible values for our parameters \mathbf{w} .

So how can we solve such an optimization problem?

We first use the gradient step as we did before,

$$\mathbf{z}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t) \quad (2.30)$$

This \mathbf{z}_{t+1} we obtained in the equation above may happen to be outside our valid set C , i.e. it does not satisfy our constraints.

Then the way we have to proceed is that we should try to find a valid \mathbf{w}_C in our allowed set of possible value which is closest to \mathbf{z}_{t+1} . We can express this property as,

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}_C \in C} \|\mathbf{z}_{t+1} - \mathbf{w}_C\| \quad (2.31)$$

This is called the projection step.

Suppose we want to minimize $(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$ subject to the ridge constraints $\mathbf{w}^T\mathbf{w} < R$.

Imagine after t^{th} iteration \mathbf{w}_t is inside the allowed set defined by the ball of radius R . At $(t + 1)^{th}$ iteration with the gradient step \mathbf{w}_{t+1} goes up to outside the ball towards the minimizer. Then we have to project it back into the ball, which is shown by magenta line in Figure 2.3.

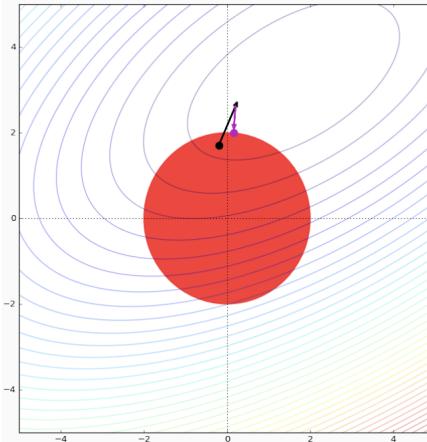


Figure 2.3: Ridge Constrained Convex Optimization

2.2.5 Second Order Methods

Gradient descent is a very general method for optimization. We can actually see gradient descent as an approach that approximates the function locally by a linear function and minimizes it instead. However, unless this linear function is constant(i.e. the gradient is 0) in which case the trajectory is the stationary point anyway, the minimum of this linear approximation will be negative infinity, which is why we have to choose carefully the step size.

In the case of Newton's method, we use the first and second-order partial derivatives. At any point, the trajectory will approximate the function by a quadratic function using the second-degree Taylor approximation. So a Newton step will take us directly to a unique stationary point of this quadratic approximation.

In optimization, the roots of $f'(i.e. \text{ solutions of } f'(\mathbf{x}) = 0)$ are stationary points(i.e. minimum/maximum/saddle points) of f , where f is a twice-differentiable function.

Newton's Method in One Dimension

- Construct a sequence of points x_1, \dots, x_n starting with an initial guess x_0
- Sequence converges towards a minimiser x^* of f using the sequence of second-order Taylor approximation of f around the iterates:

$$f(x) \approx f(x_k) + (x - x_k)f'(x_k) + \frac{1}{2}(x - x_k)^2 f''(x_k) \quad (2.32)$$

- $x_{k+1} = x^*$ is defined as the minimizer of this quadratic function.
- If f'' is positive, then the quadratic approximation is convex, and a minimizer is obtained by setting the derivative to 0:

$$0 = \frac{d(f(x_k) + (x - x_k)f'(x_k) + \frac{1}{2}(x - x_k)^2 f''(x_k))}{dx} = f'(x_k) + (x^* - x_k)f''(x_k) \quad (2.33)$$

Then,

$$x_{k+1} = x^* = x_k - f'(x_k)[f''(x_k)]^{-1} \quad (2.34)$$

Geometric Interpretation of Newton's Method

At iteration k, we fit a paraboloid to the surface of f at x_k with the same slopes and curvature as the surface at x_k and go for the extremum of that paraboloid.

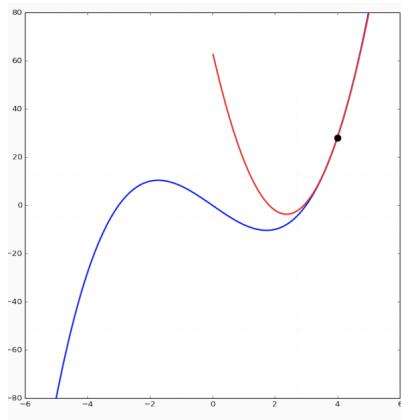


Figure 2.4: Geometric Interpretation of Newton's Method

Newton's Method in High Dimensions

- Approximate f around \mathbf{x}_k using second-order Taylor approximation:

$$f(\mathbf{x}) \approx f(\mathbf{x}_k) + \mathbf{g}_k^T(\mathbf{x} - \mathbf{x}_k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_k)^T \mathbf{H}_k (\mathbf{x} - \mathbf{x}_k) \quad (2.35)$$

Where \mathbf{H}_k is the Hessian of f .

- The gradient of f is given by,

$$\nabla_{\mathbf{x}} = \mathbf{g}_k + \mathbf{H}_k(\mathbf{x} - \mathbf{x}_k) \quad (2.36)$$

- Setting $\nabla_{\mathbf{x}} = 0$, we obtain $\mathbf{x}^* = \mathbf{x}_k - \mathbf{H}_k^{-1}\mathbf{g}_k$
- We move directly to the unique stationary point \mathbf{x}^* of f
- We repeat the above iteration with $\mathbf{x}_{k+1} = \mathbf{x}^*$

Computation and Convergence

- Computational requirements at each Newton step
 - $D + (\frac{D}{2})$ partial derivatives and inverse of Hessian
 - Instead: Compute \mathbf{x} as the solution of the system of linear equations
- $\mathbf{H}_k(\mathbf{x} - \mathbf{x}_k) = -\mathbf{g}_k \quad (2.37)$
- using factorization (e.g. Cholesky) of \mathbf{H}_k
- For convex function f
 - It converges to stationary points of the quadratic approximation
 - All stationary points are global minimum
 - Converges to unique minimizer quadratically fast if f is strongly convex
- For non-convex f
 - Stationary points may not be minimum nor in the decreasing direction of f
 - Not successful for training deep neural networks because of the abundance of saddle points for their non-convex objective functions

[..]/main.tex]subfiles

3

Classification

In machine learning, classification refers to a predictive modeling problem where a class label is predicted for a given example of input data. The input vector $\mathbf{x} = [x_1, \dots, x_D]^T$ consists of categorical and continuous features. The output y is a category,

$$y \in \{1, \dots, C\} \quad (3.1)$$

In classification problems, two main approaches are called the generative approach and the discriminative approach.

A discriminative setting is to model the conditional distribution of the output y given the input \mathbf{x} and the model parameters $\boldsymbol{\theta}$, symbolically, $p(y | \mathbf{x}, \boldsymbol{\theta})$.

A generative setting is to, however, model the full joint distribution of the input \mathbf{x} and output y given the model parameters $\boldsymbol{\theta}$, symbolically, $p(\mathbf{x}, y | \boldsymbol{\theta})$.

3.1 Generative Models for Classification

Given generative model $p(\mathbf{x}, y | \boldsymbol{\theta})$ representing the joint distribution of \mathbf{x} and y , for a new input \mathbf{x}_{new} , the conditional distribution for y , according to the Bayes rule, is given as,

$$p(y = c | \mathbf{x}_{new}, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_{new} | y = c, \boldsymbol{\theta})p(y = c | \boldsymbol{\theta})}{p(\mathbf{x}_{new} | \boldsymbol{\theta})} \quad (3.2)$$

Where $c \in \{1, \dots, C\}$ are the possible categories(i.e. labels) for y .

The numerator is the joint probability distribution $p(\mathbf{x}_{new}, y = c | \boldsymbol{\theta})$ and the denominator is the marginal distribution $p(\mathbf{x}_{new} | \boldsymbol{\theta})$.

According to the law of total probability,

$$p(\mathbf{x}_{new} | \boldsymbol{\theta}) = \sum_{c'=1}^C p(\mathbf{x}_{new} | y = c', \boldsymbol{\theta})p(y = c' | \boldsymbol{\theta}) \quad (3.3)$$

Thus, the conditional distribution for y according to the Bayes rule can be rephrased as,

$$p(y = c \mid \mathbf{x}_{new}, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_{new} \mid y = c, \boldsymbol{\theta})p(y = c \mid \boldsymbol{\theta})}{\sum_{c'=1}^C p(\mathbf{x}_{new} \mid y = c', \boldsymbol{\theta})p(y = c' \mid \boldsymbol{\theta})} \quad (3.4)$$

To predict the most-likely class:

$$\hat{y} = \arg \max_c p(y = c \mid \mathbf{x}_{new}, \boldsymbol{\theta}) \quad (3.5)$$

Suppose we have the output y only parameterized by $\boldsymbol{\pi}$. The joint distribution in our generative model can be expressed as,

$$p(\mathbf{x}, y \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = p(y \mid \boldsymbol{\theta}, \boldsymbol{\pi})p(\mathbf{x} \mid y, \boldsymbol{\theta}, \boldsymbol{\pi}) = p(y \mid \boldsymbol{\pi})p(\mathbf{x} \mid y, \boldsymbol{\theta}) \quad (3.6)$$

We use parameters π_c such that $\sum_c \pi_c = 1$ and model the marginal distribution of output y as,

$$p(y = c \mid \boldsymbol{\pi}) = \pi_c \quad (3.7)$$

For each class $c = 1, \dots, C$, the class-conditional distributions of input \mathbf{x} given the class label y parameterized by $\boldsymbol{\theta}$ is modeled as $p(\mathbf{x} \mid y = c, \boldsymbol{\theta}_c)$.

To define the model, we need to estimate the parameters $\pi_c, \boldsymbol{\theta}_c$ for each $c \in \{1, \dots, C\}$.

Assume we have the dataset $\mathcal{D} = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N))$, where each vector $\mathbf{x}_i \in \mathbb{R}^D$. Each data point is drawn independently from the same distribution $p(\mathbf{x}, y)$. Let N_c be the number of data points with $y_i = c$ so that $\sum_{c=1}^C N_c = N$.

The probability for a single data point (\mathbf{x}_i, y_i) is given as,

$$p(\mathbf{x}_i, y_i \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = p(y_i \mid \boldsymbol{\pi})p(\mathbf{x}_i \mid y_i, \boldsymbol{\theta}) = \prod_{c=1}^C \pi_c^{\mathbb{1}(y_i=c)} p(\mathbf{x}_i \mid y_i, \boldsymbol{\theta}) \quad (3.8)$$

Where $\mathbb{1}$ is the indicator function, i.e. $\mathbb{1} = 1$ for $y_i = c$ and $\mathbb{1} = 0$ otherwise.

The likelihood of the data is given by,

$$p(\mathcal{D} \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = \prod_{i=1}^N p(\mathbf{x}_i, y_i \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = \prod_{i=1}^N \left(\prod_{c=1}^C \pi_c^{\mathbb{1}(y_i=c)} p(\mathbf{x}_i \mid y_i, \boldsymbol{\theta}) \right) \quad (3.9)$$

The log-likelihood is given by,

$$\log p(\mathcal{D} \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = \sum_{c=1}^C N_c \log \pi_c + \sum_{i=1}^N \log p(\mathbf{x}_i \mid y_i, \boldsymbol{\theta}) \quad (3.10)$$

Where N_c is the number of data points falls in category $y_i = c$.

The first term $\sum_{c=1}^C N_c \log \pi_c$ is only dependent on parameter $\boldsymbol{\pi}$ while the second term $\sum_{i=1}^N \log p(\mathbf{x}_i | y_i, \boldsymbol{\theta})$ is independent of $\boldsymbol{\pi}$.

Let us first estimate for $\boldsymbol{\pi}$. We get the constrained optimization problem,

$$\begin{aligned} & \text{maximize} \quad \sum_{c=1}^C N_c \log \pi_c \\ & \text{s.t.} \quad \sum_{c=1}^C \pi_c - 1 = 0 \end{aligned} \tag{3.11}$$

This problem can be solved by the method of Lagrange multipliers given by,

$$\Lambda(\mathbf{p}\mathbf{i}, \lambda) = \sum_{c=1}^C N_c \log \pi_c - \lambda \left(\sum_{c=1}^C \pi_c - 1 \right) \tag{3.12}$$

Setting the partial derivatives with respect to $\mathbf{p}\mathbf{i}$ and λ to zero,

$$\begin{aligned} \frac{\partial \Lambda(\boldsymbol{\pi}, \lambda)}{\partial \pi_c} &= \frac{N_c}{\pi_c} - \lambda = 0 \\ \frac{\partial \Lambda(\boldsymbol{\pi}, \lambda)}{\partial \lambda} &= \sum_{c=1}^C \pi_c - 1 = 0 \end{aligned} \tag{3.13}$$

we obtain,

$$\begin{aligned} \pi_c &= \frac{N_c}{\lambda} \\ \lambda &= \sum_{c=1}^C N_c = N \end{aligned} \tag{3.14}$$

Thus, we get the estimation,

$$\pi_c = \frac{N_c}{N} \tag{3.15}$$

That is to say, for generative classification models, the marginal distribution $p(y | \boldsymbol{\pi})$ over the classes is the empirical distribution over the classes.

We next should estimate the parameters $\boldsymbol{\theta}$ of the model to choose the suitable model for the class-conditional distribution $p(\mathbf{x} | y, \boldsymbol{\theta})$.

3.1.1 Naïve Bayes Classifier

The Naïve Bayes classifier is based on Bayes rule and a set of conditional independence assumptions.

Definition 3.1 Given three sets of random variables X, Y and Z , we say X is conditionally independent of Y given Z if and only,

$$\forall (i, j, k) p(X = x_i | Y = y_j, Z = z_k) = p(X = x_i | Z = z_k) \tag{3.16}$$

Assume that the features are conditionally independent given the class label c . The class-conditional distribution for data point (\mathbf{x}_i, y_i) becomes,

$$p(\mathbf{x}_i \mid y_i = c, \boldsymbol{\theta}) = p(\mathbf{x}_i \mid y_i = c, \boldsymbol{\theta}_c) = \prod_{j=1}^D p(x_{ij} \mid \boldsymbol{\theta}_{jc}) \quad (3.17)$$

Consider, for example, the case where the input \mathbf{x}_i has two features x_{i1} and x_{i2} . In this case,

$$\begin{aligned} p(\mathbf{x}_i \mid y_i = c, \boldsymbol{\theta}_c) &= p(x_{i1}, x_{i2} \mid y = c, \boldsymbol{\theta}_c) \\ &= p(x_{i1} \mid x_{i2}, y = c, \boldsymbol{\theta}_c)p(x_{i2} \mid y = c, \boldsymbol{\theta}_c) \\ &= p(x_{i1} \mid y = c, \boldsymbol{\theta}_c)p(x_{i2} \mid y = c, \boldsymbol{\theta}_c) \end{aligned} \quad (3.18)$$

Now the joint distribution given the parameters $\boldsymbol{\theta}$ and $\boldsymbol{\pi}$ becomes,

$$\begin{aligned} p(\mathbf{x}_i, y_i \mid \boldsymbol{\theta}, \boldsymbol{\pi}) &= p(y_i \mid \boldsymbol{\pi})p(\mathbf{x}_i \mid y_i, \boldsymbol{\theta}) \\ &= \prod_{c=1}^C \pi_c^{\mathbb{1}(y_i=c)} p(\mathbf{x}_i \mid y_i, \boldsymbol{\theta}) \\ &= \prod_{c=1}^C \pi_c^{\mathbb{1}(y_i=c)} \prod_{c=1}^C \prod_{j=1}^D p(x_{ij} \mid \boldsymbol{\theta}_{jc})^{\mathbb{1}(y_i=c)} \end{aligned} \quad (3.19)$$

The likelihood of the data is given by,

$$p(\mathcal{D} \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = \prod_{i=1}^N p(\mathbf{x}_i, y_i \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = \prod_{i=1}^N \left(\prod_{c=1}^C \pi_c^{\mathbb{1}(y_i=c)} \prod_{c=1}^C \prod_{j=1}^D p(x_{ij} \mid \boldsymbol{\theta}_{jc})^{\mathbb{1}(y_i=c)} \right) \quad (3.20)$$

The log-likelihood is given by,

$$\log p(\mathcal{D} \mid \boldsymbol{\theta}, \boldsymbol{\pi}) = \sum_{c=1}^C N_c \log \pi_c + \sum_{c=1}^C \sum_{j=1}^D \sum_{i:y_i=c} \log p(x_{ij} \mid \boldsymbol{\theta}_{jc}) \quad (3.21)$$

Since we model for individual features with naïve Bayes classifier, it is easy to mix and match different models for different features.

For example, if feature x_j is real-valued(i.e. $x_j \in \mathbb{R}$), we can use Gaussian model $\theta_{jc} = (\mu_{jc}, \sigma_{jc}^2)$ where maximum likelihood estimation for θ_{jc} is given by the empirical mean and variance. If feature x_j is categorical with categories $1, \dots, k$. We can use multinoulli model: $x_j = \ell$ with probability $\theta_{jc,\ell}$ such that $\sum_{\ell=1}^k \theta_{jc,\ell} = 1$, where maximum likelihood estimation for θ_{jc} is obtained similarly to that for π_c . Also, there is no need to convert categorical variables to real-valued vectors this way.

With the naïve Bayes conditional independence assumption, naïve Bayes has at most $C \cdot D$ parameters $(\theta_{jc})_{j \in [D], c \in [C]}$. However, without the conditional independence assumption, we have to assign a probability for each of

the 2^D possible feature vectors. Thus, we have $C \cdot 2^D$ parameters which leads to either overfitting if $N \cdot 2^D$ or computationally prohibitive if $N \gg 2^D$. The 'naïve' assumption breaks the curse of dimensionality and avoids overfitting.

With the naïve Bayes classifier, we can also deal with missing data for prediction. Recall the prediction rule in a generative model,

$$p(y = c | \mathbf{x}_{new}, \boldsymbol{\theta}) = \frac{\pi_c \prod_{j=1}^D p(x_j | y = c, \boldsymbol{\theta}_{jc})}{\sum_{c'=1}^C \prod_{j=1}^D p(x_j | y = c', \boldsymbol{\theta}_{jc}) p(y = c' | \boldsymbol{\theta})} \quad (3.22)$$

Supposing our data point \mathbf{x}_{new} is missing its first element x_1 , we can simply skip it,

$$p(y = c | \mathbf{x}_{new}, \boldsymbol{\theta}) = \frac{\pi_c \prod_{j=2}^D p(x_j | y = c, \boldsymbol{\theta}_{jc})}{\sum_{c'=1}^C \prod_{j=2}^D p(x_j | y = c', \boldsymbol{\theta}_{jc}) p(y = c' | \boldsymbol{\theta})} \quad (3.23)$$

This can be done for other generative models, but marginalization requires summation/integration.

Gaussian Discriminant Analysis

Recall the form of the joint distribution in a generative model,

$$p(\mathbf{x}, y | \boldsymbol{\theta}, \boldsymbol{\pi}) = p(y | \boldsymbol{\pi}) p(\mathbf{x} | y, \boldsymbol{\theta}) \quad (3.24)$$

For classes, we use parameters π_c such that $\sum_{c=1}^C \pi_c = 1$,

$$p(y = c | \boldsymbol{\pi}) = \pi_c \quad (3.25)$$

If we model the class-conditional density for class $c \in \{1, \dots, C\}$ as a multivariate normal distribution with mean $\boldsymbol{\mu}_c$ and covariance matrix $\boldsymbol{\Sigma}_c$,

$$p(\mathbf{x} | y = c, \boldsymbol{\theta}_c) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad (3.26)$$

the model is given as,

$$p(\mathbf{x}, y | \boldsymbol{\theta}, \boldsymbol{\pi}) = \prod_{c=1}^C \pi_c^{\mathbb{1}(y=c)} \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad (3.27)$$

Note that the term 'discriminant' here refers to the shape of the boundaries that discriminate between the classes and it is still a generative model.

The model here is more suitable than naïve Bayes classifier in some cases such as predicting zebra or giraffe based on animal height and weight.

