
Elements of Statistical Learning

ESLR Notes

Notes by

DANIEL RUIZ

LAST UPDATED: NOVEMBER 2020

Notation:

We'll typically denote input variable by the symbol X . If X is a vector, its components can be accessed by subscripts X_j .

Quantitative outputs will be denoted by Y , and qualitative outputs by G .

We use uppercase letters X , Y or G when referring to the generic aspects of a variable. Observed values are written in lowercase; hence the i th observed value of X is written as x_i (x_i is again a scalar or vector).

Matrices are represented by bold uppercase letters. An example of this would be a set of N input p -vectors x_i , $i=1,\dots,N$ (each x_i is a vector of p components) would be represented by $N \times p$ matrix \mathbf{X} .

We use a hat to signify a prediction and withhold a hat to signify a measurement. Hence, if Y denotes a measured value, \hat{Y} would represent our predicted value.

Contents

2	Overview of Supervised Learning	5
2.1	Introduction	5
2.2	Variable Types and Terminology	5
2.2.1	Two Simple Approaches to Prediction: Least Squares and Nearest Neighbors	6
2.2.2	Linear Models and Least Squares	6
2.2.3	Nearest-Neighbour Methods	7
2.3	Statistical Decision Theory	8
2.4	Local Methods in High Dimensions	13
2.5	Statistical Models, Supervised Learning and Function Approximation	16
2.5.1	A Statistical Model for the Joint Distribution $\Pr(X,Y)$	16
2.5.2	Function Approximation	17
2.6	Structured Regression Models	19
2.6.1	Difficulty of the Problem	19
2.7	Classes of Restricted Estimators	20
2.7.1	Roughness Penalty and Bayesian Methods	20
2.7.2	Kernel Methods and Local Regression	20
2.7.3	Basis Functions and Dictionary Methods	21
2.8	Model Selection and the Bias-Variance Tradeoff	22
3	Linear Methods for Regression	25
3.1	Introduction	25
3.2	Linear Regression Models and Least Squares	25
3.3	Subset Selection	31
3.3.1	Best-Subset Selection	31
3.3.2	Forward- and Backward-Stepwise Selection	32
3.3.3	Forward-Stagewise Regression	32
3.4	Shrinkage Methods	33
3.4.1	Ridge Regression	33
3.4.2	The Lasso	36
3.5	Principal Component Analysis	36
4	Linear Methods for Classification	41
4.1	Introduction	41
4.2	Linear Regression of an Indicator Matrix	42
4.3	Linear Discriminant Analysis	43
A	Probability Theory	45
A.1	Constructing a Probability Space	45
A.2	Standard Definitions and Properties	46
A.3	Distributions and Densities	49
A.4	Expectations	53
A.5	Jointly Distributed Random Variables	56
A.6	Expectations and the Central Limit Theorem	62

B Linear Algebra**65**

Chapter 2

Overview of Supervised Learning

2.1 Introduction

Definition 2.1: Inputs

These are the set of variables that are measured / preset. In the statistical and pattern recognition literature, this also goes by the name of predictors and features. Classically, they are called the independent variables.

Definition 2.2: Output

The inputs influence one or more of these. These are the results. Similarly, output are also referred to as the responses or classically as the dependent variables.

Definition 2.3: Supervised Learning

Supervised learning is the process of using inputs to predict the value of outputs.

2.2 Variable Types and Terminology

Definition 2.4: Categorical / Discrete Variables

Outputs will generally vary in their nature. An example of this is the qualitative character of outputs such as eye color assuming values in a finite set $\mathcal{G} = \{\text{blue, brown, green}\}$. Such a class isn't equipped with any explicit ordering and often descriptive labels rather than numbers are used to denote classes. Qualitative variables are also referred to as categorical or discrete variables as well as factors.

Definition 2.5: Regression and Classification

Distinction in output type has led to a naming convention among prediction tasks. We use the term regression to generally refer to predicting quantitative outputs whereas classification is reserved for qualitative outputs.

