

Data Science Algorithms

Peter Meleney

"All models are wrong, but some are useful."

-George E.P. Box

Data Science Algorithms This is a repository for what I learn about algorithms I use for data science. Each section should contain a conceptual overview of the algorithm and the basic math behind its implementation. I also intend to include examples in R and Python [1].

Created: March 4, 2016

Last Updated: April 7, 2016

Contents

I	Unsupervised Algorithms I: Dimensionality Reduction	6
1	Principle Component Analysis	7
1.1	Mathematics	7
II	Unsupervised Algorithms II: Clustering and Networks	8
2	K-Means Clustering	9
2.1	K-Means Algorithm	9
2.2	Requirements	10
2.3	Choosing K	10
2.3.1	The "Elbow" Method	10
2.3.2	Silhouette Method	10
2.4	Assessing Goodness of Fit	11
2.4.1	Error Rate	11
2.4.2	The Bayes Classifier	11
2.5	Example Applications ***ToDo***	12
2.5.1	Applications in Python ***ToDo***	12
2.5.2	Applications in R ***ToDo***	12
3	K-Nearest Neighbors	13
3.1	KNN Algorithm	13
3.1.1	Distance Weighted KNN	13
3.1.2	Instance Weighted KNN ***ToDo***	14
3.1.3	Attribute Weighted KNN ***ToDo***	14
3.2	Requirements	14
3.3	Choosing K	14
3.4	Applications of KNN ***ToDo***	14
3.4.1	Applications in Python ***ToDo***	14
3.4.2	Applications in R ***ToDo***	14

<i>CONTENTS</i>	3
4 Outlier Detection	15
4.1 ***To Do***	15
III Supervised Algorithms I: Regression	16
5 Linear Regression	17
5.1 Introduction	17
5.2 Assumptions	18
5.3 Simple Linear Regression	18
5.3.1 Least Squares Criterion	19
5.3.2 Derivation of β_0	19
5.3.3 Derivation of β_1	20
5.4 Multiple Linear Regression	21
5.5 Pearson Correlation Coefficient	22
5.6 Goodness of Fit	23
5.6.1 Adjusted R^2 for Multiple Linear Regression	23
5.7 Outliers	24
5.8 Extrapolation	24
5.9 Errors of Estimated Coefficients	24
5.9.1 Prediction Bands	24
5.9.2 Confidence Intervals	24
5.10 Checking Assumptions Using Graphs	24
5.11 Model Selection Strategies	24
5.11.1 The R^2_{adj} Approach	24
5.11.2 The p-Value Approach	24
5.12 Quick Methods and Shortcuts	24
IV Supervised Algorithms I: Classification	25
6 Logistic Regression	26
6.1 Introduction	26
6.2 Mathematics	27
V Supervised Algorithms III: Algorithms for Classification and Regression	29
7 Naïve Bayes	30
7.1 Probability	30
7.1.1 Tree Diagrams	30

<i>CONTENTS</i>	4
7.2 Bayes' Thorem	32
7.3 Naïve Bayes Applied	32
7.3.1 High-Level Class Probabilities	33
7.3.2 Rare Words	33
8 Support Vector Machines	34
8.1 ***To Do***	34
9 Random Forest	35
9.1 ***To Do***	35
10 Ensemble Methods	36
10.1 ***To Do***	36
VI Time Series	37
11 Holt-Winters Forecasting	38
11.1 ***To Do***	38
VII Advanced Algorithms	39
12 Neural Networks	40
12.1 ***To Do***	40
VIII Supporting Information	41
13 Distance Metrics	42
13.1 Euclidean Distance	42
13.2 Manhattan Distance (Hamming Distance)	42
13.3 Cosine Distance	43
13.4 Jaccard Similarity	43
14 Application Program Interface	44
14.1 Introduction	44
15 Set Theory, Relational Algebra, and Structured Query Language	45
15.1 Introduction	45
15.2 Set Theory	45
15.2.1 Basic Concepts and Notation	46
15.3 Structured Query Language	46

15.3.1	Best Practices	47
15.3.2	Entity Relationship Diagrams	47
15.3.3	Relational Schema	47
15.3.4	Database Metadata	47
15.3.5	Basic Queries	47
15.3.6	Aggregation Functions ***ToDo***	49
15.3.7	Joins ***ToDo***	49
15.3.8	Subqueries ***ToDo***	49
15.3.9	Logical Expressions ***ToDo***	49

Part I

Unsupervised Algorithms I: Dimensionality Reduction

Chapter 1

Principle Component Analysis

Principle component analysis (PCA) is a popular approach for deriving a lower-dimensional set of features from a high-dimensional set. The PCA algorithm identifies and preserves p directions of highest variance in the data and destroys the $m - p$ directions of lowest variability.

This technique has two primary applications. Firstly, it can reduce the complexity of data input to another algorithm, especially if the full feature set is too large to be processed in a timely manner. Secondly it can be used as an unsupervised technique in its own right. The similarity of direction in the new PCA space may be treated as a more or less quantitative measure of feature similarity. Additionally, and of particular interest is the case of setting $p = 2$, making optimized binary diagrams possible.