Given data $\mathcal{D} = ((\mathbf{x}_1), y_1), \dots, ((\mathbf{x}_N), y_N)$, the log-likelihood is given by,

$$\log p(\mathcal{D} | \boldsymbol{\theta}) = \sum_{c=1}^C N_c \log \pi_c + \sum_{c=1}^C \left(\sum_{i:y_i=c} \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \right) \quad (3.28)$$

As is in the case of naïve Bayes classifier, we have

$$\pi_c = \frac{N_c}{N} \quad (3.29)$$

For parameters $\boldsymbol{\mu}_c$ and $\boldsymbol{\Sigma}_c$, we have:

$$\begin{aligned} \hat{\boldsymbol{\mu}}_c &= \frac{1}{N_c} \sum_{i:y_i=c} \mathbf{x}_i \\ \hat{\boldsymbol{\Sigma}}_c &= \frac{1}{N_c} \sum_{i:y_i=c} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_c)^T \end{aligned} \quad (3.30)$$

We won't go into the details of the maximum likelihood estimation of these parameters. If interested, see section 4.1 from Murphy's book.

3.1.2 Quadratic Discriminant Analysis

The prediction rule for the generative model is,

$$p(y = c \mid \mathbf{x}_{new}, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_{new} \mid y = c, \boldsymbol{\theta})p(y = c \mid \boldsymbol{\theta})}{\sum_{c'=1}^C p(\mathbf{x}_{new} \mid y = c', \boldsymbol{\theta})p(y = c' \mid \boldsymbol{\theta})} \quad (3.31)$$

When the densities $p(\mathbf{x} \mid y = c, \boldsymbol{\theta}_c)$ are multivariate normal distributions,

$$p(y = c \mid \mathbf{x}_{new}, \boldsymbol{\theta}) = \frac{\pi_c |2\pi\boldsymbol{\Sigma}_c|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1}(\mathbf{x}-\boldsymbol{\mu}_c)}}{\sum_{c'=1}^C \pi_{c'} |2\pi\boldsymbol{\Sigma}_{c'}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{c'})^T \boldsymbol{\Sigma}_{c'}^{-1}(\mathbf{x}-\boldsymbol{\mu}_{c'})}} \quad (3.32)$$

The denominator is the same for all classes $c' \in \{1, \dots, C\}$.

The boundaries between two classes c_1 and c_2 is given by the data points \mathbf{x} with,

$$p(y = c_1 \mid \mathbf{x}, \boldsymbol{\theta}) = p(y = c_2 \mid \mathbf{x}, \boldsymbol{\theta}) \quad (3.33)$$

This means,

$$\frac{\pi_{c_1} |2\pi\boldsymbol{\Sigma}_{c_1}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{c_1})^T \boldsymbol{\Sigma}_{c_1}^{-1}(\mathbf{x}-\boldsymbol{\mu}_{c_1})}}{\pi_{c_2} |2\pi\boldsymbol{\Sigma}_{c_2}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{c_2})^T \boldsymbol{\Sigma}_{c_2}^{-1}(\mathbf{x}-\boldsymbol{\mu}_{c_2})}} = 1 \quad (3.34)$$

By taking the log, we obtain,

$$\frac{1}{2}((\mathbf{x} - \boldsymbol{\mu}_{c_2})^T \boldsymbol{\Sigma}_{c_2}^{-1}(\mathbf{x} - \boldsymbol{\mu}_{c_2}) - (\mathbf{x} - \boldsymbol{\mu}_{c_1})^T \boldsymbol{\Sigma}_{c_1}^{-1}(\mathbf{x} - \boldsymbol{\mu}_{c_1})) = \log \frac{\pi_{c_2} |2\pi\boldsymbol{\Sigma}_{c_2}|^{-\frac{1}{2}}}{\pi_{c_1} |2\pi\boldsymbol{\Sigma}_{c_1}|^{-\frac{1}{2}}} \quad (3.35)$$

Boundaries are given by quadratic functions and thus piece-wise quadratic curves as is shown in Figure 3.1.

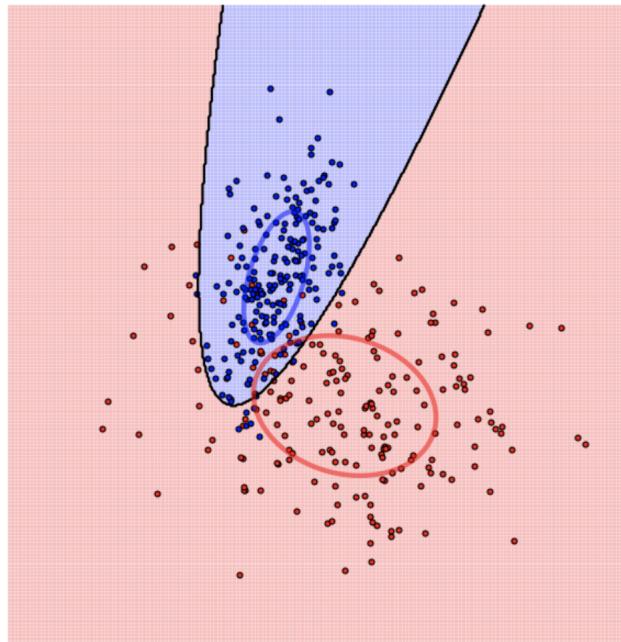


Figure 3.1: Quadratic Decision Boundaries

In quadratic discriminant analysis, the number of parameters is roughly $C \cdot D^2$. In high dimensions, this can lead to overfitting. How can we solve this problem?

We can use weight tying or parameter sharing to make the quadratic discriminant analysis problem degenerates to the linear discriminant analysis problem.

Moreover, we can use diagonal covariance matrices which corresponds to naïve Bayes. Since entries in diagonal matrices which are not along the diagonals is 0, the covariance between any two distinct features is 0, meaning that there is no correlation. This is exactly what happen with naïve Bayes - we make the conditional independence assumption that given the class labels, the features will be independent of each other.

In fact, the naïve Bayes is a special case of linear discriminant analysis, which is a special case of quadratic discriminant analysis.

Sometimes, we can use a discriminative classifier such as logistic regression and add regularization if needed.

3.1.3 Linear Discriminant Analysis

If the covariance matrix are shared or tied across different classes, denoted as Σ , our generative model is given by,

$$\begin{aligned}
 p(y = c \mid \mathbf{x}, \boldsymbol{\theta}) &= \frac{\pi_c |2\pi\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_c)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_c)}}{\sum_{c'=1}^C \pi_{c'} |2\pi\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{c'})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_{c'})}} \\
 &= \frac{\pi_c e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_c)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_c)}}{\sum_{c'=1}^C \pi_{c'} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_{c'})^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_{c'})}} \\
 &= \frac{e^{\boldsymbol{\mu}_c^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \Sigma^{-1} \boldsymbol{\mu}_c + \log \pi_c} \cdot e^{-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}}}{\sum_{c'=1}^C e^{\boldsymbol{\mu}_{c'}^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_{c'}^T \Sigma^{-1} \boldsymbol{\mu}_{c'} + \log \pi_{c'}} \cdot e^{-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}}} \\
 &= \frac{e^{\boldsymbol{\mu}_c^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^T \Sigma^{-1} \boldsymbol{\mu}_c + \log \pi_c}}{\sum_{c'=1}^C e^{\boldsymbol{\mu}_{c'}^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_{c'}^T \Sigma^{-1} \boldsymbol{\mu}_{c'} + \log \pi_{c'}}}
 \end{aligned} \tag{3.36}$$

If we let $\boldsymbol{\beta}_c = \boldsymbol{\mu}_c$ and $\gamma_c = -\frac{1}{2}\boldsymbol{\mu}_c^T \Sigma^{-1} \boldsymbol{\mu}_c + \log \pi_c$,

$$p(y = c \mid \mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\boldsymbol{\beta}_c^T \mathbf{x} + \gamma_c}}{\sum_{c'=1}^C e^{\boldsymbol{\beta}_{c'}^T \mathbf{x} + \gamma_{c'}}} = softmax(\boldsymbol{\eta})_c \tag{3.37}$$

Where $\boldsymbol{\eta} = [\boldsymbol{\beta}_1^T \mathbf{x} + \gamma_1, \dots, \boldsymbol{\beta}_C^T \mathbf{x} + \gamma_C]$.

The boundaries between two classes c_1 and c_2 is given by the data points \mathbf{x} with,

$$p(y = c_1 \mid \mathbf{x}, \boldsymbol{\theta}) = p(y = c_2 \mid \mathbf{x}, \boldsymbol{\theta}) \tag{3.38}$$

This means,

$$\frac{e^{\boldsymbol{\beta}_{c_1}^T \mathbf{x} + \gamma_{c_1}}}{e^{\boldsymbol{\beta}_{c_2}^T \mathbf{x} + \gamma_{c_2}}} = 1 \tag{3.39}$$

By taking the log, we obtain,

$$\boldsymbol{\beta}_{c_1}^T \mathbf{x} + \gamma_{c_1} = \boldsymbol{\beta}_{c_2}^T \mathbf{x} + \gamma_{c_2} \tag{3.40}$$

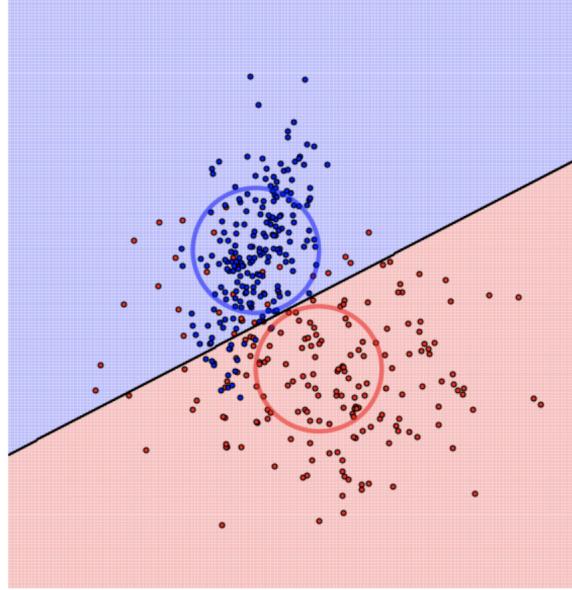


Figure 3.2: Linear Decision Boundaries

Boundaries are given by linear functions and thus piece-wise linear curves as is shown in Figure 3.2.

3.1.4 Softmax Function

Softmax function is given by,

$$\text{softmax}([z_1, \dots, z_n])_i = \frac{e^{z_i}}{\sum_{j \in [n]} e^{z_j}} \quad (3.41)$$

It normalizes the vector values into a probability distribution consisting of probabilities proportional to their exponentials.

Softmax can be seen as a smooth approximation of the argmax function. It is commonly used in statistical mechanics as the Boltzmann distribution.

Softmax has the property that large inputs generates large probabilities. It is translation invariant(i.e. multiplying all values by the same quantity does not change the ranking) but not scale invariant(i.e. de-emphasises the max value for values $\in [0, 1]$)

Softmax is often used as the last activation function of a neural network to normalise the network output to a probability distribution over predicted output classes.

3.1.5 Two-Class Linear Decision Boundaries and Sigmoid Function

When we have only 2 classes $y \in [0, 1]$,

$$\begin{aligned}
 p(y = 1 \mid \mathbf{x}, \boldsymbol{\theta}) &= \frac{e^{\boldsymbol{\beta}_1^T \mathbf{x} + \gamma_1}}{e^{\boldsymbol{\beta}_1^T \mathbf{x} + \gamma_1} + e^{\boldsymbol{\beta}_0^T \mathbf{x} + \gamma_0}} \\
 &= \frac{1}{1 + \frac{e^{\boldsymbol{\beta}_1^T \mathbf{x} + \gamma_1}}{e^{\boldsymbol{\beta}_0^T \mathbf{x} + \gamma_0}}} \\
 &= \frac{1}{1 + e^{-(\boldsymbol{\beta}_1^T - \boldsymbol{\beta}_0^T) \mathbf{x} + (\gamma_1 - \gamma_0)}} \\
 &= \text{sigmoid}((\boldsymbol{\beta}_1^T - \boldsymbol{\beta}_0^T) \mathbf{x} + (\gamma_1 - \gamma_0))
 \end{aligned} \tag{3.42}$$

The sigmoid function is defined as:

$$\text{sigmoid}(t) = \frac{1}{1 + e^{-t}} = \text{softmax}([0, t]) \tag{3.43}$$

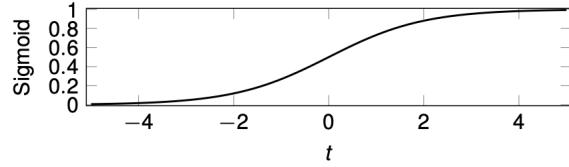


Figure 3.3: Sigmoid Function

In general, sigmoid functions refer to functions with S-shaped curves such as hyperbolic tangent and logistic(Equation 3.43) which is a special case of softmax for 1D axis in 2D space.

Sigmoid function is differentiable and has non-negative derivative at each point. It is monotonic. For $t < 0$ it is convex and for $t > 0$ it is concave.

3.2 Logistic Regression

Logistic regression is a discriminative classification method, i.e. it models the conditional distribution over the categorical output y given the input \mathbf{x} and parameters \mathbf{w} , symbolically, $p(y | \mathbf{w}, \mathbf{x})$.

3.2.1 Binary-Class Logistic Regression

In the case of binary classification, it is quite convenient to use the Bernoulli random variable $X \in \{0, 1\}$.

If X is a random variable with Bernoulli distribution, i.e. $X \sim \text{Bernoulli}(\theta)$, $\theta \in [0, 1]$ where,

$$X = \begin{cases} 1 & \text{with probability } \theta \\ 0 & \text{with probability } 1 - \theta \end{cases}$$

The probability mass function of this distribution is given as,

$$\begin{aligned} p(1 | \theta) &= \theta \\ p(0 | \theta) &= 1 - \theta \end{aligned} \tag{3.44}$$

More succinctly, it can be rephrased as,

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x} \tag{3.45}$$

for $x \in \{0, 1\}$.

Given input \mathbf{x} and models with parameters \mathbf{w} producing a value $f(\mathbf{x}, \mathbf{w}) \in [0, 1]$, we model the binary class label as,

$$y \sim \text{Bernoulli}(f(\mathbf{x}, \mathbf{w})) \tag{3.46}$$

For logistic regression, it builds up on a linear model composed with a sigmoid function $\sigma(t)$,

$$p(y | \mathbf{w}, \mathbf{x}) = \text{Bernoulli}(\sigma(\mathbf{w} \cdot \mathbf{x})) \tag{3.47}$$

Suppose we have estimated the model parameters $\mathbf{w} \in \mathbb{R}^D$, for a new data point \mathbf{x}_{new} , the model gives the probability,

$$p(y_{new} = 1 | \mathbf{x}_{new}, \mathbf{w}) = \sigma(\mathbf{w} \cdot \mathbf{x}_{new}) \tag{3.48}$$

Recall that the sigma function $\sigma : \mathbb{R} \rightarrow (0, 1)$ is defined by,

$$\sigma(t) = \frac{1}{1 + e^{-t}} \tag{3.49}$$

Thus,

$$p(y_{new} = 1 | \mathbf{x}_{new}, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}_{new}}} \tag{3.50}$$

Let's now take a look at how to compute the decision boundaries. For any given parameter θ of the Bernoulli distribution,

$$p(y = 1 \mid \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}} = \theta \quad (3.51)$$

Thus,

$$e^{-\mathbf{w} \cdot \mathbf{x}} = \frac{1}{\theta} - 1 \quad (3.52)$$

Subtracting 1 from the denominator,

$$\frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}} - 1} = \frac{1}{e^{-\mathbf{w} \cdot \mathbf{x}}} = \frac{1}{\frac{1}{\theta} - 1} = \frac{\theta}{1 - \theta} \quad (3.53)$$

Taking log on both sides,

$$\log 1 - \log e^{-\mathbf{w} \cdot \mathbf{x}} = \mathbf{w} \cdot \mathbf{x} = \log \frac{\theta}{1 - \theta} \quad (3.54)$$

And we obtain the hyper-plane,

$$\mathbf{w} \cdot \mathbf{x} = \log \frac{\theta}{1 - \theta} \quad (3.55)$$

Since $\sigma(t) \geq \frac{1}{2}$ for $\forall t \geq 0$, in order to make a prediction, we can simply use a threshold at $\frac{1}{2}$ given as,

$$\hat{y}_{new} = \mathbb{1}(\sigma(\mathbf{w} \cdot \mathbf{x}_{new}) \geq \frac{1}{2}) = \mathbb{1}(\mathbf{w} \cdot \mathbf{x}_{new} \geq 0) \quad (3.56)$$

The decision boundary is obtained when,

$$p(y = 1 \mid \mathbf{x}, \mathbf{w}) = p(y = 0 \mid \mathbf{x}, \mathbf{w}) = \frac{1}{2} \quad (3.57)$$

i.e. when,

$$\mathbf{w} \cdot \mathbf{x} = 0 \quad (3.58)$$

The decision boundary is shown in Figure 3.4.

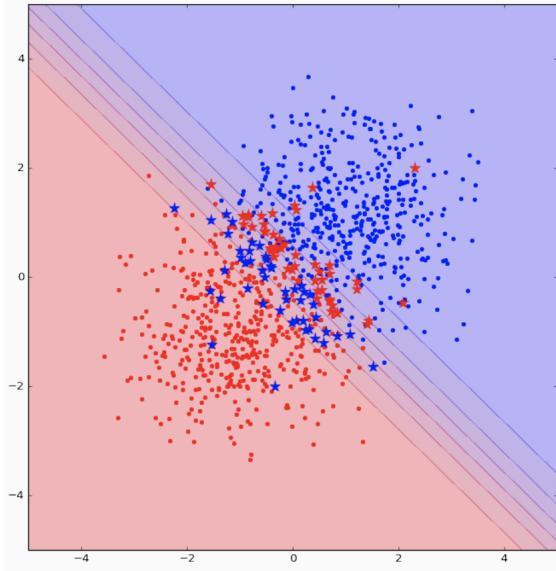


Figure 3.4: Logistic Regression Boundaries

As is shown in Figure 3.4, contour lines from bottom left to top right are for $\theta = 0.15, 0.3, 0.45, 0.6, 0.75, 0.9$ each respectively and the starred points represent mistakes made by the classifier.

Now, how can we estimate the parameters \mathbf{w} for our logistic regression classifier?

Assume we have the dataset $\mathcal{D} = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N))$, where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \{0, 1\}$.

The likelihood of the observing data, given the model parameters \mathbf{w} is given by,

$$\begin{aligned} p(\mathbf{y} | \mathbf{X}, \mathbf{w}) &= \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}_i)^{y_i} (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))^{1-y_i} \\ &= \prod_{i=1}^N \mu_i^{y_i} (1 - \mu_i)^{1-y_i} \end{aligned} \tag{3.59}$$

where $\mu_i = \sigma(\mathbf{w}^T \mathbf{x}_i)$.

The negative log likelihood is given by,

$$NLL(\mathbf{y} | \mathbf{X}, \mathbf{w}) = - \sum_{i=1}^N y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i) \tag{3.60}$$

which is the cross-entropy between y_i and u_i for $y_i \in \{0, 1\}$.

The gradient of the negative log likelihood with respect to \mathbf{w} is,

$$\nabla_{\mathbf{w}} NLL(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \sum_{i=1}^N \mathbf{x}_i (\mu_i - y_i) = \mathbf{X}^T (\boldsymbol{\mu} - \mathbf{y}) \quad (3.61)$$

The Hessian of the NLL is given as,

$$\mathbf{H} = \mathbf{X}^T \mathbf{S} \mathbf{X} \quad (3.62)$$

Where \mathbf{S} is the diagonal matrix with entries $s_{ii} = \mu_i(1 - \mu_i)$,

$$\mathbf{S} = \begin{bmatrix} s_{11} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & s_{nn} \end{bmatrix} = \begin{bmatrix} \mu_1(1 - \mu_1) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \mu_n(1 - \mu_n) \end{bmatrix} \quad (3.63)$$

For $\forall \mathbf{z} \in \mathbb{R}^D$,

$$\begin{aligned} \mathbf{z}^T \mathbf{H} \mathbf{z} &= \mathbf{z}^T \mathbf{X}^T \mathbf{S} \mathbf{X} \mathbf{z} \\ &= (\mathbf{X} \mathbf{z})^T \mathbf{S} \mathbf{X} \mathbf{z} \\ &= \sum_{i=1}^N s_{ii} [\mathbf{X} \mathbf{z}]_i^2 \end{aligned} \quad (3.64)$$

Since $s_{ii} = \mu_i(1 - \mu_i) > 0$ for $\mu_i \in (0, 1)$,

$$\mathbf{z}^T \mathbf{H} \mathbf{z} = \sum_{i=1}^N s_{ii} [\mathbf{X} \mathbf{z}]_i^2 > 0 \quad (3.65)$$

the Hessian of NLL is positive definite. Thus, NLL is convex and we can use convex optimization method to minimise NLL.