Definition 2.6: Ordered Categorical

A third variable type that deals with notions such as *small*, *medium* and *large*, with an ordering between the values, but no metric is appropriate (difference between medium and small need not be same as between large and medium).

A cursory description of the learning task is: Given the value of an input vector X , make a good prediction of the output Y , denoted by \hat{Y} . Hence, if Y takes on values in \mathbb{R} , then so should \hat{Y} . Follows analogously for categorical outputs.

Definition 2.7: Training Data

To construct prediction rules, we require data and often a lot of it. Hence, we suppose that we have available a set of measurements (x_i, y_i) or (x_i, g_i) , $i = 1, \dots, N$, which is known as the *training data*.

2.2.1 Two Simple Approaches to Prediction: Least Squares and Nearest Neighbors

Brief statement on *Least Squares* and *K-nearest Neighbours*:

- *Least Squares*: Large assumption about the structure but stable. Possibly may give inaccurate predictions.
- *K-nearest*: Mild structural assumptions with predictions often accurate but possibly unstable.

2.2.2 Linear Models and Least Squares**Definition 2.8: Linear Model and Bias**

Given some vector of inputs $X^T = (X_1, X_2, \dots, X_p)$, we predict the output Y through the model

$$\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^p X_j \hat{\beta}_j \quad (2.1)$$

$\hat{\beta}_0$ is the *intercept*, which is also known as the *bias* in machine learning. We note that (2.1) is a general expression for a linear map, taking $X \mapsto \hat{Y}$. It's convenient to absorb the 1 into the X vector through a redefinition. Hence, defining the vector $\hat{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$, we can alternatively express our linear model (2.1) as:

$$\hat{Y} = X^T \hat{\beta} \quad (2.2)$$

Definition 2.9: Hyperplane in Input-Output Space

In the $(p+1)$ -dimensional input-output space S , (X, \hat{Y}) would represent a hyperplane. That is,

$$(X, \hat{Y}) = \{(x, y) \in S \mid x \in X, y = X^T \beta\} \quad (2.3)$$

Proposition 2.1: Gradient of the Linear Model

Suppose we view the function over p -dimensional input space by $f(X) = X^T \beta$. Then the gradient $f'(X)$ is

given by:

$$f'(X) = \beta \quad (2.4)$$

We note that β is the vector in input space that points in the direction of steepest ascent. This fact is observed from considering the directional derivative. A standard exercise is showing that it is maximized in the direction of the gradient.

Definition 2.10: Residual Sum of Squares for Linear Model, RSS

Suppose that we have a set of N data points (i.e pairing of observed input-outputs (x, y)). This set would be given by $\{(x_i, y_i) | i \in \mathbb{Z}_N\}$. We define the residual sum of squares as follows:

$$RSS(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2 \quad (2.5)$$

We note that generically x_i and β are p -vectors, hence the notation. Notice that RSS essentially measures the level of deviation from our observed value y_i with what would be predicted by the linear model $x_i^T \beta$ for some β . The goal therefore becomes in finding a suitable choice for β that minimizes RSS.

Proposition 2.2: RSS Minimization

Given the residual sum of squares, $RSS(\beta)$ as defined Def 2.10 and provided that $\mathbf{X}^T \mathbf{X}$ is invertible, then RSS is minimized when β takes on the following value:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y \quad (2.6)$$

Where we have defined \mathbf{X} as the $N \times p$ matrix with x_i forming the i^{th} row of the matrix. In essence, $\mathbf{X}_{ij} = (x_i)_j$. Hence, the fitted value at the i^{th} input x_i is given by $\hat{y}_i = \hat{y}(x_i) = x_i^T \hat{\beta}$. Insofar as computation goes, we have the identity:

$$\mathbf{X}^T \mathbf{X} = \sum_k x_k x_k^T \quad (2.7)$$

2.2.3 Nearest-Neighbour Methods

This method aims to utilize the observations in the training set that are closest to some fixed point x in the input space. In essence, if we wanted to predict the output of some point in our input space, we search the immediate vicinity so as to establish whether the local configuration of observed data would favour a prediction at the desired point to go in a particular direction. Suppose that there were two outputs of {blue, red} and I selected some point completely surrounded by blue data. This method is constructed so as to predict the output to be blue as well.