1.1 Mathematics

Given a data set X , with n observations and m features, and user-defined p , the desired number of principle components, PCA computes ϕ

$$Z^{[n \times p]} = \phi^{[p \times m]} X^{[n \times m]} \quad (1.1)$$

Part II

Unsupervised Algorithms II: Clustering and Networks

Chapter 2

K-Means Clustering

K-Means Clustering (KMC) is a greedy learner that partitions a data set into k clusters. The algorithm minimizes the Mean Square Distance between all the data points and the centroid of the group to which they belong. This process is not robust, meaning that the algorithm can fall into a local minimum and never reach a global minimum.

Modifying K-Means Clustering to K-Medians Clustering can be advantageous in some circumstances. What does it mean to have bought 0.723 of an item? Sometimes this is not a meaningful value for a cluster center.

2.1 K-Means Algorithm

1. Scale features to standardize distance.
2. Initialize k centers
3. Assign each point to its nearest center
4. Calculate the centroid of each group of points
5. Move the centers to their respective centroids
6. Repeat steps 3-5 until the centers no longer move with each iteration, this is your final set of groups.

2.2 Requirements

All features should be scaled before K-Means Clustering is applied.

2.3 Choosing K

2.3.1 The "Elbow" Method

This method involves performing KMC for a range of K, usually up to \sqrt{n} or $\ln(n)$, then graphing the mean square distance from the points to their respective centers. Usually somewhere in the graph a distinct change in slope will occur, from strongly negative to weakly negative. This elbow, or very near it will probably represent an optimal choice of K [2].

2.3.2 Silhouette Method

If clusters are properly assigned, then data points within a cluster should be much closer to each other than they are to data points in other clusters. The silhouette can be calculated for each point. In plain English: it is the average distance to points in the same cluster minus the average distance to points in the nearest (not same) cluster, divided by the maximum of the two. This value will always be [1,-1]. If the value is positive, the point is nearer points in its own group, if negative, it is nearer points in the other group [1].

A distance function most appropriate to the situation may be chosen to maximize the validity of the model.

$$\frac{\text{Avg}[\text{dist}(x, x_j)] - \text{Avg}[\text{dist}(x, x_i)]}{\max\{\text{Avg}[\text{dist}(x, x_j)], \text{Avg}[\text{dist}(x, x_i)]\}}$$

$$\frac{\frac{1}{n_{k_N}} \sum_{j=0}^{n_{k_N}} \text{dist}(x, x_j) - \frac{1}{n_k} \sum_{i=0}^{n_k} \text{dist}(x, x_i)}{\max(\frac{1}{n_{k_N}} \sum_{j=0}^{n_{k_N}} \text{dist}(x, x_j) - \frac{1}{n_k} \sum_{i=0}^{n_k} \text{dist}(x, x_i))} \quad (2.1)$$

i - points in cluster k

j - points in the cluster nearest k

n_k - number of data points in cluster k

n_{k_N} : number of data points in the cluster nearest k

k should then be selected that minimizes the average silhouette of all the data points in the set.

2.4 Assessing Goodness of Fit

2.4.1 Error Rate

The most common approach for quantifying the accuracy of a K-Means model is the error rate:

$$\frac{1}{n} \sum_{i=0}^n I(y_i \neq \hat{y}_i) \quad (2.2)$$

Where $I(y_i \neq \hat{y}_i)$ is an *indicator variable* which equals 1 if $y_i \neq \hat{y}_i$ and 0 if $y_i = \hat{y}_i$. Equation 2.2 is referred to as the **training error** when it is computed on the training data set, and the **test error** when computed on a test set.

2.4.2 The Bayes Classifier

It is possible to show that the equation 2.2 is minimized when each observation is assigned to the most likely class.

$$\text{Max}(P(y = j|X = x_o)) \quad (2.3)$$

This very simple classification rule is called a **Bayes Classifier**. In a two-class problem it can be restated:

$$P(y = j|X = x_o) > 0.5 \quad (2.4)$$

The boundary between two classes at which the probability of a hypothetical point is equal to 50% is called the **Bayes decision boundary**, and is a theoretical optimal solution to binary classification problems.

A data scientist need not choose the Bayes classifier as the decision boundary. Any value (0,1) may be chosen, the accuracy will be lower, but it may be desirable to exchange accuracy for greater precision or recall.

2.5 Example Applications ***ToDo***

2.5.1 Applications in Python ***ToDo***

This is in pcr font, and this text is ForestGreen.

2.5.2 Applications in R ***ToDo***

Chapter 3

K-Nearest Neighbors

3.1 KNN Algorithm

The K-Nearest Neighbors (KNN) algorithm is usually used for classification, but can also be used for regression. Given a positive integer K , and a test observation x_0 , the KNN classifier identifies the K nearest labeled observations to x_0 , denoted η_0 , and computes the probability that x_0 belongs to class j based on the distribution of those labeled points η_0 . KNN can apply weights based on distance, or not. In its simplest, unweighted form, the algorithm identifies η_0 and assigns x_0 to whichever class has the most points represented in η_0 .

$$P(Y = j|X = x_0) = \frac{1}{K} \sum_{i \in \eta_0} I(y_i = j) \quad (3.1)$$

3.1.1 Distance Weighted KNN

KNN can be written so that the distance from x_0 to each point in η_0 is considered by the algorithm. The distance metric may be any of those from chapter 13.

$$P(Y = j|X = x_0) = \frac{1}{K} \sum_{i \in \eta_0} \frac{I(y_i = j)}{Dist(x_0, y_i)^2} \quad (3.2)$$

3.1.2 Instance Weighted KNN ***ToDo***

3.1.3 Attribute Weighted KNN ***ToDo***

3.2 Requirements

All features should be scaled prior to initiating the KNN algorithm.