For small number D of dimensions, we can apply Newton's method to estimate \mathbf{w} .

Let \mathbf{w}_t be the parameters after t Newton steps. The gradient and the Hessian are given by,

$$\begin{aligned} \mathbf{g}_t &= \mathbf{X}^T (\boldsymbol{\mu}_t - \mathbf{y}) = -\mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}_t) \\ \mathbf{H}_t &= \mathbf{X}^T \mathbf{S}_t \mathbf{X} \end{aligned} \quad (3.66)$$

The Newton updating rule is given by,

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \mathbf{H}_t^{-1} \mathbf{g}_t \\ &= \mathbf{w}_t + (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}_t) \\ &= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{S}_t \mathbf{X}) \mathbf{w}_t + (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}_t) \\ &= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T \mathbf{S}_t^T [\mathbf{X} \mathbf{w}_t + \mathbf{S}_t^{-1} (\mathbf{y} - \boldsymbol{\mu}_t)] \\ &= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T \mathbf{S}_t^T \mathbf{z}_t \\ &= (\mathbf{X}^T \mathbf{S}_t^{\frac{1}{2}} \mathbf{S}_t^{\frac{1}{2}} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{S}_t^{\frac{1}{2}} \mathbf{S}_t^{\frac{1}{2}} \mathbf{z}_t \end{aligned} \quad (3.67)$$

Where $\mathbf{z}_t = \mathbf{X}\mathbf{w}_t + \mathbf{S}_t^{-1}(\mathbf{y} - \boldsymbol{\mu}_t)$.

Since the transpose of a diagonal matrix is itself,

$$\mathbf{w}_{t+1} = ((\mathbf{S}_t^{\frac{1}{2}}\mathbf{X})^T \mathbf{S}_t^{\frac{1}{2}}\mathbf{X})^{-1} (\mathbf{S}_t^{\frac{1}{2}}\mathbf{X})^T \mathbf{S}_t^{\frac{1}{2}}\mathbf{z}_t \quad (3.68)$$

Let $\tilde{\mathbf{X}} = \mathbf{S}_t^{\frac{1}{2}}\mathbf{X}$ and $\tilde{\mathbf{y}} = \mathbf{S}_t^{\frac{1}{2}}\mathbf{z}_t$, the equation above can be rephrased as,

$$\mathbf{w}_{t+1} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}} \tilde{\mathbf{y}} \quad (3.69)$$

The updating weight \mathbf{w}_{t+1} is in the same form as the weight minimising the least square objective $\mathcal{L}(\mathbf{w}) = \sum_i (\mathbf{x}_i^T \mathbf{w} - y_i)^2$.

Thus, we can reconstruct the Newton step as the problem finding \mathbf{w} that minimize the objective,

$$\begin{aligned} \mathcal{L}(\mathbf{w}) &= \sum_i (\tilde{\mathbf{x}}_i^T \mathbf{w} - \tilde{y}_i)^2 \\ &= \sum_i [\mathbf{S}_{t,ii}^{\frac{1}{2}}(\mathbf{x}_i - z_{t,i})]^2 \\ &= \sum_i \mathbf{S}_{t,ii}(\mathbf{x}_i - z_{t,i})^2 \\ &= \mathbf{S}_t(\mathbf{X}\mathbf{w}_t - \mathbf{z}_t)^T(\mathbf{X}\mathbf{w}_t - \mathbf{z}_t) \end{aligned} \quad (3.70)$$

The objective above is called the weighted least squares of which we can use to compute \mathbf{w}_{t+1} at each Newton step. The optimization method here is called the iteratively re-weighted least squares method.

- Each step requires re-weighting of the residual by a new diagonal matrix \mathbf{S}_t
- Each step uses a new vector \mathbf{z}_t which depends on \mathbf{w}_t
- We proceed iteratively, one Newton step after another

3.2.2 Multi-Class Logistic Regression

Consider now $C > 2$ classes: $y \in \{1, \dots, C\}$. There are parameters \mathbf{w}_c for every class c which forms a parameter matrix $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_C] \in \mathbb{R}^{D \times C}$.

The multi-class logistic regression model is given by,

$$\begin{aligned} p(y = c \mid \mathbf{x}, \mathbf{W}) &= \frac{e^{\mathbf{w}_c \mathbf{x}}}{\sum_{c'=1}^C e^{\mathbf{w}_{c'} \mathbf{x}}} \\ &= \text{softmax}([\mathbf{w}_1 \mathbf{x}, \dots, \mathbf{w}_C \mathbf{x}]) \end{aligned} \quad (3.71)$$

The parameter estimation for multi-class logistic regression model is the same as in the case of binary logistic regression model. The negative log-likelihood(NLL) is still convex.

Actually, we can see binary logistic regression as a special case of multi-class logistic regression where $\mathbf{W} = [\mathbf{w}_0, \mathbf{w}_1] \in \mathbb{R}^{D \times 2}$,

$$\text{softmax}([\mathbf{w}_0, \mathbf{w}_1]) = \frac{e^{\mathbf{w}_1 \mathbf{x}}}{e^{\mathbf{w}_1 \mathbf{x}} + e^{\mathbf{w}_0 \mathbf{x}}} = \sigma((\mathbf{w}_1 - \mathbf{w}_0) \mathbf{x}) \quad (3.72)$$

3.3 Multi-class Classification and Measuring Performance

3.3.1 Multi-Class Classification

Multi-class classification refers to the problem where the number of classes $C > 2$. In practice, the following approaches are common.

The first one is called the One-vs-One approach. It trains $\frac{C(C-1)}{2}$ classifiers and each training procedure only uses on average $\frac{2}{C}$ of the training data. For a new input, the output is the most common classification for the classifiers.

The second approach is called the One-vs-Rest classifiers. It trains C different classifiers, i.e. one class vs the rest $C - 1$ classes and each training procedure only uses the entire training data. It is typical for classifiers that yield class membership probability or scores. For a new input, it picks the class with the largest probability or score.

For example, if we have three classes c_1, c_2 and c_3 , the One-vs-One approach trains $\frac{3 \times (3-1)}{2} = 3$ classifiers which discriminate c_1 from c_2 , c_1 from c_3 , and c_2 from c_3 each respectively. The output of a new input is the majority vote of the prediction labels obtained on all classifiers. The One-vs-Rest approach, however, trains 3 classifiers: c_1 against c_2 and c_3 ; c_2 against c_1 and c_3 ; and c_3 against c_1 and c_2 . The output of a new input is the class with largest probability or score. If we have ties, we can break it by taking the value of $\mathbf{w}\mathbf{x}_{new} + w_0$.

A more efficient method is to reduce multi-class classification problem to binary classification problem. It divides class into pairs of disjoint subsets and trains a binary classifier to separate the subsets of each pair. Then it uses an error correcting approach to determine the class label.

The classification error is defined by the number of misclassified data points. However, in classification problems, not all mistakes are equally problematic, for example, failing to detect the medical risk v.s. inaccurately predicting the chance of the risk. How can we solve this problem?

In logistic regression, we use threshold 0.5 to label a data point positive. If we want very few false positives, we can raise the threshold to 0.9, i.e. predicting something as positive only if it were 90% sure.

In generative models, we have decision boundaries when $\frac{p(y=1|\mathbf{x}_{new}, \mathbf{w})}{p(y=0|\mathbf{x}_{new}, \mathbf{w})} = 1$, i.e. all errors are treated equally. However, we can change the ratio if one type of the error is more costly than the other.

3.3.2 Measuring Performance Binary Classification

When measuring performance for binary classification, we use confusion matrix defined as,

And we define the following,

Prediction	Actual Labels	
	1	0
1	TN(True Positive)	FP(False Positive)
0	FN(False Negative)	TN(True Negative)

Table 3.1: Confusion Matrix

- True Positive Rate, i.e. Sensitivity, which is the ratio of true positive to actual positive, give as,

$$TPR = \frac{TP}{TP + FN} \quad (3.73)$$

- False Positive Rate, i.e. Fall-out, which is the ratio of false positive to actual negative, give as,

$$FPR = \frac{FP}{FP + TN} \quad (3.74)$$

- True Negative Rate, i.e. Specificity, which is the ratio of true negative to actual negative, give as,

$$TNR = \frac{TN}{FP + TN} = 1 - FPR \quad (3.75)$$

- Precision which is the ratio of true positives to predicted positives, given as,

$$P = \frac{TP}{TP + FP} \quad (3.76)$$

- Accuracy, given as,

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN} \quad (3.77)$$

In practice, instead of analyzing many confusion matrices of different thresholds, we use receiver operating curves, or ROC curve, which gives an intuitive compact representation of all of them.

A receiver operating curve, or ROC curve, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied.

Usually, the x-axis and y-axis is FPR and TPR each respectively. It gives the relationship between the proportion of actual positive labels correctly classified to the proportion of actual negative labels incorrectly classified.

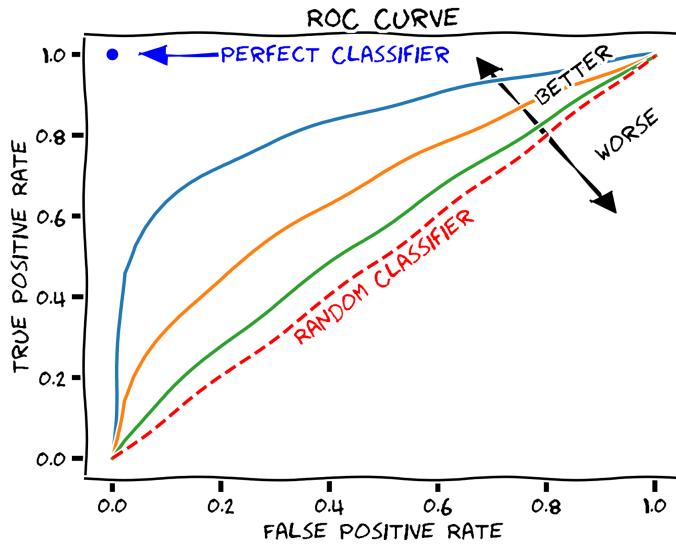


Figure 3.5: Receiver Operating Curve

The Area Under the Curve (AUC) is the measure of the ability of a classifier to distinguish between classes and is used as a summary of the ROC curve. The higher the AUC, the better the performance of the model at distinguishing between the positive and negative classes, i.e. the model is more preferred.

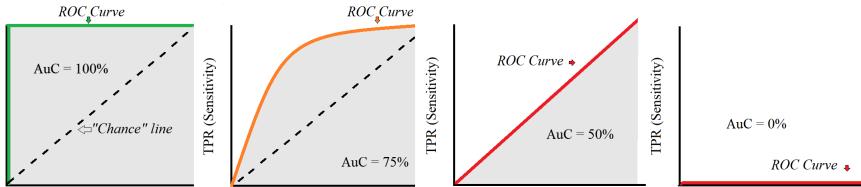


Figure 3.6: Area Under the Curve

Another metric beyond ROC curves is to replace FPR with Precision. Precision gives the proportion of positive results that were correctly classified. Precision is preferred over FPR in scenarios where there were lots of negative samples. Since precision does not include True Negatives, so it will not be affected by imbalance. For example, if we are studying a rare events, e.g., a rare disease. There are many more samples that do not observe the event than those that observe it.

3.3.3 Measuring Performance Multi-Class Classification

For multi-class classification, we generalise the confusion matrix defined as,

Prediction	Actual Labels			
	1	2	...	K
1	N ₁₁	N ₁₂	...	N _{1K}
2	N ₂₁	N ₂₂	...	N _{2K}
⋮	⋮	⋮	⋮	⋮
K	N _{K1}	N _{K2}	...	N _{KK}

Table 3.2: Confusion Matrix

Where $N_{i,j}$ is the items of class j in the dataset that are predicted to be of class i .

A good multi-class classifier should have large diagonal entries and small off-diagonal entries, i.e. such that,

$$\arg \max_{\theta} \text{tr}(\mathbf{C}) \quad (3.78)$$

Where $\text{tr}(\mathbf{C})$ is the trace of the confusion matrix.

3.4 Support Vector Machines

SVM is a popular discriminative model for classification. However, it is not a natural probabilistic interpretation.

Recall binary classification, our goal is to : find a linear separator. Suppose our data is linearly separable if there exists a linear separator that classifies all points correctly, as is shown in Figure 3.7. Which separator should be picked?

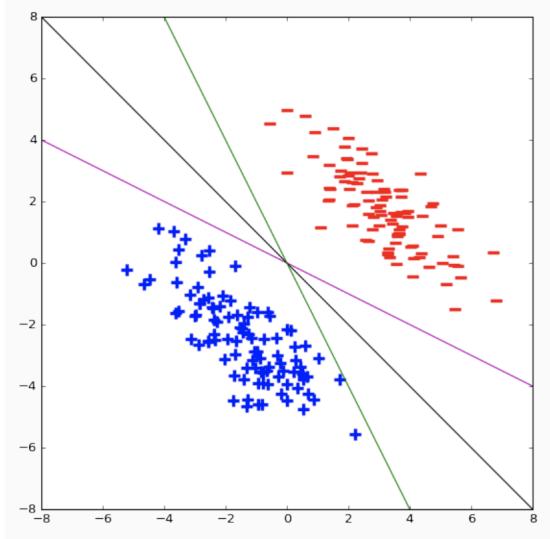


Figure 3.7: Linear Separators

3.4.1 Maximum Margin Principle

The margin for a data point refers to its distance to the separating hyperplane.

The maximum margin principle is to pick the separating boundary that maximises the smallest margin, i.e. the least distance between data and boundary. That is, to maximise the distance of the closest point from the decision boundary. The points that are closest to the decision boundary are called support vectors, as is shown in Figure 3.8.

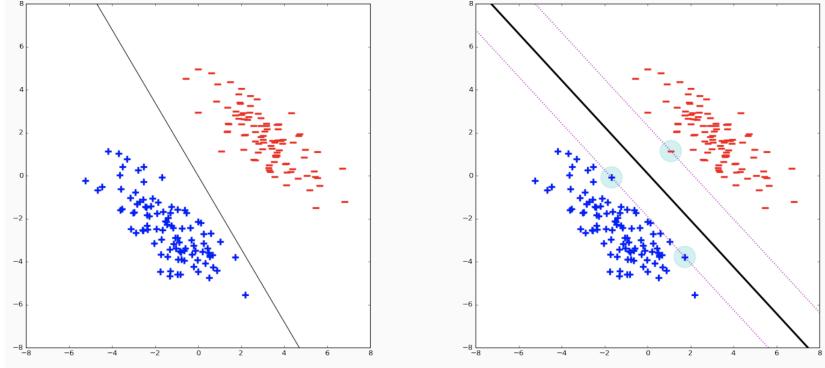


Figure 3.8: Support Vectors

Given a hyper-plane $H \equiv \mathbf{w}\mathbf{x} + w_0 = 0$ and a point $\mathbf{x} \in \mathbb{R}^D$, the distance of point \mathbf{x} from H is given by,

$$\frac{|\mathbf{w}\mathbf{x} + w_0|}{\|\mathbf{w}\|_2} \quad (3.79)$$

All points on the positive half-space created by H , labeled as $y = +1$, satisfy,

$$\mathbf{w}\mathbf{x} + w_0 > 0 \quad (3.80)$$

All points on the negative half-space created by H , labeled as $y = -1$, satisfy,

$$\mathbf{w}\mathbf{x} + w_0 < 0 \quad (3.81)$$

Our goal of maximising the distance of the closest point, i.e. support vectors, from the decision boundary can be rephrased as finding the distance $\|\mathbf{x} - \mathbf{x}^*\|$ between the point \mathbf{x}^* and the hyper-plane $\mathbf{w}\mathbf{x} + w_0 = 0$, where the point \mathbf{x} on the hyper-plane that is closest to \mathbf{x}^* gives the distance.

Equivalently, we seek for \mathbf{x} that optimises the following problem,

$$\begin{aligned} &\text{minimize } \|\mathbf{x} - \mathbf{x}^*\|_2^2 \\ &\text{s.t. } \mathbf{w}\mathbf{x} + w_0 = 0 \end{aligned} \quad (3.82)$$

We can solve this problem using the Lagrangian multiplier, given by,

$$\begin{aligned} \Lambda &= \|\mathbf{x} - \mathbf{x}^*\|_2^2 - 2\lambda(\mathbf{w}\mathbf{x} + w_0) \\ &= \|\mathbf{x}\|_2^2 - 2\mathbf{x}(\mathbf{x}^* + \lambda\mathbf{w}) + 2\lambda w_0 + \|\mathbf{x}^*\|_2^2 \end{aligned} \quad (3.83)$$

Setting its gradient,

$$\nabla_{\mathbf{x}} \Lambda(\mathbf{x}, \lambda) = 2\mathbf{x} - 2(\mathbf{x}^* + \lambda\mathbf{w}) = 0 \quad (3.84)$$

to 0, we obtained the critical point,

$$\mathbf{x} = \mathbf{x}^* + \lambda \mathbf{w} \quad (3.85)$$

By substituting \mathbf{x} into the hyper-plane equation, we obtain,

$$\lambda = -\frac{\mathbf{w}\mathbf{x}^* + w_0}{\|\mathbf{w}\|_2^2} \quad (3.86)$$

Then, the distance between \mathbf{x}^* and \mathbf{x} is given by,

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}^*\|_2 &= \|\mathbf{x}^* + \lambda \mathbf{w} - \mathbf{x}^*\|_2 \\ &= \|\lambda \mathbf{w}\|_2 \\ &= |\lambda| \|\mathbf{w}\|_2 \\ &= \frac{|\mathbf{w}\mathbf{x}^* + w_0|}{\|\mathbf{w}\|_2} \end{aligned} \quad (3.87)$$

Assume that the dataset $\mathcal{D} = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N))$ is linearly separable. For $y_i = +1$ we have $\mathbf{w}\mathbf{x}_i + w_0 > 0$ and for $y_i = -1$ we have $\mathbf{w}\mathbf{x}_i + w_0 < 0$. Thus, most compactly, for all points in the dataset \mathcal{D} we have,

$$y_i(\mathbf{w}\mathbf{x}_i + w_0) > 0 \quad (3.88)$$

Since the dataset \mathcal{D} is finite, we have a small enough number ϵ that satisfies,

$$y_i(\mathbf{w}\mathbf{x}_i + w_0) \geq \epsilon \quad (3.89)$$

Dividing ϵ on both sides, we obtain,

$$y_i\left(\frac{\mathbf{w}}{\epsilon}\mathbf{x}_i + \frac{w_0}{\epsilon}\right) \geq 1 \quad (3.90)$$

Let $\hat{\mathbf{w}} = \frac{\mathbf{w}}{\epsilon}$ and $\hat{w}_0 = \frac{w_0}{\epsilon}$, we have,

$$y_i(\hat{\mathbf{w}}\mathbf{x}_i + \hat{w}_0) \geq 1 \quad (3.91)$$

To simplify the notation, we have,

$$y_i(\mathbf{w}\mathbf{x}_i + w_0) \geq 1 \quad (3.92)$$

Recall that the margin of data point \mathbf{x}^* to hyper-plane is $\frac{|\mathbf{w}\mathbf{x}^* + w_0|}{\|\mathbf{w}\|_2}$. Thus, we have,

$$\frac{|\mathbf{w}\mathbf{x}^* + w_0|}{\|\mathbf{w}\|_2} = \frac{y^*(\mathbf{w}\mathbf{x}^* + w_0)}{\|\mathbf{w}\|_2} \geq \frac{1}{\|\mathbf{w}\|_2} \quad (3.93)$$

Since the maximum margin principle is to pick the boundary that maximises the margin, we have the optimization problem given by,

$$\begin{aligned} & \text{maximize} \quad \frac{1}{\|\mathbf{w}\|_2} \\ & \text{s.t. } y_i(\mathbf{w}\mathbf{x}_i + w_0) \geq 1, \quad i \in [N] \end{aligned} \tag{3.94}$$

Equivalently, the optimization problem can be rephrased as,

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2}\|\mathbf{w}\|_2^2 \\ & \text{s.t. } y_i(\mathbf{w}\mathbf{x}_i + w_0) \geq 1, \quad i \in [N] \end{aligned} \tag{3.95}$$

The constant factor $\frac{1}{2}$ is added without loss of generality.