Definition 2.11: K-Nearest Neighbours

Let (x_i, y_i) represent an input-output pair from our training sample (observed data). Then the k -nearest neighbours fit for \hat{Y} is defined as follows:

$$\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i \quad (2.8)$$

Where $N_k(x)$ denotes the k closest points in the training sample lying in the neighbourhood of x . We note that invoking a notion of closeness means we have assumed a metric space. Typically, we will assume the

metric to be the Euclidean distance.

2.3 Statistical Decision Theory

Definition 2.12: Loss Function

We aim to define a function $f(X)$ for predicting Y , given values of the input X . Hence, we require the notion of a loss function $L(Y, f(X))$ for penalizing errors in prediction.

Definition 2.13: Squared Error Loss Function

The most common and convenient loss function is squared error loss:

$$L(Y, f(X)) = (Y - f(X))^2 \quad (2.9)$$

Proposition 2.3

Let $X \in \mathbb{R}^p$ denote a real valued random input vector (of dimension p), and $Y \in \mathbb{R}$ a real valued random output vector with joint distribution $Pr(X, Y)$. The criterion for choosing f will now be demonstrated. We define the expected (squared) prediction error by

$$EPE(f) = E[(Y - f(X))^2] \quad (2.10)$$

where the expectation is with respect to the joint distribution. We will show that the solution to the function f that minimizes EPE is given by the regression function $f(x) = E[Y|x = x]$.

Proof. We first define the joint density $pr(x, y)$ to satisfy $Pr(x, y) = \int_{-\infty}^x \int_{-\infty}^y pr(u, v) du dv$. We also note from the definition of conditional density, we have that $pr(x, y) = pr(y|x)pr(x)$ where $pr(x)$ denotes the marginal density of the random variable X . We therefore have

$$EPE(f) = \int \int [y - f(x)]^2 pr(x, y) dx dy \quad (2.11)$$

$$= \int \left(\int [y - f(x)]^2 pr(y|x) dy \right) pr(x) dx \quad (2.12)$$

We observe that EPE is a functional of f and can therefore be subject to variational techniques. For convenience, we define $EPE(f) = \int \mathcal{L}_{EPE}[f] dx$ where \mathcal{L}_{EPE} is explicitly defined by

$$\mathcal{L}_{EPE}(f) = pr(x) \int [y - f(x)]^2 pr(y|x) dy. \quad (2.13)$$

Observe that EPE is bounded below since it's argument is positive definite and therefore extremizing should yield a minimization solution^a. Observe that the Euler-Lagrange equation for this system is necessarily:

$$\frac{\delta}{\delta f} EPE = \frac{\partial \mathcal{L}_{EPE}}{\partial f} = 0 \quad (2.14)$$

We now compute the partial derivative of \mathcal{L}_{EPE} with respect to f :

$$\frac{\partial \mathcal{L}_{EPE}}{\partial f} = 2 pr(x) \int [y - f(x)] pr(y|x) dy \quad (2.15)$$

$$= 2 pr(x) \left(\int y pr(y|x) dy - f(x) \int pr(y|x) dy \right) \quad (2.16)$$

$$= 2 pr(x) (E[Y|X = x] - f(x)), \quad (2.17)$$

where we have used the fact that $\int pr(y|x) dy = 1$. Hence, we can therefore conclude that

$$f(x) = E[Y|X = x] \quad (2.18)$$

□

^aNote that $pr(x)$ and $pr(y|x)$ are probability densities, which are always positive definite.