3.3 Choosing K

The choice of K has a drastic effect on the KNN algorithm. Choosing $K < K_{opt}$ will result in underfitting the decision boundary, while choosing $K > K_{opt}$ will result in overfitting. We expect the training error to monotonically decrease as $K \rightarrow \infty$, however the test error will not monotonically decrease. Instead the test error will decrease as $K = 1 \rightarrow K_{opt}$ and then increase again as $K_{opt} \rightarrow \infty$.

Choosing K_{opt} can be achieved by minimizing the error rate of the cross validation data subset, or through k-fold cross validation ***To Do: Add ref to CV Chapter***. Be careful not to optimize to the test set, as that set will no longer give you a proper analysis of model accuracy.

3.4 Applications of KNN ***ToDo***

3.4.1 Applications in Python ***ToDo***

3.4.2 Applications in R ***ToDo***

Chapter 4

Outlier Detection

4.1 ***To Do***

Part III

Supervised Algorithms I: Regression

Chapter 5

Linear Regression

5.1 Introduction

Linear regression is the most common supervised learning technique because it is straightforward to apply and relatively easy to interpret. Linear regression has several drawbacks however. Many systems do not follow a linear model, the most common violations of the assumptions of linear regression occur when features interact or when error terms are auto-correlated. Unlike several more advanced algorithmic approaches, linear regression is particularly sensitive to violations of its assumptions.

Linear regression is usually solved to minimize the **sum of squared errors** (SSE), confusingly also called the **residual sum of squares** (RSS). When applying this **cost function**, the solution to simple linear regression is as follows:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad (5.1)$$

$$\hat{\beta}_1 = \frac{\sum_{i=0}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=0}^n x_i - \bar{x}} \quad (5.2)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (5.3)$$

The solution to multiple linear regression, when \mathbf{y} is the solution to is somewhat more involved, and is better treated as a linear algebra problem. The solution is:

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y} \quad (5.4)$$

Derivations of these solutions are below.

5.2 Assumptions

Assumptions 1-3 are critical when applying linear regression, the others are required if linear regression is to be BLUE (best linear unbiased estimator), but their violation may not invalidate the model entirely.

1. The objective can be estimated by a linear combination of the independent features.
2. Features do not interact or are additive only.
3. Errors are independent (no auto-correlation).
4. Outliers and influential points are removed.
5. Errors are homoskedastic.
6. Errors are normally distributed.

5.3 Simple Linear Regression

Linear regression assumes that the objective is a linear combination of the input features. Simple linear regression, when the objective y is a function of only one independent feature x , has the form:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad (5.5)$$

Where ϵ_i represents the error between the linear relationship and the actual value of y_i . The linear regression algorithm approximates equation 5.5 by enforcing the relationship below:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad (5.6)$$

Where \hat{y}_i is the estimated output based on the linear relationship. The actual output y_i must include an error term e_i because the relationship is not perfect.

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + e_i \quad (5.7)$$

Subtracting equation 5.7 from equation 5.6 gives us the relationship:

$$e_i = \hat{y}_i - y_i \quad (5.8)$$

This remainder, e_i is commonly used in the **cost function**. The most common cost function is the sum of the square of the cost function, called the **least squares criterion**.

5.3.1 Least Squares Criterion

The least squares criterion minimizes the **residual sum of squares** (RSS):

$$RSS = \sum_{i=0}^n e_i^2 \quad (5.9)$$

$$RSS = \sum_{i=0}^n (y_i - \hat{y}_i)^2 \quad (5.10)$$

$$RSS = \sum_{i=0}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad (5.11)$$

We can minimize the RSS cost function by setting $\frac{\partial RSS}{\partial \hat{\beta}_1} = 0$ and $\frac{\partial RSS}{\partial \hat{\beta}_0} = 0$ and solving for β_1 and β_0 respectively.

5.3.2 Derivation of β_0

$$\begin{aligned} \frac{\partial RSS}{\partial \hat{\beta}_0} = 0 &= \frac{\partial}{\partial \hat{\beta}_0} \left[\sum_{i=0}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \right] \\ 0 &= \sum_{i=0}^n \left[\frac{\partial}{\partial \hat{\beta}_0} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \right] \\ 0 &= -2 \sum_{i=0}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \\ 0 &= \sum_{i=0}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \\ 0 &= \sum_{i=0}^n y_i - \sum_{i=0}^n \hat{\beta}_0 - \sum_{i=0}^n \hat{\beta}_1 x_i \end{aligned}$$

$$\begin{aligned}
n\beta_0 &= \sum_{i=0}^n y_i - \hat{\beta}_1 \sum_{i=0}^n x_i \\
\beta_0 &= \frac{\sum_{i=0}^n y_i}{n} - \hat{\beta}_1 \frac{\sum_{i=0}^n x_i}{n} \\
\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}
\end{aligned} \tag{5.12}$$