We can easily observe that the objective is a convex quadratic function and the constraints are convex linear functions. Also, the feasible set as defined by the constraints is convex as well. Thus, our optimisation problem is convex quadratic which can be solved using generic convex optimisation methods.

If our data is linearly-separable, we will find the SVM classifier with no classification error on the training set.

However, if our data is not linearly-separable, as is shown in Figure 3.9, the quadratic program from the previous slides has no feasible solution.

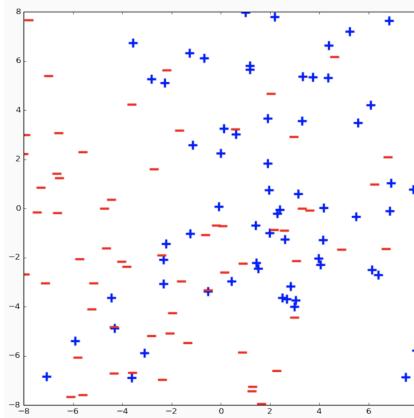


Figure 3.9: Linearly Non-separable Data

We can solve this problem by adding a relaxation term to the previous optimisation problem and find a separator that makes the least mistakes on the training error. However, minimising the number of misclassifications is NP-hard. Equivalently, we can use a proxy for least mistakes, that is to satisfy as many of the N constraints as possible.

Alternatively, we allow all constraints to be satisfied with some slack, given by,

$$\begin{aligned} \text{minimize } & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \zeta_i \\ \text{s.t. } & y_i(\mathbf{w}\mathbf{x}_i + w_0) \geq 1 - \zeta_i, \quad i \in [N] \\ & \zeta_i \geq 0 \quad i \in [N] \end{aligned} \tag{3.96}$$

Feasible solution to the optimization problem always exists thanks to slack variables ζ_i and ζ_i can be as large, i.e. constraints can be violated as much, as necessary.

Why we have $\zeta_i \geq 0$?

Assume $\zeta_i \ll 0$ and \mathbf{x}_i is correctly classified by a huge margin. This would compensate the penalty incurred on misclassified points. That is to say, positive ζ_i which is induced by misclassified points and negative ζ_i which is induced by correctly classified points will cancel each other out. Thus, we need $\zeta_i \geq 0$ to ensure there is no bonus for correct classification by a huge margin.

3.4.2 Hinge Loss Optimization

The optimal solution to the optimization problem above must satisfy either $\zeta_i = 0$, which is the ideal case, where no slack is needed to classify \mathbf{x}_i or that ζ_i is the minimum which satisfies the constraint $\zeta_i = 1 - y_i(\mathbf{w}\mathbf{x}_i + w_0) \geq 0$. Equivalently,

$$\zeta_i = \max(0, 1 - y_i(\mathbf{w}\mathbf{x}_i + w_0)) \stackrel{\text{def}}{=} \ell_{\text{hinge}}(\mathbf{w}, w_0; \mathbf{x}_i, y_i) \tag{3.97}$$

Our SVM formulation becomes equivalent to minimising the objective function,

$$\mathcal{L}_{SVM} = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \ell_{\text{hinge}}(\mathbf{w}, w_0; \mathbf{x}_i, y_i) \tag{3.98}$$

which is the Hinge loss with ℓ_2 regularization.

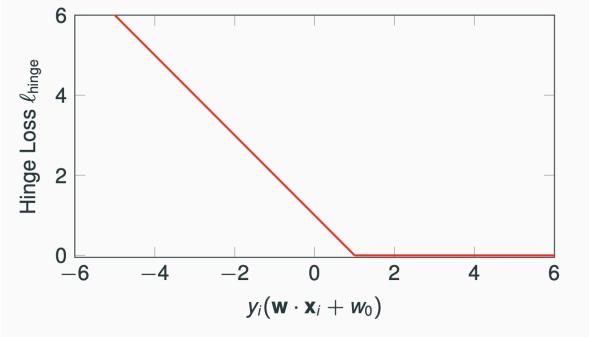


Figure 3.10: Hinge Loss

Recall the the logistic loss function, we can see some similarities between them as is shown in Figure 3.11.

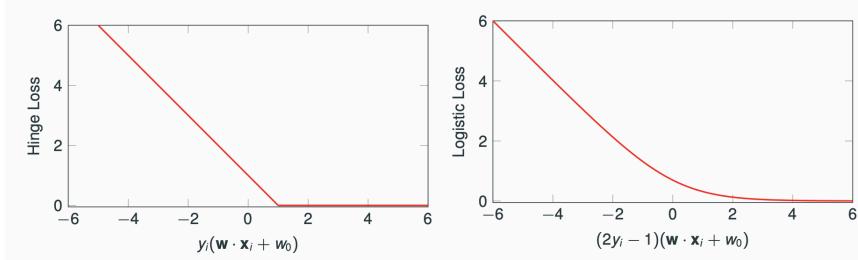


Figure 3.11: Hinge Loss v.s. Logistic Loss

3.4.3 Dual SVM Formulation

Constrained Optimisation with Inequalities

Recall the primary form of the optimization problem

$$\begin{aligned} \min \quad & F(\mathbf{z}) \\ \text{s.t.} \quad & g_i(\mathbf{z}) \geq 0, \quad i = 1, \dots, m \\ & h_j(\mathbf{z}) = 0, \quad j = 1, \dots, \ell \end{aligned} .$$

The Lagrange is

$$\Lambda(\mathbf{z}; \alpha, \mu) = F(\mathbf{z}) - \sum_{i=1}^m \alpha_i g_i(\mathbf{z}) - \sum_{j=1}^{\ell} \mu_j h_j(\mathbf{z})$$

For convex problems, Karush-Kuhn-Tucker (KKT) conditions:

- Dual feasibility:

$$\alpha_i > 0 \quad \text{for } i = 1, \dots, m$$

- Primal feasibility:

$$g_i(\mathbf{z}) \geq 0 \quad \text{for } i = 1, \dots, m$$

$$h_j(\mathbf{z}) = 0 \quad \text{for } j = 1, \dots, \ell$$

- Complementary slackness

$$\alpha_i g_i(\mathbf{z}) = 0 \quad \text{for } i = 1, \dots, m$$

are necessary and sufficient for a critical point of Λ to be the minimum of the original constrained optimisation.

For non-convex problems: KKT are necessary but not sufficient.

Dual SVM Formulation: The Linearly Separable Case

Recall our optimization problem in SVM.

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{s.t. } y_i(\mathbf{w}\mathbf{x}_i + w_0) \geq 1, \quad i \in [N] \end{aligned}$$

The Lagrange function is

$$\Lambda(\mathbf{w}; w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i (y_i(\mathbf{w}\mathbf{x}_i + w_0) - 1)$$

and its gradients with respect to the primal variables are

$$\begin{aligned} \frac{\partial \Lambda}{\partial w_0} &= - \sum_{i=1}^N \alpha_i y_i \\ \nabla_{\mathbf{w}} \Lambda &= \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \end{aligned}$$

At optimality, these gradients are 0, yielding

$$\begin{aligned} \sum_{i=1}^N \alpha_i y_i &= 0 \\ \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i &= 0 \end{aligned}$$

Optimal \mathbf{w} obtained as

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

At optimality

$$y_i(\mathbf{w}\mathbf{x}_i + w_0) = 1 \quad \text{for } i \in [N]$$

We can obtain w_0 using the above equality constraint for any i :

$$y_i^2 \left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x}_i + w_0 \right) = y_i$$

Since $y_i^2 = (\pm 1)^2 = 1$,

$$\sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x}_i + w_0 = y_i$$

Thus,

$$w_0 = y_i - \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x}_i$$

It is numerically more stable to have w_0 the average over all possible values:

$$w_0 = \frac{1}{N} \sum_{i=1}^N \left(y_i - \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x}_i \right)$$

Next we plug the optimal \mathbf{w} and constraint

$$\sum_{i=1}^N \alpha_i y_i = 0$$

into Lagrangian.

$$\begin{aligned} \Lambda(\mathbf{w}; w_0, \alpha) &= \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^N \alpha_i (y_i(\mathbf{w}\mathbf{x}_i + w_0) - 1) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{w} \cdot \mathbf{x}_i - \sum_{i=1}^N \alpha_i y_i w_0 + \sum_{i=1}^N \alpha_i \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j (\mathbf{x}_i \cdot \mathbf{x}_j) + \sum_{i=1}^N \alpha_i \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \\ &= g(\alpha) \end{aligned} \tag{3.99}$$

To find critical points of Λ satisfying the KKT conditions, it is sufficient to find the critical points of g that satisfy the constraints:

$$\alpha_i \geq 0$$

and

$$\sum_{i=1}^N \alpha_i y_i = 0$$

Thus, the dual form SVM optimization problem becomes:

$$\begin{aligned} & \text{maximize} \quad \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j (\mathbf{x}_i \mathbf{x}_j) \\ & \text{s.t.} \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \geq 0 \quad \text{for } i \in [N] \end{aligned}$$

This is quadratic concave problem with N variables and very simple box constraints plus one zero-sum constraint.

There are several reasons why SVM dual formulation is preferred.

- $D \approx D$ or D : Basis expansion to capture non-linear discriminating boundaries
- Natural kernelisation: Replace dot product $\mathbf{x}_i \mathbf{x}_j$ by kernel $\kappa(\mathbf{x}_i, \mathbf{x}_j)$
- The number of non-zero variables α_i can be much smaller in practice. The complementary slackness conditions in KKT

$$\alpha_i (y_i (\mathbf{w} \mathbf{x}_i + w_0) - 1) = 0$$

are satisfied when $\alpha_i = 0$ or $y_i (\mathbf{w} \mathbf{x}_i + w_0) = 1$. Thus, only the points \mathbf{x}_i with $\alpha_i > 0$ contribute to the solution \mathbf{w}

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

That is to say such vectors \mathbf{x}_i form the support of the solution, i.e. are support vectors. Such vectors \mathbf{x}_i are points with least margin.

Dual SVM Formulation: The Non-Linearly Separable Case

Recall the primary objective in the non-linearly separable case.

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \zeta_i \\ & \text{s.t.} \quad y_i (\mathbf{w} \mathbf{x}_i + w_0) \geq 1 - \zeta_i, \quad i \in [N] \\ & \quad \zeta_i \geq 0 \quad i \in [N] \end{aligned} \tag{3.100}$$

The Lagrange is

$$\Lambda(\mathbf{w}; w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \zeta_i - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w} \mathbf{x}_i + w_0) - (1 - \zeta_i)) - \sum_{i=1}^N \mu_i \zeta_i$$

The derivatives with respect to \mathbf{w} , w_0 and ζ_i are:

$$\nabla_{\mathbf{w}} \Lambda = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial \Lambda}{\partial w_0} = - \sum_{i=1}^N \alpha_i y_i$$

$$\frac{\partial \Lambda}{\partial \zeta_i} = C - \alpha_i - \mu_i$$

For (KKT) dual feasibility constraints, we require $\alpha \geq 0$ and $\mu_i \geq 0$.

Setting the derivatives to 0, substituting the resulting expressions in Λ (and simplifying), we get the optimization problem:

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \\ & \text{s.t.} && \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C \quad \text{for } i \in [N] \end{aligned}$$

Making Predictions using SVM Dual

Recall the optimal solution for \mathbf{w} expressed using the dual variables $\boldsymbol{\alpha}$

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

The bias w_0 is also expressible using $\boldsymbol{\alpha}$.

Prediction on a new point \mathbf{x}_{new} requires inner products with the support vectors

$$\mathbf{w} \mathbf{x}_{new} = \sum_{i=1}^N \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}_{new})$$

We can as well use blackbox access to a function $\kappa(\cdot, \cdot)$ that maps two inputs \mathbf{x} , \mathbf{x}' to their inner product $\mathbf{x} \cdot \mathbf{x}'$. This is a kernel function

3.4.4 Kernel Methods

A function κ is a kernel that computes the dot product for some expansion ϕ :

$$\kappa(\mathbf{x}', \mathbf{x}) : \chi \times \chi \rightarrow \mathbb{R}, \quad \kappa(\mathbf{x}', \mathbf{x}) = \phi(\mathbf{x}') \phi(\mathbf{x})$$

$\kappa(\mathbf{x}', \mathbf{x})$ is some measure of similarity between \mathbf{x} and \mathbf{x}' .

Mercer Kernels

The gram matrix of a kernel is defined as,

$$\mathbf{K} = \begin{pmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \kappa(\mathbf{x}_1, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_N) \\ \kappa(\mathbf{x}_2, \mathbf{x}_1) & \kappa(\mathbf{x}_2, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_1) & \kappa(\mathbf{x}_N, \mathbf{x}_2) & \cdots & \kappa(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix} \quad (1)$$

A kernel is Mercer or positive definite kernel if the Gram matrix is always positive semi-definite.

Kernel

$$\kappa(\mathbf{x}', \mathbf{x}) = \mathbf{x}'^T \mathbf{x}$$

is a Mercer kernel.

Kernel Engineering

We can build kernels using simpler kernels as building blocks.

Given kernels κ_1, κ_2 , the following are kernels:

- $\kappa(\mathbf{x}', \mathbf{x}) = c\kappa_1(\mathbf{x}', \mathbf{x})$ where $c > 0$
- $\kappa(\mathbf{x}', \mathbf{x}) = f(\mathbf{x}')\kappa_1(\mathbf{x}', \mathbf{x})f(\mathbf{x})$ where f is a function.
- $\kappa(\mathbf{x}', \mathbf{x}) = q[\kappa_1(\mathbf{x}', \mathbf{x})]$ where q is a polynomial with non-negative coefficients.
- $\kappa(\mathbf{x}', \mathbf{x}) = e^{\kappa_1(\mathbf{x}', \mathbf{x})}$
- $\kappa(\mathbf{x}', \mathbf{x}) = \kappa_1(\mathbf{x}', \mathbf{x})\kappa_2(\mathbf{x}', \mathbf{x})$
- $\kappa(\mathbf{x}', \mathbf{x}) = \kappa_a(\mathbf{x}'_a, \mathbf{x}_a)\kappa_b(\mathbf{x}'_b, \mathbf{x}_b)$ where $\mathbf{x}' = (\mathbf{x}'_a, \mathbf{x}'_b)$, $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, κ_a, κ_b are kernels.
- $\kappa(\mathbf{x}', \mathbf{x}) = \kappa_a(\mathbf{x}'_a, \mathbf{x}_a) + \kappa_b(\mathbf{x}'_b, \mathbf{x}_b)$
- $\kappa(\mathbf{x}', \mathbf{x}) = \mathbf{x}'^T \mathbf{A} \mathbf{x}$ where \mathbf{A} is a symmetric positive semi-definite matrix.

For example, we can use kernel engineering to show that RBF kernel κ_{RBF} is kernel by:

$$\begin{aligned} \kappa_{RBF} &= e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}} = e^{\frac{-\mathbf{x}^T \mathbf{x}}{2\sigma^2}} e^{\frac{\mathbf{x}^T \mathbf{x}'}{\sigma^2}} e^{\frac{-\mathbf{x}'^T \mathbf{x}'}{2\sigma^2}} \\ &= f(\mathbf{x}) e^{c\kappa_1(\mathbf{x}', \mathbf{x})} f(\mathbf{x}') \end{aligned} \quad (3.101)$$

4

Neural Network

4.1 Feed-forward Neural Networks

A perceptron is an algorithm for supervised learning of binary classifiers. It has an input of a feature vector \mathbf{x} and a label y and outputs $\text{sign}(\mathbf{w}\mathbf{x} + b)$, where \mathbf{w} is the weight parameter vector and b is the bias.

An artificial neuron is the generalisation of the perceptron. For any activation function $f : \mathbb{R} \rightarrow \mathbb{R}$, it has an input of a feature vector \mathbf{x} and a label y and outputs $f(\mathbf{w}\mathbf{x} + b)$.

An artificial neuron in a neural network computes a linear function of its input and is then composed with a non-linear activation function, as is shown in Figure 4.1.

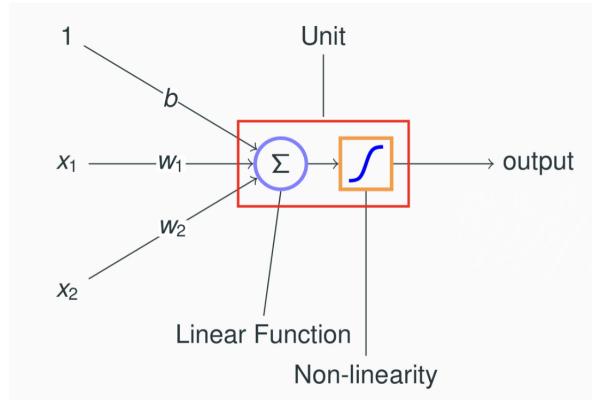


Figure 4.1: Neural Network with One Artificial Neuron

The output for regression is $\hat{y} = \mathbb{E}[y | \mathbf{x}, \mathbf{W}, \mathbf{b}]$ and the output for classification is $\hat{y} = Pr(y = 1 | \mathbf{x}, \mathbf{w}, b)$.

For logistic regression, the non-linear activation function is the sigmoid

function,

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (4.1)$$

and the separating surface is linear.

An artificial neural network refers to the composition of artificial neurons into a network, where the artificial neurons have continuous activation functions and the network has a well-chosen loss function that, for any input-output, is differentiable.

Artificial neural networks are typically composed of input layers, hidden layers and output layers. A hidden layer is located between the input and output of the algorithm, in which the function applies weights to the inputs and directs them through an activation function as the output. In short, the hidden layers perform nonlinear transformations of the inputs entered into the network.

In theory, the neural networks can compute any function that a computer can.

A feed-forward network is a network with no recurrent connections. It is an important distinction because in a feed forward network the gradient is clearly defined and computable through backpropagation (i.e. chain rule).

A fully-connected network, or a fully-connected layer in a network, is one such that neurons are organised in a layered directed-acyclic graph, i.e. no loops. This, for example, contrasts with convolutional layers, where each output neuron depends on a subset of the input neurons.

Figure 3.2 shows a fully-connected feed-forward network.

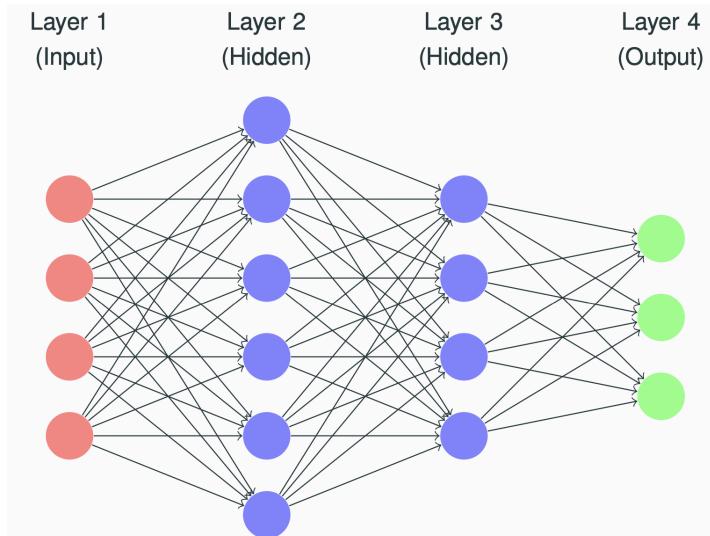


Figure 4.2: Fully-connected Feed-forward Network

4.1.1 Multi-layer Perceptrons

Multi-layer perceptron(MLP) is a type of feed-forward networks. Though it is called multi-layer perceptron, the units composed of it are not necessarily perceptrons.

Figure 4.3 shows a three-layered MLP. It first maps the feature vector \mathbf{x} to non-linear features and then applies logistic regression to compute \hat{y} .