Proposition 2.4

An alternative procedure to Proposition 2.3 involves a pointwise minimization, but this will turn out to be an effectively equivalent proof. Let $X \in \mathbb{R}^p$ denote a real valued random input vector (of dimension p), and $Y \in \mathbb{R}$ a real valued random output vector with joint distribution $\Pr(X, Y)$. We define the *expected (squared) prediction error* by

$$\begin{aligned} EPE(f) &= E[(Y - f(X))^2] \\ &= \int \left(\int (y - f(x))^2 pr(y|x) dy \right) pr(x) dx \\ &= \int E_{Y|X}[(y - f(x))^2 | x] pr(x) dx \\ &= E_X[E_{Y|X}[(Y - f(X))^2 | X]]. \end{aligned} \quad (2.19)$$

Then, it's sufficient to minimize this pointwise, which is essentially what was performed in Proposition 2.3. In essence, the solution to minimizing EPE is given by

$$f(x) = \underset{c \in \mathbb{R}}{\operatorname{argmin}} E_{Y|X}[(Y - c)^2 | X = x]. \quad (2.20)$$

One can observe this being true as $E_{Y|X}$ and $pr(x)$ are positive definite. The expectation, E_X thereby preserves this ordering. The solution is therefore given by

$$f(x) = E[Y|X = x] \quad (2.21)$$

Proof. We define $g(c, x) := E_{Y|X}[(Y - c)^2 | X = x]$. We observe that g is clearly continuous and differentiable in c . It's also positive definite and thereby has a minima. Simple differentiation gives us

$$\frac{\partial g}{\partial c} = 2 \int (y - c) pr(y|x) dy. \quad (2.22)$$

Hence,

$$c_{min} = \int y pr(y|x) dy, \quad (2.23)$$

which satisfies $0 = \frac{\partial g}{\partial c}(c_{min})$, thereby setting $f(x) = E[Y|X = x]$. □

Nearest Neighbours Implementation

Nearest neighbours tries to directly implement this recipe with training data. In essence, we essentially ask for the average of all y'_i s with input $x_i = x$. One can settle for

$$\hat{f}(x) = \operatorname{Ave}(y_i | x_i \in N_k(x)), \quad (2.24)$$

where Ave denotes average and $N_k(x)$ is a neighbourhood containing k points in \mathcal{T} closest to x . We note that there are two approximation being used here:

1. The Expectation is approximated by an averaging over sample data.
2. Conditioning at a point is relaxed to conditioning on some region close to the target point.

Under mild regularity conditions of $\Pr(x,y)$, one can show that as $N, k \rightarrow \infty$ such that $k/N \rightarrow 0$, then $\hat{f}(x) \rightarrow E(Y|X = x)$.

Proposition 2.5

Let $X \in \mathbb{R}^p$ denote a real valued random input vector (of dimension p), and $Y \in \mathbb{R}$ a real valued random output vector with joint distribution $\Pr(X, Y)$. Suppose that we approximate $f(x)$ to be a linear function, given below by

$$f(x) \approx x^T \beta \quad (2.25)$$

Then, EPE is given by $EPE(\beta) = E[(Y - X^T \beta)^2]$. We will show that the solution to β is given by $\beta = (E[XX^T])^{-1}E[XY]$.

Proof. Similarly, we express our expectation as before

$$EPE(\beta) = \int \int [y - x^T \beta] pr(x, y) dx dy \quad (2.26)$$

We aim to minimize as before, this time noting that EPE is regular function of β . Hence, we aim to solve for $\nabla(EPE) = 0$ (set the gradient to zero). We compute

$$\frac{\partial}{\partial \beta_i} EPE(\beta) = 2 \int \int (y - x^T \beta) x_i dx dy \quad (2.27)$$

Hence, the gradient is given below by (Where we have assumed our space to be Euclidean)

$$\nabla(EPE) = 2 \int \int (y - x^T \beta) x dx dy \quad (2.28)$$

We now press onwards to solve for β

$$0 = \int \int yx dx dy - \int \int xx^T \beta dx dy \quad (2.29)$$

$$= E[XY] - \left(\int \int xx^T dx dy \right) \beta \quad (2.30)$$

$$= E[XY] - E[XX^T] \beta \quad (2.31)$$

$$(2.32)$$

We emphasize that $x^T \beta$ is a scalar and we can therefore move the x term across it. We have also moved β outside the integral since it has no x or y dependence. Note that $E[XX^T]$ is a matrix valued expectation (in particular, it is the cross-correlation matrix of \mathbf{X} with itself) which we assume is invertible and therefore yields the unique solution given by

$$\beta = (E[XX^T])^{-1} E[XY] \quad (2.33)$$

□

Least Squares and K-nearest Neighbours

Both Least-Squares and k -nearest neighbours end up approximating conditional expectations by averages.