5.3.3 Derivation of β_1

$$\begin{aligned}
\frac{\partial RSS}{\partial \hat{\beta}_1} &= 0 = \frac{\partial}{\partial \hat{\beta}_1} \left[\sum_{i=0}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \right] \\
0 &= \sum_{i=0}^n \left[\frac{\partial}{\partial \hat{\beta}_1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \right] \\
0 &= -2 \sum_{i=0}^n \left(\frac{\partial}{\partial \hat{\beta}_1} \hat{\beta}_1 x_i \right) (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \\
0 &= -2 \sum_{i=0}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \\
0 &= \sum_{i=0}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \\
0 &= \sum_{i=0}^n (y_i x_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2) \\
0 &= \sum_{i=0}^n (y_i x_i - (\bar{y} - \hat{\beta}_1 \bar{x}) x_i - \hat{\beta}_1 x_i^2) \\
0 &= \sum_{i=0}^n (y_i x_i - \bar{y} x_i - \hat{\beta}_1 \bar{x} x_i - \hat{\beta}_1 x_i^2) \\
0 &= \sum_{i=0}^n (y_i x_i - \bar{y} x_i) - \hat{\beta}_1 \sum_{i=0}^n (\bar{x} x_i - x_i^2) \\
\hat{\beta}_1 &= \frac{\sum_{i=0}^n (y_i x_i - \bar{y} x_i)}{\sum_{i=0}^n (\bar{x} x_i - x_i^2)}
\end{aligned}$$

Equation 5.13 can be further simplified into a more convenient form if we note that:

$$\sum_{i=0}^n (\bar{x}\bar{y} - y_i\bar{x}) = 0 \quad (5.13)$$

$$\sum_{i=0}^n (\bar{x}^2 - x_i\bar{x}) = 0 \quad (5.14)$$

Then:

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=0}^n (y_i x_i - \bar{y} \bar{x}_i) + \sum_{i=0}^n (\bar{x} \bar{y} - y_i \bar{x})}{\sum_{i=0}^n (\bar{x} x_i - x_i^2) + \sum_{i=0}^n (\bar{x}^2 - x_i \bar{x})} \\ \hat{\beta}_1 &= \frac{\sum_{i=0}^n (y_i x_i - \bar{y} \bar{x}_i + \bar{x} \bar{y} - y_i \bar{x})}{\sum_{i=0}^n (\bar{x} x_i - x_i^2 + \bar{x}^2 - x_i \bar{x})} \\ \hat{\beta}_1 &= \frac{\sum_{i=0}^n (y_i x_i - \bar{y} \bar{x}_i + \bar{x} \bar{y} - y_i \bar{x})}{\sum_{i=0}^n (\bar{x} x_i - x_i^2 + \bar{x}^2 - x_i \bar{x})} \end{aligned}$$

$$\hat{\beta}_1 = \frac{\sum_{i=0}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=0}^n x_i - \bar{x}} \quad (5.15)$$

$$\hat{\beta}_1 = \frac{Cov(x, y)}{Var(x)} \quad (5.16)$$

5.4 Multiple Linear Regression

Deriving the solution to the multiple linear regression trained on the **least squares criterion** is very similar to the derivation for simple linear regression we just went through in section 5.3, and the solution takes significantly fewer steps. All coefficients β_0 through β_m can be solved simultaneously, using linear algebra.

The basic equation for multiple linear regression and its remainder are very similar to equation their simple linear regression equivalents: 5.5:

$$\mathbf{y} = X\beta + \epsilon \quad (5.17)$$

$$\begin{aligned}
\mathbf{e} &= \mathbf{y} - X\hat{\beta} \\
SSE &= \mathbf{e}^T \mathbf{e} \\
SSE &= (\mathbf{y} - X\hat{\beta})^T (\mathbf{y} - X\hat{\beta}) \\
SSE &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T X\hat{\beta} - \mathbf{y}^T X^T \hat{\beta} + \hat{\beta}^T X^T X \hat{\beta} \\
SSE &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X\hat{\beta} + \hat{\beta}^T X^T X \hat{\beta} \\
\frac{\partial SSE}{\partial \hat{\beta}} &= \frac{\partial}{\partial \hat{\beta}} (\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X\hat{\beta} + \hat{\beta}^T X^T X \hat{\beta}) \\
0 &= -2X^T \mathbf{y} + 2X^T X \hat{\beta} \\
X^T X \hat{\beta} &= X^T \mathbf{y} \\
\hat{\beta} &= (X^T X)^{-1} X^T \mathbf{y}
\end{aligned}$$

5.5 Pearson Correlation Coefficient

The **Pearson correlation coefficient** is a measure of the linear correlation of two variables. Usually denoted R , the property is equal to the ratio of the covariance to the products of the standard deviations.

$$R = \frac{Cov(x, y)}{\sigma_x \sigma_y} \quad (5.18)$$

$$= \frac{\sum_{i=0}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=0}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=0}^n (y_i - \bar{y})^2}} \quad (5.19)$$

$$= \frac{1}{n-1} \sum_{i=0}^n \frac{x_i - \bar{x}}{\sigma_x} \frac{y_i - \bar{y}}{\sigma_y} \quad (5.20)$$

There are many formulas for R , all of them equivalent, use whichever is most convenient.

$$R \in [-1, 1] \quad (5.21)$$

R is always within the range $[-1, 1]$, where -1 represents perfect negative correlation, and 1 represents perfect positive correlation.