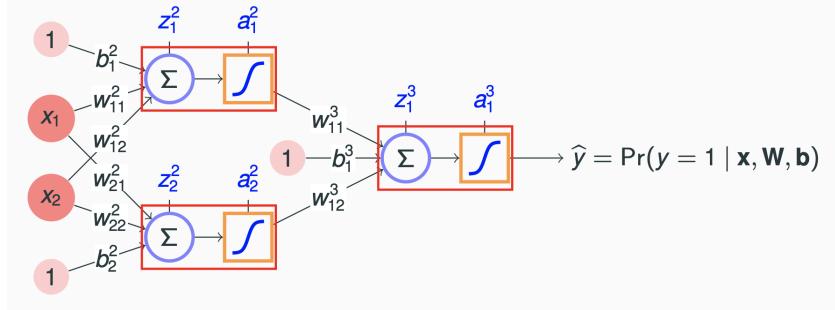


Figure 4.3: Three-Layered MLP

Let's first look at some notations used in nerual networks.

- w_{ij}^ℓ : weight for the connection to neuron i in layer ℓ from neuron j in layer $\ell - 1$
- b_i^ℓ : bias term for neuron i in layer in layer ℓ
- z_i^ℓ : pre-activation for neuron i in layer in layer ℓ , i.e. sum of weighted outputs from layer $\ell - 1$
- a_i^ℓ : activation for neuron i in layer in layer ℓ , i.e. non-linear function of pre-activation z_i^ℓ

Consider the three-layered MLP shown in Figure 4.3

In layer 1, i.e. the input layer, we have pre-activation = activation = input feature vector, given as,

$$\mathbf{a}^1 = \mathbf{z}^1 = \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad (4.2)$$

In layer 2, i.e. the hidden layer, we have the weight matrix:

$$\mathbf{W}^2 = \begin{pmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \end{pmatrix} \quad (4.3)$$

bias vector:

$$\mathbf{b}^2 = \begin{pmatrix} b_1^2 \\ b_2^2 \end{pmatrix} \quad (4.4)$$

pre-activation:

$$\mathbf{z}^2 = \begin{pmatrix} z_1^2 \\ z_2^2 \end{pmatrix} = \mathbf{W}^2 \mathbf{a}^1 + \mathbf{b}^2 = \begin{pmatrix} w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2 \\ w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2 \end{pmatrix} \quad (4.5)$$

and activation:

$$\mathbf{a}^2 = \begin{pmatrix} a_1^2 \\ a_2^2 \end{pmatrix} = \tanh(\mathbf{z}^2) = \begin{pmatrix} \tanh(w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2) \\ \tanh(w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2) \end{pmatrix} \quad (4.6)$$

In layer 3, i.e. the output layer, we have the weight matrix:

$$\mathbf{W}^3 = \begin{pmatrix} w_{11}^3 & w_{12}^3 \end{pmatrix} \quad (4.7)$$

bias vector:

$$\mathbf{b}^3 = \begin{pmatrix} b_1^3 \end{pmatrix} \quad (4.8)$$

preactivation:

$$\mathbf{z}^3 = \begin{pmatrix} z_1^3 \end{pmatrix} = \mathbf{W}^3 \mathbf{a}^2 + \mathbf{b}^3 = \begin{pmatrix} w_{11}^3 a_1^2 + w_{12}^3 a_2^2 + b_1^3 \end{pmatrix} \quad (4.9)$$

and activation:

$$\mathbf{a}^3 = \begin{pmatrix} a_1^3 \end{pmatrix} = \sigma(\mathbf{z}^3) = \begin{pmatrix} \sigma(w_{11}^3 a_1^2 + w_{12}^3 a_2^2 + b_1^3) \end{pmatrix} \quad (4.10)$$

The three-layered MLP gives the output:

$$\hat{y} = \mathbf{a}^3 = \sigma(\mathbf{z}^3) \quad (4.11)$$

Suppose we have the data shown in Figure 4.4 and apply logistic regression to classifier these data points into two classes: red and blue points.

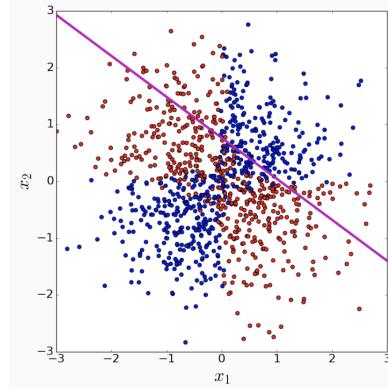
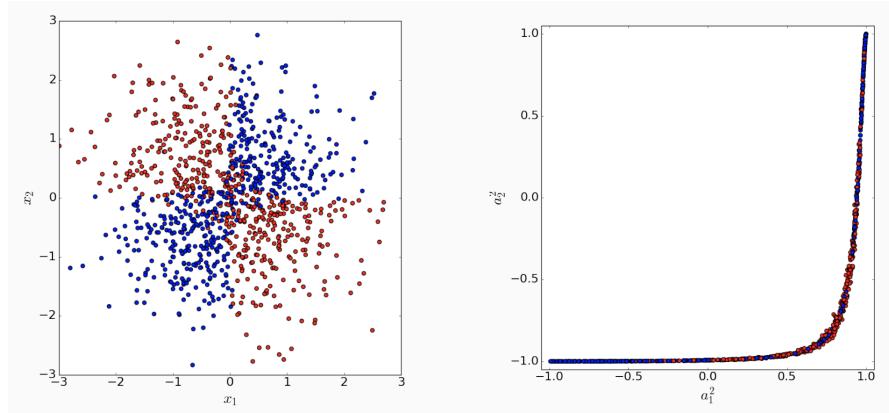


Figure 4.4: Data with Logistic Regression

We can see that the data is not linearly-separable, i.e. no linear separator can separate the blue from the red. Thus, the accuracy of logistic regression $\approx 50\%$.

If we apply our three-layered MPL to classify the data and scatter plot the input feature and the activation of the hidden layer, i.e. layer 2, the results are shown in Figure 4.5.

Figure 4.5: Input Feature \mathbf{x} v.s. The Activation of the Hidden Layer \mathbf{a}^2

We can see that what the hidden do is to Mapping input features \mathbf{x} to non-linear features \mathbf{a}^2 , where

$$\mathbf{a}^2 = \begin{pmatrix} a_1^2 \\ a_2^2 \end{pmatrix} = \tanh(\mathbf{W}^2 \mathbf{a}^1 + \mathbf{b}^2) = \begin{pmatrix} \tanh(w_{11}^2 x_1 + w_{12}^2 x_2 + b_1^2) \\ \tanh(w_{21}^2 x_1 + w_{22}^2 x_2 + b_2^2) \end{pmatrix} \quad (4.12)$$

The non-linear features \mathbf{a}^2 after the mapping becomes linearly-separable. We can find a linear separator to separate them, as is shown in Figure 4.6.

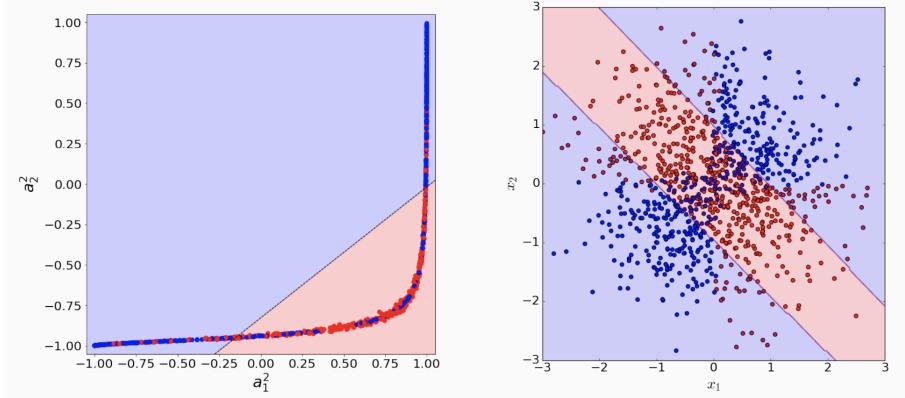
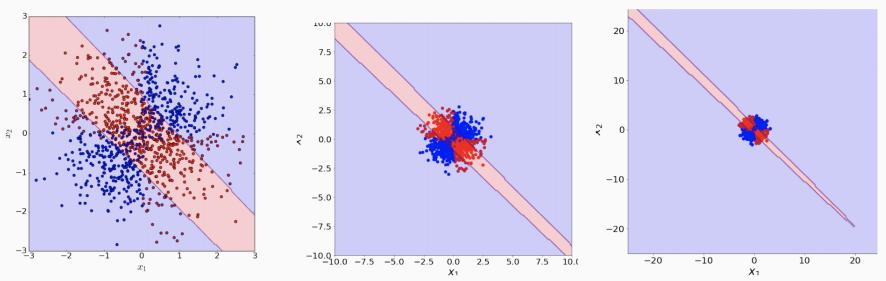


Figure 4.6: Data after Mapping with Linear Separator

We can easily see that the linear separator in the (a_1, a_2) plane becomes non-linear in the (x_1, x_2) plane.

Figure 4.7: Separator in (x_1, x_2) Plane

4.1.2 Back-Propagation

Now let's consider how to learn the MLP parameters.

We first compute the derivatives of the loss function with respect to all parameters of the MLP for each data point. Then, we average the derivatives over the training data points. Since this can be computationally very expensive in practice, we commonly only average over a mini-batch instead of the entire dataset. Finally, we perform a gradient descent step to update the model parameters.

Recall the chain rule for multivariate calculus.

Suppose we have functions: $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$, $g : \mathbb{R}^k \rightarrow \mathbb{R}^m$, $h = f \circ g : \mathbb{R}^n \rightarrow \mathbb{R}^m$. For $\mathbf{x} \in \mathbb{R}^n$, let $\mathbf{z} = f(\mathbf{x})$,

$$\frac{\partial h}{\partial \mathbf{x}} = \frac{\partial h}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \quad (4.13)$$

where $\frac{\partial h}{\partial \mathbf{x}}$ is an $m \times n$ Jacobian matrix, $\frac{\partial h}{\partial \mathbf{z}}$ is an $m \times k$ Jacobian matrix and $\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ is an $k \times n$ Jacobian matrix.

Figure 4.8 shows an example of a three-layered MLP.

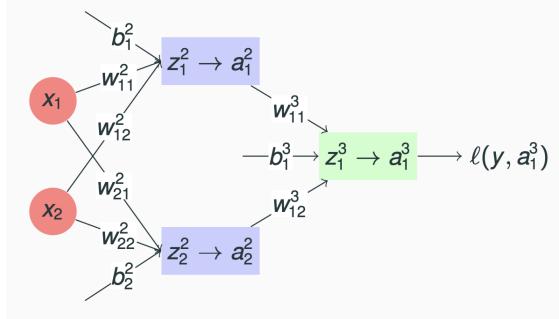


Figure 4.8: Separator in (x_1, x_2) Plane

We want the derivatives of the loss function with respect to the weights $\frac{\partial \ell}{\partial \mathbf{W}^\ell}$ and bias $\frac{\partial \ell}{\partial \mathbf{b}^\ell}$ in each layer, i.e. $\frac{\partial \ell}{\partial \mathbf{W}^2}$, $\frac{\partial \ell}{\partial \mathbf{b}^2}$, $\frac{\partial \ell}{\partial \mathbf{W}^3}$, and $\frac{\partial \ell}{\partial \mathbf{b}^3}$. It suffices to compute the derivatives of the loss function with respect to the activation $\frac{\partial \ell}{\partial \mathbf{z}^\ell}$ at each layer and then derive $\frac{\partial \ell}{\partial \mathbf{W}^\ell}$ and $\frac{\partial \ell}{\partial \mathbf{b}^\ell}$ from $\frac{\partial \ell}{\partial \mathbf{z}^\ell}$ using the chain rule.

Recall for a vector $\mathbf{z} \in \mathbb{R}^n$ and a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the derivative $\frac{\partial f}{\partial \mathbf{z}} \in \mathbb{R}^n$ is a row vector defined as,

$$\frac{\partial f}{\partial \mathbf{z}} = \left(\frac{\partial f}{\partial z_1} \quad \dots \quad \frac{\partial f}{\partial z_n} \right) \quad (4.14)$$

For a matrix $\mathbf{W} \in \mathbb{R}^{n \times m}$ and a function $f : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$, the derivative $\frac{\partial f}{\partial \mathbf{W}} \in \mathbb{R}; n \times m$ is given as,

$$\frac{\partial f}{\partial \mathbf{W}} = \begin{pmatrix} \frac{\partial f}{\partial w_{11}} & \dots & \frac{\partial f}{\partial w_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial w_{n1}} & \dots & \frac{\partial f}{\partial w_{nm}} \end{pmatrix}_{n \times m} \quad (4.15)$$

Let's now compute the derivatives in our three-layered MLP.

The derivatives of the loss function with respect to the weights in layer 2 is given by,

$$\frac{\partial \ell}{\partial \mathbf{W}^2} = \begin{pmatrix} \frac{\partial \ell}{\partial w_{11}^2} & \frac{\partial \ell}{\partial w_{12}^2} \\ \frac{\partial \ell}{\partial w_{21}^2} & \frac{\partial \ell}{\partial w_{22}^2} \end{pmatrix}_{2 \times 2} \quad (4.16)$$

Assume that $\frac{\partial \ell}{\partial \mathbf{z}^2}$ is known. Using the chain rule, we have:

$$\frac{\partial \ell}{\partial w_{ij}^2} = \frac{\partial \ell}{\partial \mathbf{z}_i^2} \frac{\partial \mathbf{z}_i^2}{\partial w_{ij}^2} \quad (4.17)$$

Recall that,

$$z_i^2 = \sum_j w_{ij}^2 x_j + b_i^2, \quad i, j \in \{1, 2\} \quad (4.18)$$

Thus,

$$\frac{\partial \ell}{\partial w_{ij}^2} = \frac{\partial \ell}{\partial z_i^2} x_j \quad (4.19)$$

Then,

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{W}^2} &= \begin{pmatrix} \frac{\partial \ell}{\partial z_1^2} x_1 & \frac{\partial \ell}{\partial z_2^2} x_2 \\ \frac{\partial \ell}{\partial z_2^2} x_1 & \frac{\partial \ell}{\partial z_2^2} x_2 \end{pmatrix}_{2 \times 2} \\ &= \left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \begin{pmatrix} \frac{\partial \ell}{\partial z_1^2} & \frac{\partial \ell}{\partial z_2^2} \end{pmatrix} \right)^T \quad (4.20) \\ &= \left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \frac{\partial \ell}{\partial \mathbf{z}^2} \right)^T \end{aligned}$$

Since $\begin{pmatrix} x_1 & x_2 \end{pmatrix}^T = \mathbf{a}^1$, the equation above can be rephrased as,

$$\frac{\partial \ell}{\partial \mathbf{W}^2} = (\mathbf{a}^1 \frac{\partial \ell}{\partial \mathbf{z}^2})^T \quad (4.21)$$

The derivatives of the loss function with respect to the bias in layer 2 is given by,

$$\frac{\partial \ell}{\partial \mathbf{b}^2} = \begin{pmatrix} b_1^2 & b_2^2 \end{pmatrix} = \begin{pmatrix} \frac{\partial \ell}{\partial b_1^2} & \frac{\partial \ell}{\partial b_2^2} \end{pmatrix} \quad (4.22)$$

Since,

$$\frac{\partial \ell}{\partial b_i^2} = \frac{\partial \ell}{\partial z_i^2} \times 1 \quad (4.23)$$

we have,

$$\frac{\partial \ell}{\partial \mathbf{b}^2} = \begin{pmatrix} \frac{\partial \ell}{\partial z_1^2} & \frac{\partial \ell}{\partial z_2^2} \end{pmatrix} = \frac{\partial \ell}{\partial \mathbf{z}^2} \quad (4.24)$$

Similarly, we can compute the derivatives of the loss function with respect to the weights and bias in layer 3 given by,

$$\frac{\partial \ell}{\partial \mathbf{W}^3} = (\mathbf{a}^2 \frac{\partial \ell}{\partial \mathbf{z}^3})^T \quad (4.25)$$

$$\frac{\partial \ell}{\partial \mathbf{b}^3} = \frac{\partial \ell}{\partial \mathbf{z}^3} \quad (4.26)$$

each respectively.

Suppose we use the cross-entropy loss function in our example,

$$\ell(\mathbf{a}^3, y) = -[y \log a_1^3 + (1 - y) \log(1 - a_1^3)] \quad (4.27)$$

Then we have,

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{a}^3} &= \left(\frac{\partial \ell}{\partial a_1^3} \right) \\ &= \left(-[\frac{y}{a_1^3} + \frac{1-y}{1-a_1^3}(-1)] \right) \\ &= \left(\frac{a_1^3 - y}{a_1^3(1-a_1^3)} \right) \end{aligned} \quad (4.28)$$

Our choice of the activation function is $a_1^3 = \sigma(z_1^3)$. Thus,

$$\frac{\partial \mathbf{a}^3}{\partial \mathbf{z}^3} = \left(\sigma'(z_1^3) \right) = \left(a_1^3(1 - a_1^3) \right) \quad (4.29)$$

Then,

$$\frac{\partial \ell}{\partial \mathbf{z}^3} = \frac{\partial \ell}{\partial \mathbf{a}^3} \frac{\partial \mathbf{a}^3}{\partial \mathbf{z}^3} = \left(a_1^3 - y \right) \quad (4.30)$$

Now we compute $\frac{\partial \ell}{\partial \mathbf{z}^2}$.

Since,

$$z_i^3 = w_{ij}^3 a_j^2 + b_i^3 \quad (4.31)$$

we have,

$$\frac{\partial z_i^3}{\partial a_j^2} = w_{ij}^3 \quad (4.32)$$

Thus,

$$\frac{\partial \mathbf{z}^3}{\partial \mathbf{a}^2} = \mathbf{W}^3 \quad (4.33)$$

Recall for function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, the derivative $\frac{\partial f}{\partial \mathbf{z}} \in \mathbb{R}^{m \times n}$ is given as,

$$\frac{\partial \mathbf{f}}{\partial \mathbf{z}} = \begin{pmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & \frac{\partial f_1}{\partial z_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial z_1} & \cdots & \frac{\partial f_m}{\partial z_n} \end{pmatrix}_{m \times n} \quad (4.34)$$

Since $a_i^2 = \tanh z_i^2$, we have

$$\frac{\partial \mathbf{a}^2}{\partial \mathbf{z}^2} = \begin{pmatrix} \frac{\partial a_1^2}{\partial z_1^2} & \frac{\partial a_1^2}{\partial z_2^2} \\ \frac{\partial a_2^2}{\partial z_1^2} & \frac{\partial a_2^2}{\partial z_2^2} \end{pmatrix} = \begin{pmatrix} 1 - (a_1^2)^2 & 0 \\ 0 & 1 - (a_1^2)^2 \end{pmatrix} \quad (4.35)$$

Thus,

$$\frac{\partial \ell}{\partial z^2} = \frac{\partial \ell}{\partial z^3} \frac{\partial z^3}{\partial a^2} \frac{\partial a^2}{\partial z^2} = (a_1^3 - y) W^3 \begin{pmatrix} 1 - (a_1^2)^2 & 0 \\ 0 & 1 - (a_1^2)^2 \end{pmatrix} \quad (4.36)$$

More generally, consider that all layers are fully connected, i.e. every neuron in layer $\ell - 1$ has a connection with every unit in layer ℓ , and each layer consists of a linear function and non-linear activation.

We have the following forward equations:

$$\mathbf{a}^1 = \mathbf{x}_{input} \quad (4.37)$$

$$\mathbf{z}^\ell = \mathbf{W}^\ell \mathbf{a}^{\ell-1} + \mathbf{b}^\ell \quad (4.38)$$

$$\mathbf{a}^\ell = f_\ell(\mathbf{z}^\ell) \quad (4.39)$$

where f_ℓ may map \mathbf{z}^ℓ to \mathbf{a}^ℓ element-wise or vector-wise.

$$\ell(\mathbf{a}^L, y) = - \sum_{c=1}^C \mathbb{1}(y = c) \log a_c^L \quad (4.40)$$

where ℓ is the cross entropy loss function.

If there are n_L (output) neurons in layer L , then $\frac{\partial \ell}{\partial a^L}$ and $\frac{\partial \ell}{\partial z^L}$ are row vectors with n_L elements and $\frac{\partial a^L}{\partial z^L}$ is the $n_L \times n_L$ Jacobian matrix given by,

$$\begin{pmatrix} \frac{\partial a_1^L}{\partial z_1^L} & \cdots & \frac{\partial a_1^L}{\partial z_{n_L}^L} \\ \vdots & \ddots & \vdots \\ \frac{\partial a_{n_L}^L}{\partial z_1^L} & \cdots & \frac{\partial a_{n_L}^L}{\partial z_{n_L}^L} \end{pmatrix}$$

If f_L is applied element-wise, e.g., sigmoid, then this matrix is diagonal.

To compute the derivatives with respect to the network parameters, we have the following back-propagation equations:

$$\frac{\partial \ell}{\partial z^L} = \frac{\partial \ell}{\partial a^L} \frac{\partial a^L}{\partial z^L} \quad (4.41)$$

where layer L is the last layer, i.e. output layer.

$$\frac{\partial \ell}{\partial z^\ell} = \frac{\partial \ell}{\partial z^{\ell+1}} \mathbf{W}^{\ell+1} \frac{\partial a^{\ell+1}}{\partial z^\ell} \quad (4.42)$$

$$\frac{\partial \ell}{\partial \mathbf{W}^\ell} = (\mathbf{a}^{\ell-1} \frac{\partial \ell}{\partial z^\ell})^T \quad (4.43)$$

$$\frac{\partial \ell}{\partial \mathbf{b}^\ell} = 1 \times \frac{\partial \ell}{\partial z^\ell} = \frac{\partial \ell}{\partial z^\ell} \quad (4.44)$$

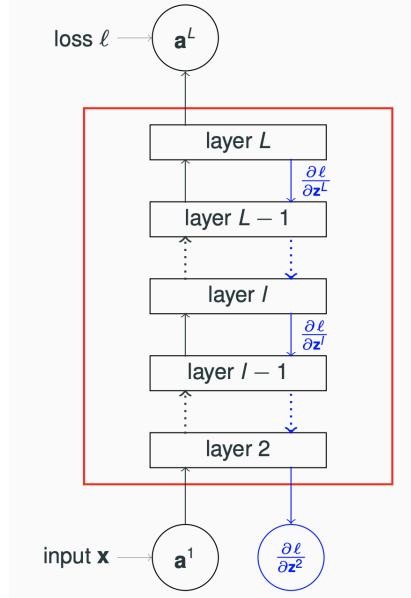


Figure 4.9: Flowchart of Back-Propagation

The running time to compute the gradient for a single data point is as many matrix multiplications as there are fully connected layers and performed twice during forward and backward pass.

The space requirement is big enough to store vectors \mathbf{a}^ℓ , \mathbf{z}^ℓ , and $\frac{\partial \ell}{\partial z^\ell}$ for each layer.

We can process multiple examples together if training on minibatch. In this case, we need to make sure that all parameters fit in GPU memory.

4.2 Training Neural Networks

4.2.1 Initialising Weights and Biases

Initialising is important when minimising non-convex functions. We may get very different results depending on where we start the optimisation.

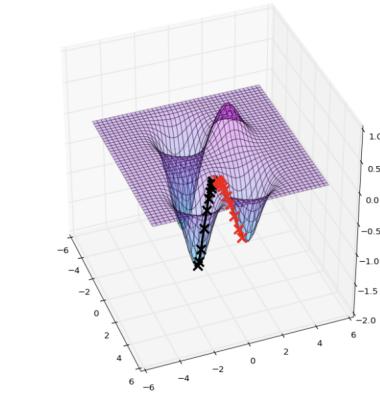


Figure 4.10: Non-Convex Optimization Surface

Weight initialisation for sigmoid/ReLU units:

- Suppose there are D weights w_1, \dots, w_D .
- Draw w_i randomly from $\mathcal{N}(0, \frac{1}{D})$.

Bias initialisation:

- For sigmoid, use a random value around 0.
- For ReLU, use a small positive constant

4.2.2 Saturation and Vanishing Gradient

Saturation happens when pre-activations are in range where the activation function is very flat.

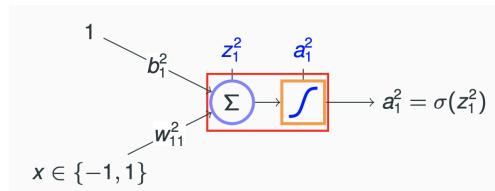


Figure 4.11: NOT Gate

For example, consider the NOT gate shown in Figure 4.10. If we use the squared loss function $\ell(a_1^2, y) = (a_1^2 - y)^2$, then,

$$\frac{\partial \ell}{\partial z_1^2} = 2(a_1^2 - y) \frac{\partial a_1^2}{\partial z_1^2} = 2(a_1^2 - y)\sigma'(z_1^2) \quad (4.45)$$

Recall that $\sigma'(z_1^2) = \sigma(z_1^2)(1 - \sigma(z_1^2))$. If $x = -1$, $w_{11}^2 = 5$ and $b_1^2 = 0$, then $\sigma'(z_1^2) = 0$. Gradient steps may not make much progress, as the gradient is very small yet the loss function $\ell(a_1^2, y) = (a_1^2 - y)^2$ has a relatively large value.

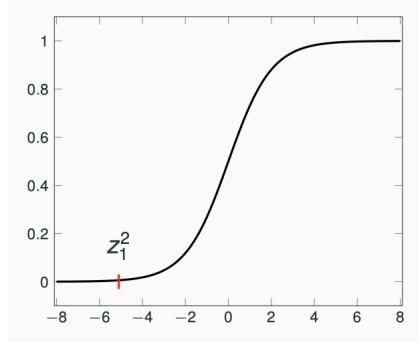


Figure 4.12: Saturation

Sometimes we can fix this problem by changing the squared loss to cross-entropy loss. Recall the cross-entropy loss $\ell(a_1^2, y) = -[y \log a_1^2 + (1 - y) \log(1 - a_1^2)]$. Thus, $\frac{\partial \ell}{\partial a_1^2} = \frac{a_1^2 - y}{a_1^2(1 - a_1^2)}$. Then $\frac{\partial \ell}{\partial z_1^2} = \frac{\partial \ell}{\partial a_1^2} \frac{\partial a_1^2}{\partial z_1^2} = a_1^2 - y$. Large gap between a_1^2 and y also means large magnitude of the gradient.

Saturation is inherent to many activation functions such as sigmoid and tanh.

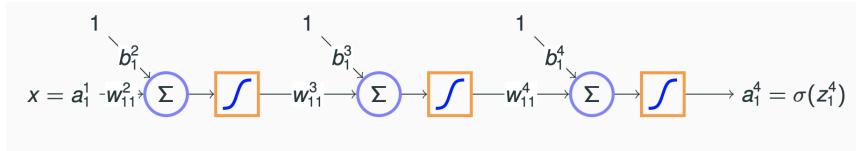


Figure 4.13: Four-layered Neural Network

For neural networks with many layers, as is shown in Figure 4.12, if we use sigmoid function as the activation function for all the hidden layers, we need to multiply many gradients $\sigma'(z_i^\ell)$ where $\sigma'(z) \in [0, \frac{1}{4}]$ when doing back-propagation. Thus, the gradient product \approx in case of many layers.

The product with weights w_{ij}^ℓ may avoid vanishing yet may explode the gradient.

4.2.3 Rectified Linear Unit

Rectified linear unit, i.e. ReLU, is defined as,

$$f(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ z & \text{if } z \geq 0 \end{cases} = \max(0, z) \quad (4.46)$$

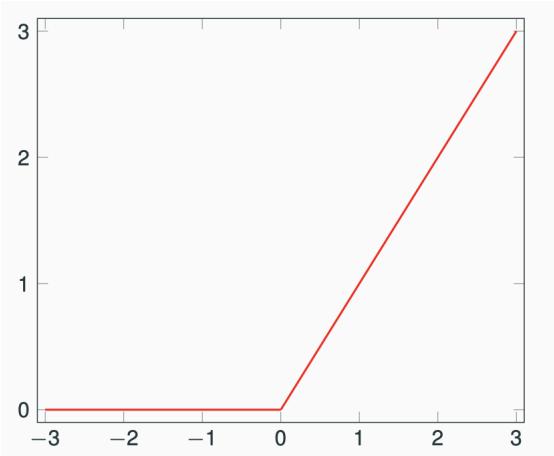


Figure 4.14: Rectified Linear Unit

The rectifier only saturates on one side and has derivative 1 for large z , so there is no saturation. To some extent, this solves the saturation and vanishing gradient problems.

Also, the neurons are sparsely activated, i.e. some might not activate at all. This is often desirable because not all neurons should fire for each input. For example, neurons that identifies animal eyes should not activate for building images. In practice, concise models often have better predictive power and less overfitting.

However, there still exists a dying ReLU problem, i.e. once a neuron gets negative, it is unlikely to recover and always outputs 0 as its pre-activation remains negative. Then those dying neurons become useless and hence a large part of network does nothing.

Possible reasons of the dying ReLU problem is large negative bias or too high learning rate. Thus, lowering the learning rate can solve it to some extent. If not, we can use leaky ReLU, where the slope for negative input is a constant.

Leaky ReLU is given as,

$$f(z) = \begin{cases} cz & \text{if } z \leq 0 \\ z & \text{if } z \geq 0 \end{cases} \quad (4.47)$$

Where c is a constant, usually chosen small.

Leaky ReLU can solves the dying ReLU problem since it does not have zero-slope parts. However, the performance of Leaky ReLU is not always superior to ReLU.

We can also use parametric ReLU where the slope for negative input is a network parameter a , given as,

$$f(z) = \begin{cases} az & \text{if } z \leq 0 \\ z & \text{if } z \geq 0 \end{cases} \quad (4.48)$$

4.2.4 Overfitting and Avoiding Overfitting

Deep neural networks have a lot of parameters. A fully connected layers with n_1, \dots, n_L units have at least $n_1n_2 + n_2n_3 + \dots + n_{L-1}n_L$ parameters. For example, MLP for digit recognition has 2 million parameters and 60,000 training images and for image detection, the famous neural net by Krizhevsky, Sutskever, Hinton(2012), has 60 million parameters and 1.2 million training images.

Beyond regularisation, we can prevent overfitting in deep neural networks by early stopping, adding data and dropout.

Early Stopping

Assume the training uses iterative optimisation methods. The idea of early stopping is to keep aside a validation set and measure performance of the classifier after each gradient step. When validation error stops decreasing, i.e. a local minimum is reached, we stop the training. In this case, we need to compute the validation error after each iteration which is computationally more expensive. Thus, we can do this every few iterations to reduce overhead.

Adding Data

Typically, getting additional data is either impossible or expensive. However, we can fake the data. For example, images can be translated/rotated slightly and their brightness can be changed, etc.



Figure 4.15: Image Data Augmentation

We can also add data by adversarial training:

- Take trained (or partially trained model)
- Create examples by modifications “imperceptible to the human eye”, but where the model fails
- Linear perturbation of non-linear models:

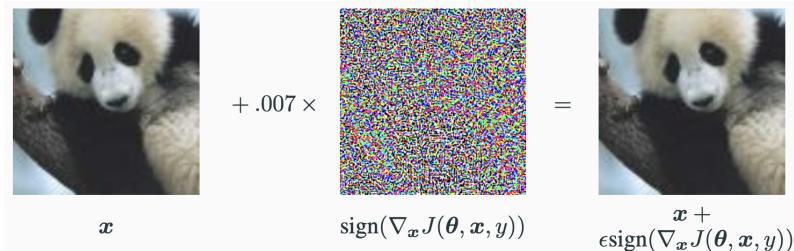


Figure 4.16: Linear Perturbation of Non-linear Models

Bagging and Dropout

Suppose we have the data set \mathcal{D} , sample $\mathcal{D}_1, \dots, \mathcal{D}_k$ of size N from \mathcal{D} with replacement. The bagging method is that we train k classifiers f_1, \dots, f_k on $\mathcal{D}_1, \dots, \mathcal{D}_k$ and predict using majority (or average if using regression).

However, this bagging approach is not practical for deep networks. A variant of bagging that works for neural networks is dropout.

The drop-out method is, during each training step, a fraction of the randomly-chosen hidden neurons are ignored and weights and biases of the others are updated in the gradient step. At each gradient update step, the neurons ignored are different. This can prevent so-called co-adaptation among different neurons. At test time, the entire network is used. All weights need to be scaled: halved if dropout rate was $1/2$.

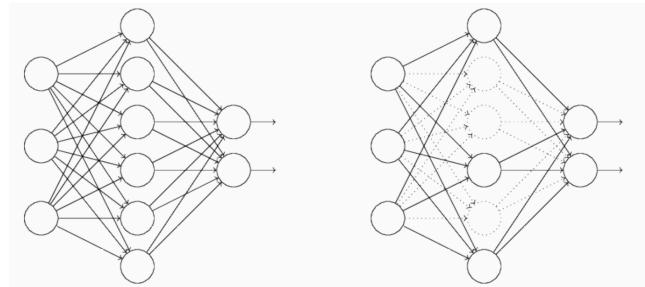


Figure 4.17: Dropout

Further ideas on avoiding overfitting include:

- Use parameter sharing (weight tying) in the model which will also reduce the computational cost.
- Exploit invariances to translation, rotation which is important to achieve high accuracy.
- Exploit locality in images, audio, text which is important to achieve high accuracy.

These are all achieved in convolutional neural networks, which will be discussed in the next section.

4.3 Convolutional Neural Networks

Convolutional neural networks, based on works of Fukushima, LeCun, Hinton (1980s), is one of the key reasons behind the recent popularity of machine learning. It makes significant advances in image, audio, signal processing and leads to highly accurate image recognition.

4.3.1 Convolution

Convolution operation on two real-valued functions. It expresses how the shape of one function is modified by the other. In neural networks, it is common to use cross-correlation instead.

A typical case is the convolution of an input 2-D image tensor I with 2-D filter (or kernel) tensor K . Tensor I maps coordinates to, eg, black/white or RGB colour code and tensor K is made up of weight parameters that are learned.

Tensor K convolving across tensor I gives the feature map $S(i, j)$, given by,

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n)K(m, n) \quad (4.49)$$

where $\sum_m \sum_n I(i + m, j + n)K(m, n)$ is the summation of the Hadamard, i.e. element-wise, products of the image and the kernel.

It is common to have a non-linear activation function applied to S .

Figure 4.18 shows an example of an input image $I \in \mathbb{R}^{3 \times 4}$ convolving with a kernel $K \in \mathbb{R}^{2 \times 2}$.

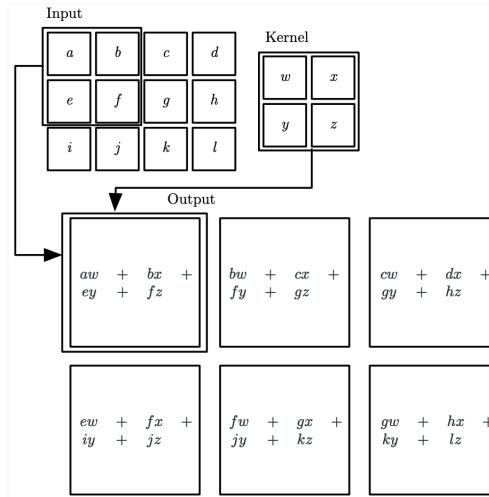


Figure 4.18: Convolution

We slide filter K over input image I one column at a time and move to the next row when finishing with one.

We compute $S(i, j)$ for $i \in \{0, 1, 2\}, j \in \{0, 1\}$, given by,

$$\begin{aligned} S(i, j) = & I(i, j)K(0, 0) + I(i + 1, j)K(1, 0) + \\ & I(i, j + 1)K(0, 1) + I(i + 1, j + 1)K(1, 1) \end{aligned} \quad (4.50)$$

In this case, we use a stride of 1, i.e. sliding in steps of 1 columns/rows.

The convolution of the image $I^{\ell-1} \in \mathbb{R}^{m_{\ell-1} \times n_{\ell-1}}$ with a kernel $K \in \mathbb{R}^{k_1 \times k_2}$ and a stride s will give an output image of $I^\ell \in \mathbb{R}^{(\frac{m_{\ell-1}-k_1}{s}+1) \times (\frac{n_{\ell-1}-k_2}{s}+1)}$.

In order to retain the original size of the input image, we will have to add columns and rows with entries of 0 so that the input image $I^{\ell-1} \in \mathbb{R}^{(m_{\ell-1}+(k_1-s)) \times (n_{\ell-1}+(k_2-s))}$, which is called padding.

Convolutional filters or kernels can detect patterns in images.

For example, to detect black-white vertical edges, we can use the filter,

$$\begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}$$

and bias $b = -1$ with ReLU activation.

The output for one application is:

$$\begin{aligned} \max(0, \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix} - 1) = \\ \max(0, -a + b - c + d - 1) \end{aligned} \quad (4.51)$$

This output is > 0 when the input matrix is

$$\begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}$$

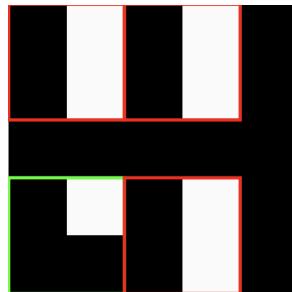


Figure 4.19: Black-White Vertical Edges Detector

To detect black L shapes, we can use the filter,

$$\begin{bmatrix} -1 & 1 \\ -1 & -1 \end{bmatrix}$$

and bias $b = 0$ with ReLU activation.

The output for one application is:

$$\begin{aligned} \max(0, \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} -1 & 1 \\ -1 & -1 \end{bmatrix}) = \\ \max(0, -a + b - c - d) \end{aligned} \quad (4.52)$$

This output is > 0 when the input matrix is

$$\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$$

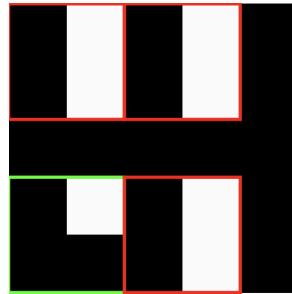


Figure 4.20: Black L Shapes Detector

Basic patterns detected by filters placed in the first layers of convnet are edges, textures, curves, objects, colors and shapes, i.e. corners, circles and squares. More sophisticated objects such as eyes and ears are detectable in later layers. In even deeper layers, things like full dogs, cats, lizards, and birds can be detected.

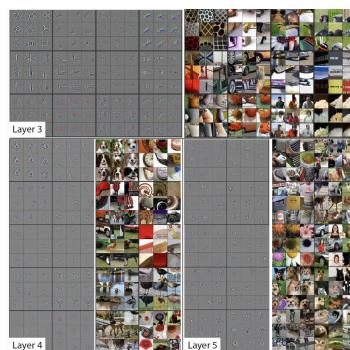


Figure 4.21: Patterns Detected

4.3.2 Max-Pooling

Max-pooling is a sample-based discretization process which generalise the model by combining several values into a single one. It decreases the chance of overfitting, reduces the number of parameters, and provides basic translation invariance and therefore more robust matching of features in the presence of small distortions.

Once we know a specific feature is in the input (its activation value is high), its exact location is less important relative to the other features.

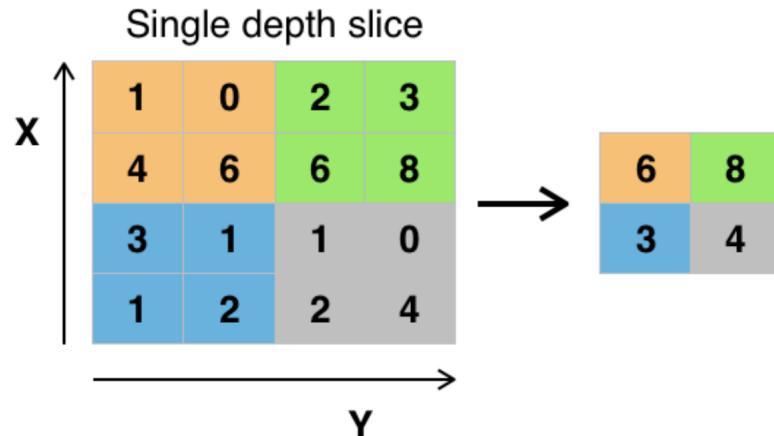


Figure 4.22: Max-Pooling

In addition to max-pooling, we also use min-pooling and average pooling which takes the minimum and average of the activation value each respectively. Average pooling is a softer version of max-pooling. For activations crossing several pooling neighbourhoods, average pooling gives a strong signal in the middle and soft at edges. It gives more information on where the feature edges are localised but might not extract good features if all not all inputs are needed.