5.6 Goodness of Fit

The **coefficient of determination** is usually used to assess how well a particular linear model fits a particular set of data and is usually reported whenever a linear regression is performed. Denoted R^2 , this number quantitatively indicates how well the regression explains the objective variable. It is equal to unity minus the fraction of unexplained variance:

$$R^2 = 1 - \frac{\sum_{i=0}^n (y_i - \hat{y}_i)^2}{\sum_{i=0}^n (y_i - \bar{y})^2} = 1 - \frac{\mathbf{e}^T \mathbf{e}}{n[\text{Var}(\mathbf{y})]} \quad (5.22)$$

If linear regression is performed correctly, the linear model should never have higher variance than simply using the mean, so:

$$0 \leq R^2 \leq 1 \quad (5.23)$$

Finally, the coefficient of determination is also equal to the square of the correlation coefficient and the square root of the coefficient of determination is equal to the absolute value of the correlation coefficient.

5.6.1 Adjusted R^2 for Multiple Linear Regression

The adjusted R^2 value (R_{adj}^2) is based on the calculation for R^2 but penalizes the model for each additional variable. In theory, it would be possible to fit any arbitrary data set of size m with a linear model using $m - 1$ terms, R_{adj}^2 accounts for this:

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n - 1}{n - k - 1} \quad (5.24)$$

Where:

n = number of samples.

k = number of predictor variables.

It is important to note that R^2 and R_{adj}^2 are not equivalent when $k = 1$, and that R_{adj}^2 is no longer the fraction of explained variance, rather it is best used as a measure of whether the addition of the latest variable to a multiple linear regression is actually a benefit to the model.

5.7 Outliers

5.8 Extrapolation

Extrapolation from a linear model should only be done with extreme caution, or not at all. The data show a correlation between two or more variables within the bounds of the data set, moving beyond those bounds may invalidate the model.

A good example of this is exhibited by springs and Hooke's Law. A spring may respond elastically within a particular range of applied force, but beyond that range (above the yield strength) the spring will begin to exhibit plastic deformation, and not return to its original shape.

5.9 Errors of Estimated Coefficients

5.9.1 Prediction Bands

5.9.2 Confidence Intervals

5.10 Checking Assumptions Using Graphs

5.11 Model Selection Strategies

5.11.1 The R_{adj}^2 Approach

5.11.2 The p-Value Approach

5.12 Quick Methods and Shortcuts

In simple linear regression, the slope $\hat{\beta}_1$ can be estimated to be:

$$\hat{\beta}_{x_1} \approx \frac{s_y}{s_{x_1}} R \quad (5.25)$$

The intercept point $\hat{\beta}_0$ can then be estimated from equation 5.12.

Part IV

Supervised Algorithms I: Classification

Chapter 6

Logistic Regression

6.1 Introduction

Logistic regression is a relatively simple method for binary classification. Multi-class logistic regression is possible, but usually avoided in favor of discriminant analysis, SVM, or some other more advanced method. Logistic regression is common because it is relatively straightforward to interpret: the function computes the natural log of the odds that a particular observation belongs to one class or the other. If we assume that the log odds are a linear function of the components of x , we can use the same approach as in Chapter 5 to train the model.

$$\ln \left[\frac{p(X)}{1 - p(X)} \right] = \beta_0 + \beta_1 x \quad (6.1)$$

The odds at which observation x is assigned to one class or the other can be specified based on circumstance (l), but the best estimator in terms of accuracy is when that limit is set to $\log \text{odds} = 0$, equivalent to $\text{odds} = 1$, and $p(X) = 0.5$. Usually a default value is specified for the special case when $p(X) = l$.

Probability vs Odds

Most people are familiar with the numerical representation of probability. If something has a 10% probability of occurring, then it has a 1/10 chance of occurring. Odds treat the same subject matter with a significantly different approach. While probability is bound numerically on the scale $[0,1]$, odds are represented numerically on a scale $[0, \infty]$. If something has a 1 in 10 chance of occurring, then the odds are 1 to 9, denoted 1:9, and calculated

numerically as $1/9 = 0.\overline{111}$. Just to reiterate, an outcome C_0 with probability $p(C_0) = 0.1$, has odds $od(C_0) = 0.\overline{111}$.

$$Od(X) = \frac{P(X)}{1 - P(x)} \quad (6.2)$$

Using odds has two considerable drawbacks: Firstly, half of the true range of odds is compressed between 0 and 1 while the other half is within all numbers 1 to ∞ . Probability has a nicer symmetry that odds lack. Secondly, and more importantly, odds in the range 0-1 may be mistaken as numerical probabilities of significantly unequal value. Data scientists use odds when working with logistic regression because it is the most directly interpretable value from the function, but I will attempt to make it abundantly clear when to interpret numbers as odds or as probabilities.