4.3.3 Some Popular Convolutional Neural Networks

LeNet-5 (1998)

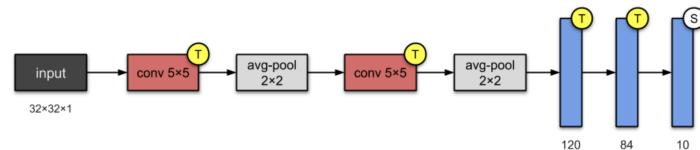


Figure 4.23: LeNet-5

- Two convolutional layers with 5×5 filters and tanh non-linear activation
- Two 2×2 average-pooling layers
- Three fully connected layers, 2 with tanh, one with softmax at end
- 60K parameters in total

AlexNet (2012)

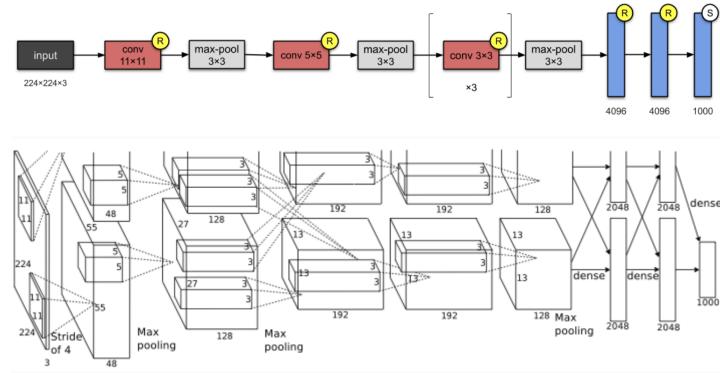


Figure 4.24: AlexNet

- Five convolutional layers with ReLU non-linear activation
- Two 3×3 max-pool layers
- Three fully connected layers, 2 with ReLU, one with softmax at end
- 60M parameters in total

VGG-16 (2014)

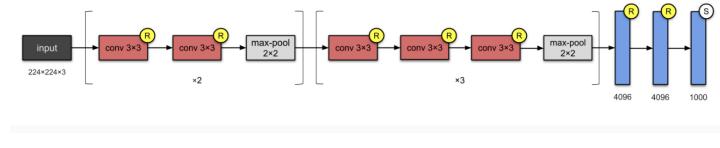


Figure 4.25: VGG-16

- Thirteen convolutional layers with 3×3 filters and ReLU non-linear activation
- Five 2×2 max-pool layers

- Three fully connected layers, 2 with ReLU, one with softmax at end
- 138M parameters in total

Inception-ResNet-V2 (2016)

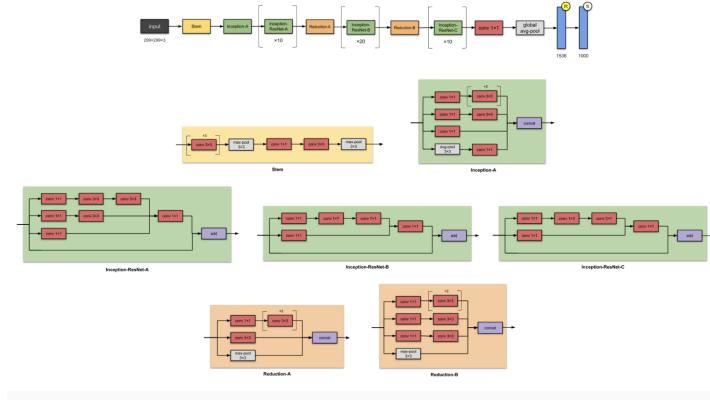


Figure 4.26: Inception-ResNet-V2

- Google-scale construction
- 56M parameters in total

4.3.4 Back-propagation in Convolutional Neural Networks

Consider convolutional layer between layer ℓ and $\ell + 1$. The output of the layer ℓ has the shape $m_\ell \times n_\ell \times F_\ell$, with each entry indexed as $a_{i,j,f}^\ell$, where $m_\ell \times n_\ell$ is the width and height of the output image each respectively and F_ℓ is the number of the filters.

For the $\ell + 1$ layer, we apply $F_{\ell+1}$ filters, each a tensor of shape $W_f \times H_f \times F_\ell$. The weights used in convolutional layers are shared as the filters convolve across the input layer. We assume no zero-padding and stride of 1 in each direction and index the entries in this convolutional layer, i.e. layer $\ell + 1$, as $z_{i',j',f'}^{\ell+1}$. Then we have,

$$z_{i',j',f'}^{\ell+1} = b^{\ell+1,f'} + \sum_{i=1}^{W_f} \sum_{j=1}^{H_f} \sum_{f=1}^{F_\ell} a_{i'+i-1,j'+j-1,f}^\ell w_{i,j,f}^{\ell+1,f'} \quad (4.53)$$

where $b^{\ell+1,f'}$ is the bias term of the f' -th filter at layer $\ell + 1$. Thus,

$$\frac{\partial z_{i',j',f'}^{\ell+1}}{\partial w_{i,j,f}^{\ell+1,f'}} = a_{i'+i-1,j'+j-1,f}^\ell \quad (4.54)$$

Then,

$$\frac{\partial \ell}{\partial w_{i,j,f}^{\ell+1,f'}} = \sum_{i',j'} \frac{\partial \ell}{\partial z_{i',j',f'}^{\ell+1}} a_{i'+i-1,j'+j-1,f}^{\ell} \quad (4.55)$$

Backpropagation also needs to compute $\frac{\partial \ell}{\partial z_{i,j,f}^{\ell}}$ using $\frac{\partial \ell}{\partial z_{i',j',f'}^{\ell+1}}$.

$$\frac{\partial z_{i',j',f'}^{\ell+1}}{a_{i,j,f}^{\ell}} = w_{i-i'+1,j-j'+1,f}^{\ell+1,f'} \quad (4.56)$$

Thus,

$$\frac{\partial \ell}{a_{i,j,f}^{\ell}} = \sum_{i',j',f'} i', j', f' \frac{\partial \ell}{\partial z_{i',j',f'}^{\ell+1}} w_{i-i'+1,j-j'+1,f}^{\ell+1,f'} \quad (4.57)$$

Then,

$$\frac{\partial \ell}{z_{i,j,f}^{\ell}} = f'(z_{i,j,f}^{\ell}) \sum_{i',j',f'} i', j', f' \frac{\partial \ell}{\partial z_{i',j',f'}^{\ell+1}} w_{i-i'+1,j-j'+1,f}^{\ell+1,f'} \quad (4.58)$$

For max-pooling layer,

$$s_{i',j'}^{\ell+1} = \max_{i,j \in \Omega(i',j')} a_{i,j}^{\ell} \quad (4.59)$$

Thus,

$$\frac{\partial s_{i',j'}^{\ell+1}}{\partial a_{i,j}^{\ell}} = \mathbb{1}((i,j) = \arg \max_{\tilde{i},\tilde{j} \in \Omega(i',j')} a_{\tilde{i},\tilde{j}}^{\ell}) \quad (4.60)$$

In a word, the partial derivative is 1 precisely for the coordinates of the max value in the set and 0 otherwise.

4.3.5 Number of Parameters in Convolutional Neural Networks

Consider the convolutional layer with c^ℓ input channels and $c^{\ell+1}$ output channels.

The number of filters is $c^\ell \times c^{\ell+1}$ and the number of bias is $c^{\ell+1}$.

Suppose we use filters, i.e. kernels, of size $k_1 \times k_2$. The number of parameters of this convolutional layer is given as ,

$$\# \text{parameters} = c^\ell \times c^{\ell+1} \times k_1 \times k_2 + c^{\ell+1} \quad (4.61)$$

5

Unsupervised Learning

5.1 Clustering

Often data can be grouped together into subsets that are coherent. However, this grouping may be subjective. It is hard to define a general framework.

There are two types of clustering algorithms: feature-based where points are represented as vectors in \mathbb{R}^D and (Dis)similarity-based where only pairwise (dis)similarities are known.

Two types of clustering methods: flat which partitions the data into k clusters and hierarchical which organises data as clusters, clusters of clusters, and so on.

5.1.1 k-Means Formulation of Clustering

The goal is to partition the data into subsets C_1, \dots, C_k where k is fixed in advance.

The quality of a partition defined by,

$$W(C) = \frac{1}{2} \sum_{j=1}^k \frac{1}{|c_j|} \sum_{i,i' \in C_j} d(\mathbf{x}_i, \mathbf{x}_{i'}) \quad (5.1)$$

If we use the Euclidean distance:

$$d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2^2$$

then the k-means objective is to minimise the sum of squares of distances to the mean within each cluster, i.e. to minimise jointly over partitions C_1, \dots, C_k and $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$, given as,

$$W(C) = \sum_{j=1}^k \sum_{i \in C_j} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|_2^2 \quad (5.2)$$

where $\boldsymbol{\mu}_j = \frac{1}{|c_j|} \sum_{i \in C_j} \mathbf{x}_i$.

This problem is NP-hard even for $k = 2$ for points in \mathbb{R}^D .

However, if we fix means $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$, finding a partition $(C_j)_{j=1}^k$ that minimises $W(C)$

$$C_j = \{i \mid \|\mathbf{x}_i - \boldsymbol{\mu}_j\|_2 = \min_{j'} \|\mathbf{x}_i - \boldsymbol{\mu}'_{j'}\|_2\}$$

is easy. Also, if we fix the clusters C_1, \dots, C_k , minimising W with respect to $(\boldsymbol{\mu}_j)_{j=1}^k$

$$\boldsymbol{\mu}_j = \frac{1}{|C_j|} \sum_{i \in C_j} \mathbf{x}_i$$

is easy.

Thus, alternative minimisation can be performed by iteratively running these assignment and update steps.

5.1.2 The k-Means Algorithm

The k-means algorithm is used to partition a given set of observations into a predefined amount of k clusters. The algorithm starts with a random set of k center-points

$$\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$$

During each update step, all observations \mathbf{x} are assigned to their nearest center-point (see equation 5.3). In the standard algorithm, only one assignment to one center is possible. If multiple centers have the same distance to the observation, a random one would be chosen.

$$S_i^{(t)} = \{\mathbf{x}_p : \|\mathbf{x}_p - \boldsymbol{\mu}_i^{(t)}\|^2 \leq \|\mathbf{x}_p - \boldsymbol{\mu}_j^{(t)}\|^2 \forall j, 1 \leq j \leq k\} \quad (5.3)$$

Afterwards, the center-points are repositioned by calculating the mean of the assigned observations to the respective center-points,

$$\boldsymbol{\mu}_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{\mathbf{x}_j \in S_i^{(t)}} \mathbf{x}_j \quad (5.4)$$

The update process reoccurs until all observations remain at the assigned center-points and therefore the center-points would not be updated anymore.

This algorithm is also called the Lloyd's algorithm.

Note that ties can be broken arbitrarily and choosing k random data points to be the initial k-means is a good idea.

Because the objective function W decreases every time a new partition is used, there are only finitely many partitions. Thus, the algorithm always converge.

Regardless of the initialization, the algorithm always terminate after at most k^N iterations where N is the number of input points and k is the number of clusters.

Convergence may be very slow in the worst-case, but typically fast on real-world instances. However, convergence is likely to be a local minimum. We can run the algorithm multiple times with random initialisation to tackle with this problem.

For the parameter k , larger k lowers the objective but might give poorer clustering, as is shown in Figure 5.1.

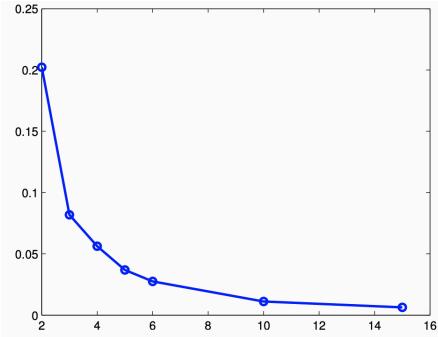


Figure 5.1: MSE on Test v.s. K

We can choose a suitable k by identifying a “kink” or “elbow” in the curve.

5.1.3 Beyond k -Means

k -medoids

In k -medoids algorithm, actual data points are assigned as cluster means rather than Euclidean averages. We can also use any distance between points beyond the squared Euclidean distance. In particular, we can use the Euclidean distance or any other ℓ_p norm

k -center

In k -center algorithm, the objective is to maximize over all dissimilarities between data and anchor. However, for k -means, it’s the sum of all the dissimilarities in each cluster. In both cases, the anchors or centres take the role of means.

5.1.4 Transforming Input Formats

Clustering algorithms may work on different data formats. The data can be given in Euclidean space or as matrix of (dis)similarities. In this section, we will discuss how to transform Euclidean distance to dissimilarity and vice versa.

The weighted dissimilarity between (real-valued) attributes is,

$$d(\mathbf{x}, \mathbf{x}') = f\left(\sum_{i=1}^D w_i d_i(x_i, x'_i)\right)$$

The simplest setting is $w_i = 1$, $d_i(x_i, x'_i) = (x_i - x'_i)^2$ and $f(z) = \sqrt{z}$, which corresponds to the Euclidean distance.

Weights allow us to emphasise features differently. If features are ordinal or categorical then define distance suitably as

$$d_i(x_i, x'_i) = 1$$

if $x_i = x'_i$ and 0 otherwise.

In practice, it may be easier to define (dis)similarity between objects than embed them in Euclidean space. For example, for DNA sequences, we can use Hamming distance as measure of dissimilarity and for text data, we can use cosine kernel as measure of similarity.

Since algorithms such as k-means require however points to be in Euclidean space, we need to transform (dis)similarity measures into Euclidean space

Multidimensional Scaling gives a way to find an embedding of the data in Euclidean space that (approximately) respects the original distance/similarity.

Assume an ideal case where we are given the dissimilarity matrix \mathbf{D} with pairwise Euclidean distances of points $\mathbf{x}_1, \dots, \mathbf{x}_N$

$$D_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$$

for $i, j \in [N]$.

Since distances are preserved under translation, rotation, reflection, etc. We cannot recover $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ exactly. However, we can aim to determine \mathbf{X} up to these transformations.

If D_{ij} is the distance between points \mathbf{x}_i and \mathbf{x}_j , then

$$\begin{aligned} D_{ij} &= \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \\ &= \mathbf{x}_i^T \mathbf{x}_i - 2\mathbf{x}_i^T \mathbf{x}_j + \mathbf{x}_j^T \mathbf{x}_j \\ &= M_{ii} - 2M_{ij} + M_{jj} \end{aligned} \tag{5.5}$$

where $\mathbf{M} = \mathbf{x}\mathbf{x}^T$ is the $N \times N$ matrix of dot products: $M_{ij} = \mathbf{x}_i^T \mathbf{x}_j$.

If

$$\sum_i \mathbf{x}_i = 0$$

\mathbf{M} can be uniquely recovered from \mathbf{D} . Otherwise, \mathbf{M} can be recovered from \mathbf{D} up to some translation.

Now we have \mathbf{M} . To obtain points $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_N$ from \mathbf{M} , we use Singular Value Decomposition of \mathbf{M} .

The Singular Value Decomposition of a matrix $\mathbf{X} \in \mathbb{R}^{D \times D}$ is

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^T$$

SVD is thin SVD if $\mathbf{U} \in \mathbb{R}^{N \times D}$ and $\mathbf{V} \in \mathbb{R}^{D \times D}$ are orthonormal matrices; $\Sigma \in \mathbb{D}^{D \times D}$ is diagonal with entries along the diagonal satisfying

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_D \geq 0$$

and $\mathbf{U}^T\mathbf{U} = \mathbf{V}^T\mathbf{V} = \mathbf{I}_D$.

SVD is full SVD $\mathbf{U} \in \mathbb{R}^{N \times N}$, $\Sigma \in \mathbb{R}^{N \times D}$ and $\mathbf{V} \in \mathbb{R}^{D \times D}$.

Eigendecomposition: if \mathbf{M} is squared and $\mathbf{U} = \mathbf{V}$

$$\mathbf{M} = \mathbf{U}\Sigma\mathbf{U}^T$$

If \mathbf{M} is square, symmetric and positive semi-definite, starting from \mathbf{M} , we can reconstruct $\hat{\mathbf{X}}$ using the eigendecomposition of \mathbf{M}

$$\mathbf{M} = \mathbf{U}\Sigma\mathbf{U}^T$$

Since Σ has non-negative diagonal entries, so we can take its square root.

By letting

$$\hat{\mathbf{X}} = \mathbf{U}\Sigma^{\frac{1}{2}}$$

, we obtain

$$\hat{\mathbf{X}}\hat{\mathbf{X}}^T = \mathbf{U}\Sigma^{\frac{1}{2}}(\mathbf{U}\Sigma^{\frac{1}{2}})^T$$

We thus find the points $\hat{\mathbf{x}}_i$ as rows of $\mathbf{U}\Sigma^{\frac{1}{2}}$.

In practice, we can use Mercer kernel to define the similarity and thus the matrix \mathbf{M} is positive semi-definite. Then the above derivation works.

In general, recovering the data points from dissimilarity matrix is to solve the optimisation problem with (non-convex) objective:

$$\arg \min_{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_N} \sum_{i \neq j} (\|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j\|_2^2 - D_{ij})^2$$

This objective function is called the stress function. It gives the degree to which it is not possible to find a Euclidean embedding for \mathbf{D} . There is no closed-form solution to this optimization problem. Local optimum is obtained via gradient descent.

5.1.5 Hierarchical Clustering

Hierarchical structured data exists all around us. For example, measurements of different species and individuals within species, top-level and low-level categories in news articles and country, canton, town level data are all hierarchical structured data.

There are two algorithmic strategies for clustering: agglomerative and divisive algorithms. Agglomerative algorithm refers to bottom-up method, where

clusters are formed by merging smaller clusters. The most popular agglomerative algorithm is called the Linkage algorithm. Divisive algorithm, on the other hand, refers to top-down method, where clusters are formed by splitting larger clusters.

We can visualise the clustering as a dendrogram or binary tree, as is shown in the Figure 5.2.

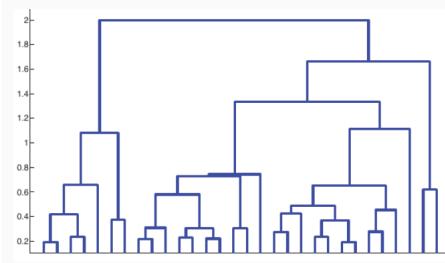


Figure 5.2: Visualizing the Clustering as a Dendrogram

Cutting the dendrogram at some level gives a partition of data.

To find hierarchical clusters we need to define dissimilarity at cluster level, not just at data points.

Suppose we have dissimilarity at data point level, e.g., $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|$.

Different ways to define dissimilarity at cluster level, say C and C' are given by

- Single Linkage

$$D(C, C') = \min_{x \in C, x' \in C'} d(x, x')$$

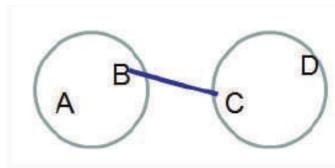


Figure 5.3: Single Linkage

- Complete Linkage

$$D(C, C') = \max_{x \in C, x' \in C'} d(x, x')$$

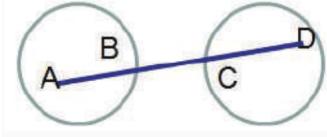


Figure 5.4: Complete Linkage

- Average Linkage

$$D(C, C') = \frac{1}{|C| \cdot |C'|} \sum_{x \in C, x' \in C'} d(x, x')$$

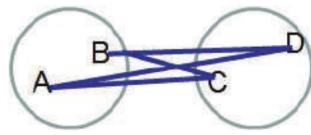


Figure 5.5: Average Linkage

Linkage-based Clustering Algorithm

The algorithm starts with randomly initialized singleton clusters

$$C_i = \{i\}$$

and clusters available for merging

$$S = \{1, \dots, N\}$$

During each update step, Pick 2 most similar clusters

$$(j, k) = \arg \min_{j, k \in S} D(j, k)$$

For some new index l , let

$$C_l = C_j \cup C_k$$

Break if

$$C_l = \{1, \dots, N\}$$

Otherwise, update

$$S := (S \setminus \{j, k\}) \cup \{l\}$$

Then update $D(i, l)$ for all $i \in S$ using desired linkage property.