6.2 Mathematics

In Chapter 5 we showed how to train models of the form:

$$f(x) = \beta_0 + \beta_1 x \quad (6.3)$$

Linear functions like these are unbound in the values they are allowed to take any real value: $f(x) \in (-\infty, \infty)$. Unlike linear functions, probability is bound within the range $[0,1]$. We want to use the same approach to fitting the function we took in 5, so we need a way to translate the bound function $P(X)$ into the unbound function $f_{lin}(x)$. We do this by calculating the odds instead of the probability, since the odds are bound on only one side $Od(X) \in [0, \infty]$. Finally, we notice that the natural logarithm (or any logarithm for that matter) transforms a number bound on $[0, \infty]$ into an unbound number:

$$\ln(x) \in [0, \infty), \text{ when } x \in [1, \infty) \quad (6.4)$$

$$\ln(x) \in (-\infty, 0], \text{ when } x \in (0, 1] \quad (6.5)$$

This is why we can set the log of the odds equal to a linear function. If solve equation 6.1 for the probability, we find that we've set the probability equal to the logistic function, also known as the logit function, and we've fitted the exponent to a linear function of the feature set X .

$$\begin{aligned} \ln \left[\frac{p(X)}{1-p(X)} \right] &= \beta_0 + \beta_1 x \\ \frac{p(X)}{1-p(X)} &= e^{\beta_0 + \beta_1 x} \\ P(X) &= e^{\beta_0 + \beta_1 x} (1 - P(X)) \\ P(X) &= e^{\beta_0 + \beta_1 x} - e^{\beta_0 + \beta_1 x} P(X) \\ P(X)(1 + e^{\beta_0 + \beta_1 x}) &= e^{\beta_0 + \beta_1 x} \end{aligned}$$

$$P(X) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} \tag{6.6}$$

Part V

Supervised Algorithms III: Algorithms for Classification and Regression

Chapter 7

Naïve Bayes

Naïve Bayes, also called Idiot Bayes, is an extremely powerful method of probability calculation, and therefore prediction. It is powerful because it is insensitive to even severe violations of the assumption of independence. Bayes' Theorem can be thought of as the logical extension of a tree diagram, so we will digress for a moment to get an intuitive understanding of Bayes[3].

7.1 Probability

Calculating **conditional probability**:

$$P(A|B) = \frac{P(A, B)}{P(B)} \quad (7.1)$$

Or equivalently, the **general multiplication rule of probabilities**[3], also known as the **chain rule of probability**[1]:

$$P(A, B) = P(A|B) * P(B) \quad (7.2)$$

7.1.1 Tree Diagrams

A common example of conditional probability is the outcome of medical tests. The probability that someone tests positive or negative for breast cancer (BC) given their situation of either being afflicted or cancer free is usually known, determined experimentally or thorough research.

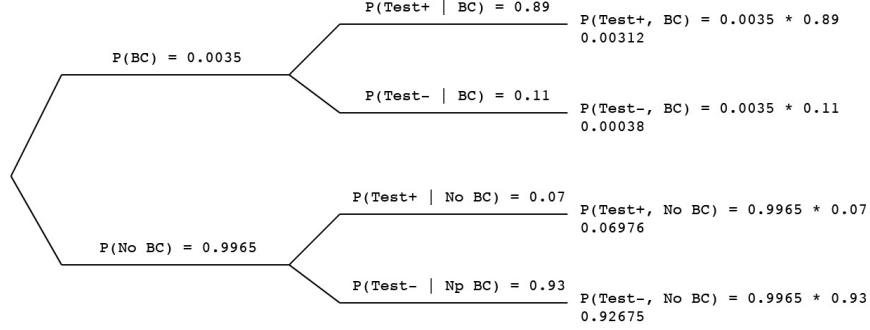


Figure 7.1: A tree diagram showing all possible outcomes of a test for breast cancer.

However, this calculation is not applicable to the real world. We don't know if someone has BC or not, only whether their test result was positive or not. Instead of calculating $P(\text{Test+}|\text{BC})$, $P(\text{Test+}|\text{NoBC})$ we want to know $P(\text{BC}|\text{Test+})$. Using Bayes' Theorem, the tree diagram above contains all the information necessary to calculate $P(\text{BC}|\text{Test+})$. The probability that someone has breast cancer given a positive test can then be calculated:

$$P(\text{BC}|\text{Test+}) = \frac{P(\text{BC}, \text{Test+})}{P(\text{Test+})} \quad (7.3)$$

$$P(\text{BC}, \text{Test+}) = P(\text{BC}) * P(\text{Test+} | \text{BC}) = 0.89 * 0.0035 = 0.00312 \quad (7.4)$$

$$P(\text{Test+}) = P(\text{Test+}, \text{BC}) + P(\text{Test+}, \text{NoBC}) \quad (7.5)$$

$$P(\text{Test+}, \text{NoBC}) = P(\text{NoBC}) * P(\text{Test+} | \text{NoBC}) = 0.9965 * 0.07 = 0.06976 \quad (7.6)$$

$$P(\text{Test+}) = 0.00312 + 0.06976 = 0.07288 \quad (7.7)$$

$$P(\text{BC}|\text{Test+}) = \frac{0.00312}{0.07288} = 0.04281, 4.3\% \quad (7.8)$$

So the probability that someone has cancer, even given a positive test result is only 4.3%. This calculation is predicated on the fact that the sum of all the potential outcomes of some conditional probability must equal unity.

$$\sum_{i=0}^n P(A_i|B) = 1 \quad (7.9)$$

Where: $A_0, A_1, A_2 \dots A_n$ represent all potential outcomes given B.