5.1.6 Spectral Clustering

K-means will typically form clusters that are spherical, elliptical and convex. To find non-convex clusters, spectral clustering is a (related) alternative that often works better.

Similarity Graph Construction

To perform spectral clustering, we need to construct k -nearest neighbour graph from data by assigning one node for every point in dataset. Then (i, j) is an edge if either i is among the k nearest neighbours of j or vice versa.

The weight of edge (i, j) , if it exists, is given by similarity measure $s_{i,j}$

$$s_{i,j} = e^{-\frac{\|x_i - x_j\|^2}{\sigma}}$$

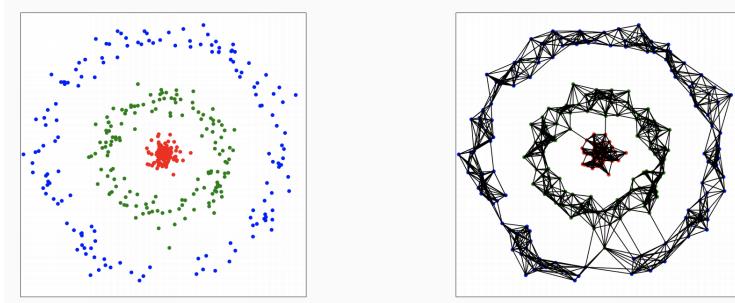


Figure 5.6: Spectral Clustering: Graph Construction

Partitioning a Graph into Two Clusters

Min-Cut Partition graph into two sets A and B such that weight of edges connecting vertices in A to vertices in B is minimum.

$$\text{cut}(A, B) := \sum_{i \in A, j \in B} s_{ij}$$

Mincut can give bad cuts (only one node on one side of the cut).

Normalized-Cut Partition graph into two sets A and B such that weight of edges connecting vertices in A to vertices in B is minimum and the size of A and B are very similar.

$$\text{Ncut}(A, B) := \text{cut}(A, B) \left(\frac{1}{\text{vol}(A)} + \frac{1}{\text{vol}(B)} \right)$$

where

$$\frac{1}{\text{vol}(A)}$$

is the sum of the edges connecting all pairs of vertices in A.

Normalized cut is typically NP-hard to compute. Relaxations of these problems give eigenvectors of Laplacian.

Spectral Clustering Algorithm

Suppose we have the weighted graph with weighted adjacency matrix \mathbf{W} as input and number k of clusters to construct.

We start with construct Laplacian

$$\mathbf{L} = \mathbf{D} - \mathbf{W}$$

where \mathbf{W} is the weighted adjacency matrix and \mathbf{D} is (diagonal) degree matrix with entries

$$D_{ii} = \sum_j w_{ij}$$

Then we find the first k (hyper-parameter) eigenvectors of the matrix \mathbf{L}

$$\mathbf{v}_1 = \mathbf{1}, \mathbf{v}_2, \dots, \mathbf{v}_k$$

Next we build the feature matrix

$$\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k] \in \mathbb{R}^{N \times k}$$

with the eigenvalues as columns.

We can interpret the rows of \mathbf{V} as new data points $\mathbf{z}_i \in \mathbb{R}^k$ and cluster the points \mathbf{z}_i with the k -means algorithm in \mathbb{R}^k .

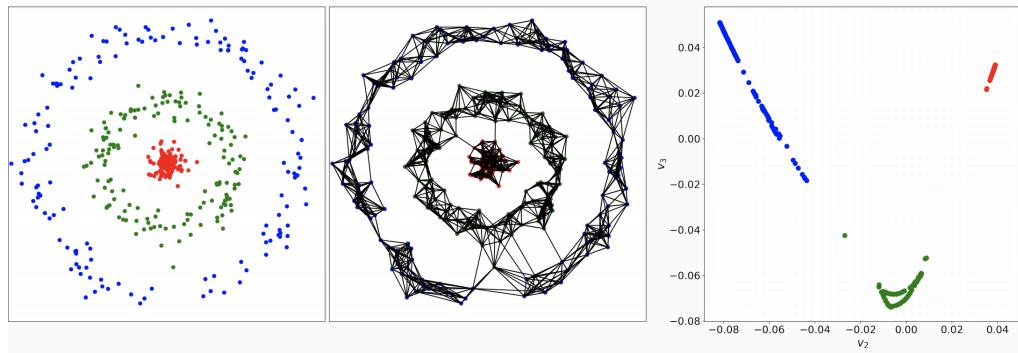


Figure 5.7: Spectral Clustering: $k = 2$

5.2 Principal Component Analysis

Real-world data may have correlated variables, i.e. when one variable changes, the other changes as well. Data synthesised by different teams may have redundancy.

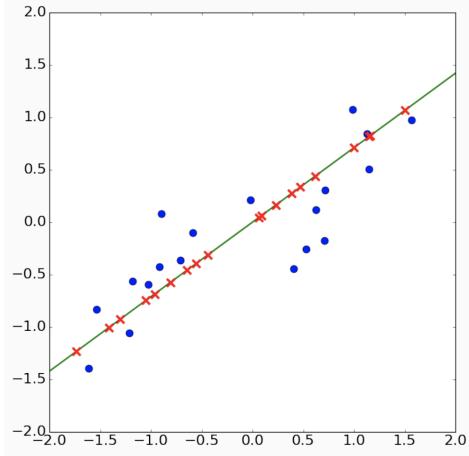


Figure 5.8: Dimensionality Reduction: Principal Component Analysis

Dimensionality reduction projects the data into a lower dimensional subspace while still capturing the essence of the original data. It is a preprocessing step before applying other learning algorithms.

As is shown in Figure 5.8, the red points are orthogonal projections of the blue points onto a vector representing the 1st principal component. The 2D points are projected onto 1D.

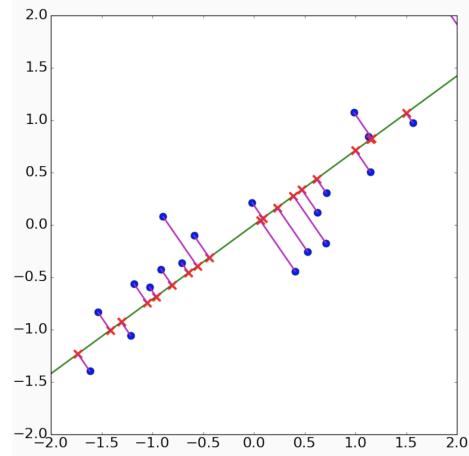


Figure 5.9: Principal Component Analysis: Reconstruction Errors

Principal Component Analysis (PCA) used for dimensionality reduction identifies a small number of directions which explain most variation in data.

The magenta lines shown in Figure 5.9 are the distances of the blue points from the first principal component. They are also called reconstruction errors.

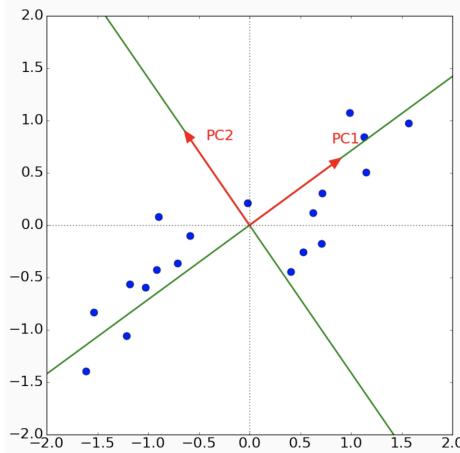


Figure 5.10: Principal Component Analysis

The second principal component is computed as the first principal component after the correlation with first principal component has been subtracted from the points. The second principal component is orthogonal to the first principal component.

5.2.1 PCA: Maximum Variance View

PCA keeps subspace with maximum variance. If we align the points with the principal components and move as much of the variance as possible using an orthogonal transformation into the first few dimensions, the values in the remaining dimensions tend to be small and may be dropped with minimal loss of information, i.e. PCA is the optimal orthogonal transformation that keeps the subspace with the largest "variance".

PCA is a linear dimensionality reduction technique searching for the directions of maximum variance in the data $(\mathbf{x}_1, \dots, \mathbf{x}_N)$. It finds a set of orthonormal vectors $(\mathbf{v}_1, \dots, \mathbf{v}_k)$ such that the first principal component (PC) \mathbf{v}_1 is the direction of largest variance, the second PC \mathbf{v}_2 is the direction of largest variance orthogonal to \mathbf{v}_1 and the i^{th} PC \mathbf{v}_i is the direction of largest variance orthogonal to $\mathbf{v}_1, \dots, \mathbf{v}_{i-1}$.

Suppose we are independent and identically distributed data

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$$

PCA can be rephrased as the following optimization problem:

$$\begin{aligned} & \text{maximize} && \|\mathbf{X}\mathbf{v}_1\|^2 \\ & \text{subject to} && \|\mathbf{v}_1\| = 1 \end{aligned}$$

where $\mathbf{v}_1 \in \mathbb{R}^D$.

The solution is

$$\begin{aligned} \mathbf{v}_1 &= \arg \max_{\mathbf{v}_1: \|\mathbf{v}_1\|=1} \mathbf{v}_1^T \mathbf{X}^T \mathbf{X} \mathbf{v}_1 \\ &= \arg \max_{\mathbf{v}_1} \frac{\mathbf{v}_1^T \mathbf{X}^T \mathbf{X} \mathbf{v}_1}{\mathbf{v}_1^T \mathbf{v}_1} \end{aligned} \quad (5.6)$$

Recall the Rayleigh quotient, for positive semi-definite \mathbf{M} , the max value of

$$R(\mathbf{M}, \mathbf{v}) = \frac{\mathbf{v}^T \mathbf{M} \mathbf{v}}{\mathbf{v}^T \mathbf{v}}$$

is the largest eigenvalue of \mathbf{M} which is attained for \mathbf{v} being the corresponding eigenvector.

Since $\mathbf{X}^T \mathbf{X}$ is positive semi-definite, the solution to our problem is the corresponding eigenvector of the largest eigenvalue of $\mathbf{X}^T \mathbf{X}$.

Let

$$\mathbf{z} = \mathbf{X}\mathbf{v}_1 = [\mathbf{x}_1^T \mathbf{v}_1, \dots, \mathbf{x}_N^T \mathbf{v}_1]$$

Thus,

$$z_i = \mathbf{x}_i^T \mathbf{v}_1$$

$\|\mathbf{X}\mathbf{v}_1\|^2$ is the variance of the projections of data onto \mathbf{v}_1 (mod constant factor $\frac{1}{N}$)

$$\|\mathbf{X}\mathbf{v}_1\|^2 = \sum_{i=1}^N z_i^2 = \mathbf{z}^T \mathbf{z} = \mathbf{v}_1^T \mathbf{X}^T \mathbf{X} \mathbf{v}_1$$

This assumes the data is centered

$$\sum_i \mathbf{x}_i = \mathbf{0}$$

We can find $\mathbf{v}_2, \dots, \mathbf{v}_k$ that are all successively orthogonal to previous directions and maximise (as yet unexplained) variance

$$\hat{\mathbf{X}}_j = \mathbf{X} - \sum_{s=1}^{j-1} \mathbf{X} \mathbf{v}_s \mathbf{v}_s^T \quad j \in [k] \quad (5.7)$$

5.2.2 PCA: Best Reconstruction View

Suppose we are independent and identically distributed data

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$$

PCA in case of one PC can be rephrased as the following optimization problem:

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 \\ & \text{subject to} \quad \|\mathbf{v}_1\| = 1 \end{aligned}$$

where \mathbf{v}_1 is the direction of projection and $\tilde{\mathbf{x}}_i$ is the mapping of the point \mathbf{x}_i given by

$$\tilde{\mathbf{x}}_i = \frac{\mathbf{v}_1 \mathbf{x}_i}{\|\mathbf{v}_1\|} \mathbf{v}_1 = (\mathbf{v}_1 \mathbf{x}_i) \mathbf{v}_1 \quad (5.8)$$

Let's now look at the equivalence between the objectives of the maximum variance view and best reconstruction view in the case of one PC.

$$\begin{aligned} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 &= \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 \\ &= \sum_{i=1}^N \|\mathbf{x}_i\|^2 - 2(\mathbf{x}_i \tilde{\mathbf{x}}_i) + \|\tilde{\mathbf{x}}_i\|^2 \\ &= \sum_{i=1}^N \|\mathbf{x}_i\|^2 - 2(\mathbf{x}_i \tilde{\mathbf{x}}_i) + (\mathbf{v}_1 \mathbf{x}_i)^2 \|\mathbf{v}_1\|^2 \quad (5.9) \\ &= \sum_{i=1}^N \|\mathbf{x}_i\|^2 - 2(\mathbf{x}_i \tilde{\mathbf{x}}_i) + (\mathbf{v}_1 \mathbf{x}_i)^2 \times 1 \\ &= \sum_{i=1}^N \|\mathbf{x}_i\|^2 - \sum_{i=1}^N (\mathbf{v}_1 \mathbf{x}_i)^2 \end{aligned}$$

Thus, our objective of the best reconstruction view in the case of one PC becomes,

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^N \|\mathbf{x}_i\|^2 - \sum_{i=1}^N (\mathbf{v}_1 \mathbf{x}_i)^2 \\ & \text{subject to} \quad \|\mathbf{v}_1\| = 1 \end{aligned}$$

Since $\sum_{i=1}^N \|\mathbf{x}_i\|^2$ is a constant given by the inputs (thus is not to be optimized), this objective is equivalent to:

$$\begin{aligned} & \text{minimize} \quad - \sum_{i=1}^N (\mathbf{v}_1 \mathbf{x}_i)^2 \\ & \text{subject to} \quad \|\mathbf{v}_1\| = 1 \end{aligned}$$

Equivalently,

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^N (\mathbf{v}_1 \mathbf{x}_i)^2 \\ & \text{subject to} && \|\mathbf{v}_1\| = 1 \end{aligned}$$

So the same \mathbf{v}_1 satisfies the two objectives of the maximum variance view and best reconstruction view in the case of one PC.

Finding Principal Components using Singular Value Decomposition

Suppose $\mathbf{V}_k \in \mathbb{R}^{D \times k}$ is such that columns of \mathbf{V}_k are orthogonal, then projection of data \mathbf{X} onto subspace is defined by

$$\mathbf{Z}_k = \mathbf{X} \mathbf{V}_k \quad (5.10)$$

PCA in case of k PCs can be rephrased as the following optimization problem:

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{V}_k \mathbf{V}_k^T \mathbf{x}_i\|^2 \\ & \text{subject to} && \|\mathbf{v}_i\| = 1 \quad i \in [k] \end{aligned}$$

Recall the SVD of \mathbf{X} :

$$\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T$$

where

- $\Sigma \in \mathbb{R}^{N \times D}$ is diagonal with singular values.

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_D \geq 0$$

- $\mathbf{V} \in \mathbb{R}^{D \times D}$ is orthonormal matrix with right singular vectors.
- $\mathbf{U} \in \mathbb{R}^{N \times N}$ is orthonormal matrix with left singular vectors.

Since $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T$,

$$\mathbf{X}^T \mathbf{X} = (\mathbf{U} \Sigma \mathbf{V}^T)^T (\mathbf{U} \Sigma \mathbf{V}^T) = \mathbf{V} \Sigma^T \mathbf{U}^T \mathbf{U} \Sigma \mathbf{V}^T = \mathbf{V} \Sigma^T \Sigma \mathbf{V}^T \quad (5.11)$$

where

$$\Sigma^T \Sigma = \Sigma^2 \in \mathbb{R}^{D \times D}$$

The eigenvectors of $\mathbf{X}^T \mathbf{X}$ are the right singular vectors \mathbf{v}_i of \mathbf{X} and the eigenvalues of $\mathbf{X}^T \mathbf{X}$ are the squares of the singular values σ_i of \mathbf{X} .

Similarly, the eigenvectors \mathbf{u}_j of $\mathbf{X}^T \mathbf{X}$ with eigenvalue σ^2 are the left singular vectors of \mathbf{X} (this can be used in case $D > N$).

Now, the PCA transformation is given as

$$\mathbf{Z} = \mathbf{X}\mathbf{V} = \mathbf{U}\Sigma\mathbf{V}^T\mathbf{V} = \mathbf{U}\Sigma \quad (5.12)$$

We can truncate \mathbf{Z} . For $\mathbf{Z}_k \in \mathbb{R}^{N \times k}$, we only need the first k largest singular values and their singular vectors

$$\mathbf{Z}_k = \mathbf{U}_k\Sigma_k = \mathbf{X}\mathbf{V}_k$$

Then the reconstruction is given by

$$\tilde{\mathbf{X}} = \mathbf{Z}_k\mathbf{V}_k^T = \mathbf{U}_k\Sigma_k\mathbf{V}_k^T \quad (5.13)$$

where $\|\cdot\|_F$ is the Frobenius norm

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{i,j}|^2}$$

The reconstruction error

$$\begin{aligned} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{V}_k\mathbf{V}_k^T\mathbf{x}_i\|^2 &= \|\mathbf{X} - \tilde{\mathbf{X}}\|_F \\ &= \|\mathbf{U}\Sigma\mathbf{V}^T - \mathbf{U}_k\Sigma_k\mathbf{V}_k^T\|_F \\ &= \left\| \sum_{j=1}^D \sigma_j \mathbf{u}_j \mathbf{v}_j^T - \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j^T \right\|_F \\ &= \sum_{j=k+1}^D \sigma_j^2 \end{aligned} \quad (5.14)$$

This corresponds to the Eckart-Young theorem: \mathbf{Z}_k is the nearest possible matrix of rank k to \mathbf{X} , i.e. the difference between the two has the smallest Frobenius norm.

Algorithm for Finding PCs (when $N > D$)

- Variant 1: Construct $\mathbf{X}^T\mathbf{X}$ in $O(D^2N)$ and its eigenvectors in $O(D^3)$.
- Variant 2: Iterative methods to get top k singular (right) vectors directly:
 - Initiate \mathbf{v}_0 to be random unit norm vector
 - Iterative Update:

$$\mathbf{v}^{t+1} = \mathbf{X}^T\mathbf{X}\mathbf{v}^t$$

$$\mathbf{v}^{t+1} = \frac{\mathbf{v}^{t+1}}{\|\mathbf{v}^{t+1}\|}$$

- Update step takes $O(ND)$ time (compute $\mathbf{X}\mathbf{v}^T$ first, then $\mathbf{X}^T(\mathbf{X}\mathbf{v}^T)$).
- This gives the singular vector corresponding to the largest singular value.
- Subsequent singular vectors obtained by choosing \mathbf{v}_0 orthogonal to previously identified singular vectors (this needs to be done at each iteration to avoid numerical errors creeping in).

Algorithm for Finding PCs (when $D \gg N$)

Constructing the matrix $\mathbf{X}\mathbf{X}^T$ takes time $O(N^2D)$. Eigenvectors of $\mathbf{X}\mathbf{X}^T$ can be computed in time $O(N^3)$.

The eigenvectors give the ‘left’ singular vectors, \mathbf{u}_i of \mathbf{X} . To obtain \mathbf{v}_i , we use the fact that

$$\mathbf{v}_i = \sigma^{-1} \mathbf{X}^T \mathbf{u}_i$$

Iterative method can be used directly as in the case when $N > D$.

5.2.3 How Many Principal Components to Pick

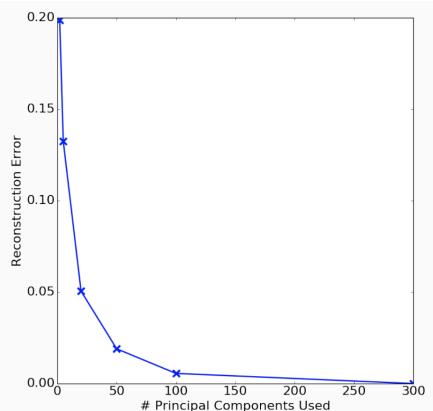


Figure 5.11: Principal Component Used v.s. Reconstruction Errors

We can choose a suitable number of PCs by identifying an ‘elbow’ in the curve of reconstruction error v.s. the number of PCs.

If SVD is computed, fix relative error threshold $0 \leq t \leq 1$ and choose k such that

$$\frac{\|\mathbf{X} - \mathbf{X}_k\|_F^2}{\|\mathbf{X}\|_F^2} \leq t$$

where $\|\cdot\|_F$ is the Frobenius norm.