7.2 Bayes' Theorem

Bayes' Theorem works exactly like equation 7.1, it takes $P(A_0|B)$, then re-normalizes that outcome to the sum of all probabilities ($A_0 \dots A_n$) given B.

$$P(A_0|B) = \frac{P(B|A_0)P(A_0)}{\sum_{i=0}^n P(B|A_i)P(A_i)} \quad (7.10)$$

7.3 Naïve Bayes Applied

Naïve Bayes commonly usually used to classify natural language. We will take the example of classifying tweets as either about an app we're making or about something else, given the words in the tweet. We will classify a tweet as being about our app whenever:

$$P(App|word_0, word_1, word_2 \dots) > P(NotApp|word_0, word_1, word_2 \dots) \quad (7.11)$$

This is called the **maximum a posteriori** rule, or the MAP rule.

Bayes' Theorem indicates that:

$$P(App|words) = \frac{P(App) * P(words|App)}{\sum_{i=0}^n P(A_i|words)P(A_i)} \quad (7.12)$$

And

$$P(NotApp|words) = \frac{P(NotApp) * P(words|NotApp)}{\sum_{i=0}^n P(A_i|words)P(A_i)} \quad (7.13)$$

Where $A_0, A_1, A_2 \dots A_n$ represents all the possible classifications of the tweet.

However, we aren't interested in the absolute probability that the tweet in question is about our app, all we're interested in is whether it's more likely that it's about our app than about something else. Therefore, since the denominators of equations 7.12 and 7.13 are the same, we can compare the numerators directly:

$$P(App) * P(words|App) > P(NotApp) * P(words|NotApp) \quad (7.14)$$

Here's the magic. If we assume that the probability of words appearing in a document are independent of one another then:

$$P(App) * P(words|App) = P(App) * \prod_{i=0}^n P(word_i|App) \quad (7.15)$$

and we need only determine whether

$$P(App) * \prod_{i=0}^n P(word_i|App) > P(NotApp) * \prod_{i=0}^n P(word_i|NotApp) \quad (7.16)$$

This is an exceptional assumption to make. Words are not independent of the word that came before it or comes after it. However, the error is on both sides of the MAP inequality, and so balances out on average.

7.3.1 High-Level Class Probabilities

$P(App)$ and $P(NotApp)$ are the a priori probabilities, i.e. if 20% of all tweets are about your app, then $P(App) = 0.2$, and $P(NotApp) = 0.8$. These values can be changed if you have an asymmetric attitude regarding **precision** and **recall**.

7.3.2 Rare Words

The approach of Naïve Bayes can be subverted if rare words (RW) appear in the test set but not in the training set. If the word hasn't been seen before then the $P(RW|App)$ and $P(RW|NotApp) = 0$. This zero is propagated through the entire product in equation 7.16 and the result is a nonsensical $0 = 0$.

The solution to this problem is called **additive smoothing** where a value of one is added to every word count, whether or not they've been seen before. Therefore words that have been seen n times will show a count of $n+1$.

Chapter 8

Support Vector Machines

8.1 ***To Do***

Chapter 9

Random Forest

9.1 ***To Do***

Chapter 10

Ensemble Methods

10.1 ***To Do***

Part VI

Time Series

Chapter 11

Holt-Winters Forecasting

11.1 ***To Do***

Part VII

Advanced Algorithms

Chapter 12

Neural Networks

12.1 ***To Do***

Part VIII

Supporting Information

Chapter 13

Distance Metrics

13.1 Euclidean Distance

This is the way we usually measure distance. It is the straight-line path in n-dimensional space.

$$Dist_{Euc} = \sqrt{\sum_{i=0}^n (p_i - q_i)^2} \quad (13.1)$$

Where $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$.

13.2 Manhattan Distance (Hamming Distance)

Manhattan distance is the grid distance between two points. If a straight-line trajectory doesn't make sense we can measure the distance as if we were driving a cab in Manhattan, and only take streets perpendicular to the coordinate axes. This is the same as measuring the legs of a right triangle when the Euclidean distance is measuring the hypotenuse.

$$Dist_{Man} = \sum_{i=0}^n (p_i - q_i) \quad (13.2)$$

Where $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$.

13.3 Cosine Distance

Cosine distance is an asymmetric distance metric. When one kind of similarity is more important than another similarity, cosine distance may be a good choice. For example if we are grouping customers then similar purchases are more important than similar non-purchases.

$$Dist_{Cos} = \frac{n_m}{\sqrt{n_0}\sqrt{n_1}} \quad (13.3)$$

Where:

n_m is the number of matched purchases.

n_0 is the total number of items purchased by customer 0.

n_1 is the total number of items purchased by customer 1.

13.4 Jaccard Similarity

Jaccard similarity is the ratio of the magnitude of $A \cap B$ to the magnitude of $A \cup B$. It is like the ratio of the size of the INNER JOIN table to the size of the FULL OUTER JOIN table.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \quad (13.4)$$

Chapter 14

Application Program Interface

14.1 Introduction

Application Program Interfaces (APIs) are a set of **routines**, **protocols**, and tools for building **software applications**. APIs specify how software components should interact with each other, and are used when programming graphical user interface (GUI) components. A good API makes it easier to develop a program by providing all the building blocks.

Chapter 15

Set Theory, Relational Algebra, and Structured Query Language

15.1 Introduction

Structured Query Language (SQL) is the most common syntax for searching databases. SQL is based on set theory and relational algebra, though in the real world, many if not most databases do not follow the tenants of these branches of academia. Much of the skill that employers are looking for in advanced SQL users is how to interact with databases that don't follow these rules with relative ease and rigor.

15.2 Set Theory

Set theory is the branch of mathematical logic that studies sets, which, informally, are collections of objects. SQL is largely based on set theory, and the best maintained databases conform to the best practices of set theory. Of the items listed below, only 1 and 2 are adhered to more or less strictly, all others can and are often violated, especially in small and/or data-naïve, companies.

1. **columns** should represent a unique category of data
2. **order of columns** shouldn't matter
3. **tables** should be the smallest logical subset of data
4. **rows** should represent a unique category of data

5. **attributes** should be unique to the table in which they reside.
6. **primary keys** should be never be repeated.

15.2.1 Basic Concepts and Notation

In	$B \in A$	If o is a member (or element) of A , then $o \in A$. Sets can also be objects, so set B can also be a member of A , $B \in A$.
Proper Subset	$B \subset A$	A proper subset means that B is contained within A and A has at least one element not in B . So $B \subset A, B \neq A$.
Subset	$B \subseteq A$	B is contained within A including the situation where $B = A$
Union	$B \cup A$	All space in B and A , including overlapping space.
Symmetric Difference	$B \triangle A$	All space in B and A not including overlapping space. The complement of intersection.
Relative Complement	$A \setminus B$	The part of A not contained within B . If $A \subset B$, $A \setminus B$ will return NULL
Intersection	$B \cap A$	The space included in both A and B . The complement to symmetric difference.
Cartesian Product	$B \times A$	All possible ordered pairs of A and B .
Power Set	2^A	All possible subsets of A including the empty set and A .

15.3 Structured Query Language

SQL is the most common database query language. It was written to be easy to read and write by humans. SQL does is case and white space agnostic, but convention is to capitalize all SQL commands and not column, table or db names.

SQL is popular because:

1. Its semantics are easy to learn
2. Can directly access data without graphical representation
3. It is easy to replicate and verify SQL queries compared to spreadsheet programs.

15.3.1 Best Practices

1. The most common mistake SQL users make is accidentally running excessively long queries. This can be avoided by running queries on subqueries that have a LIMIT imposed on them. Since inner queries are run first, and outer queries are only run on whatever is SELECT(ed) from the inner query, this technique can ensure that you don't accidentally query millions or billions of rows.
2. Use the EXPLAIN keyword prior to execution, running this command will tell you how many rows will be searched by this query, and give you some idea of the amount of time required to run. This will often catch improperly specified queries that will take much longer than expected.
3. Always select only the minimum amount of data you require. The less data you query, the faster your query will run, and the less hardware you will take. This is especially true when joining tables.
4. Making ER Diagrams and Relational Schema prior to SQL querying will greatly improve query accuracy. These diagrams are especially important before you are familiar or have memorized the database structure.
5. Write queries from the inside out. This way you ensure that your outer queries are acting on the table you expect it to.
6. Clean data prior to analysis. Duplicated rows are multiplied by JOINS, so cleaning prior to query will reduce database load.

15.3.2 Entity Relationship Diagrams

15.3.3 Relational Schema

15.3.4 Database Metadata

SHOW tables Lists the tables contained within the current db

DESCRIBE <table> Lists the columns contained within <table>

15.3.5 Basic Queries

Basic queries have a simple syntax that was written to be easy to be easily read by humans. The issue is that this syntax does not and cannot be directly translated as-is to a computer database for analysis. The order must be rearranged.

In the basic query, the first command `SELECT`, identifies the columns the user is interested in, and the third command `FROM`, identifies the table where the columns exist. It is quickly obvious that this syntax cannot be directly "read" by computers, a database system must know which table to search before it knows what columns it's looking for.

Another item to keep in mind is that `LIMIT` is the final command considered by SQL. If your query includes an aggregation, grouping, or down-selection, these will be executed prior to `LIMIT` being applied, and so the query may take a long time to execute even if only a few rows are actually being returned by the query. Look to section 15.3.1 for solutions to this issue.

Syntax Order	Runtime Order
SELECT	FROM
DISTINCT	WHERE
FROM	GROUP BY
WHERE	HAVING
GROUP BY	SELECT
HAVING	DISTINCT
ORDER BY	ORDER BY
LIMIT	LIMIT

Command	Description
SELECT	Indicates which columns are to be returned
DISTINCT(<column>)	Identifies columns that should be unique in the output. Commonly these columns are the target of aggregation.
FROM	Indicates the table that the query is acting on, this is where nested queries come from.
WHERE	Conditional that can be selected on.
GROUP BY	This is how aggregation functions know what groups to aggregate over.
HAVING	Conditional that can be applied to aggregated columns.
ORDER BY	Order the output by values in a particular column, ASC by default, DESC is optional.
LIMIT	reduces the output to a small number of rows. Can be used to make fast queries of subsets from databases, but should be careful.

15.3.6 Aggregation Functions ***ToDo***

15.3.7 Joins ***ToDo***

15.3.8 Subqueries ***ToDo***

15.3.9 Logical Expressions ***ToDo***

Bibliography

- [1] Foreman, John, *Data Smart: Using Data Science to Transform Information into Insight*, Indianapolis: Wiley, 2014. Print.
- [2] Mason Gallow, *General Assembly Data Science*, Fall 2015.
- [3] Diez, David; Barr, Christopher and Cetinkaya-Rundel, Mine. *OpenIntro Statistics*, 3rd ed. www.openintro.org, Updated: Jan. 13, 2016