- Least-Squares assumes that $f(x)$ is well approximated by a globally linear function.
- k -nearest neighbours assumes that $f(x)$ is well approximated by a locally constant function.

Additive Models

Many of the techniques that we will encounter in this book are model based, although more flexible than the rigid linear model. For instance, additive models assume that

$$f(X) = \sum_{j=1}^p f_j(X_j) \quad (2.34)$$

This retains additivity of linear model, but each coordinate function f_j is arbitrary.

Proposition 2.6

Let $X \in \mathbb{R}^p$ denote a real valued random input vector (of dimension p), and $Y \in \mathbb{R}$ a real valued random output vector with joint distribution $\Pr(X, Y)$. We now instead consider the following loss function: $L_1 := |Y - f(X)|$. We define the expected absolute prediction error by

$$EPE(f) = E[|Y - f(X)|]. \quad (2.35)$$

The minimizing solution is given by

$$\hat{f}(x) = \text{median}(Y|X = x) \quad (2.36)$$

Proof. We expand the integral as usual:

$$EPE(f) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |y - f(x)| \Pr_{X,Y}(x, y) \, dx \, dy \quad (2.37)$$

$$= \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} |y - f(x)| \Pr_{Y|X}(y|x) \, dy \right) \Pr_X(x) \, dx. \quad (2.38)$$

It's sufficient to minimize point-wise. Hence, our solution is encoded as follows:

$$f(x) = \underset{c \in \mathbb{R}}{\operatorname{argmin}} \int_{-\infty}^{\infty} |y - c| \Pr_{Y|X}(y|x) \, dy. \quad (2.39)$$

We define a new function $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$g(c, x) = \int_{-\infty}^{\infty} |y - c| \Pr_{Y|X}(y|x) \, dy \quad (2.40)$$

$$= \int_{-\infty}^c (c - y) \Pr_{Y|X}(y|x) \, dy + \int_c^{\infty} (y - c) \Pr_{Y|X}(y|x) \, dy. \quad (2.41)$$

We compute the partial w.r.t c :

$$\frac{\partial g}{\partial c} = \int_{-\infty}^c \Pr_{Y|X}(y|x) \, dy - \int_c^{\infty} \Pr_{Y|X}(y|x) \, dy \quad (2.42)$$

$$= \int_{-\infty}^c \Pr_{Y|X}(y|x) \, dy - \left[1 - \int_{-\infty}^c \Pr_{Y|X}(y|x) \, dy \right] \quad (2.43)$$

$$= 2 \int_{-\infty}^c \Pr_{Y|X}(y|x) \, dy - 1 \quad (2.44)$$

Hence, $\frac{\partial g}{\partial c} = 0$ when we have

$$\int_{-\infty}^c pr_{Y|X}(y|x) = \frac{1}{2}. \quad (2.45)$$

We note that g is positive definite, hence when c satisfies the above relation, we have a minimizing solution. We note that by definition, the median is the value m satisfying $Pr(Y \leq m|X = x) = 1/2$. Hence,

$$f(x) = \text{median}(Y|X = x) \quad (2.46)$$

□

Categorical Variable G

If the output is instead a categorical variable G , then we have to use a different loss function for penalizing prediction errors. We denote \mathcal{G} as the set of all possible classes. Hence, an estimate \hat{G} will assume values in \mathcal{G} . Our loss function could be represented by a $K \times K$ matrix \mathbf{L} , where $K = \text{card}(\mathcal{G})$ (Cardinality of \mathcal{G}).

We would therefore construct \mathbf{L} such that $\mathbf{L}(k, l)$ is the penalty for classifying an observation belonging to class \mathcal{G}_k as \mathcal{G}_l . Hence, we would want \mathbf{L} to be zero along the diagonal and non-negative elsewhere.

Definition 2.14: Zero-One Loss Function

The zero-one loss function is a loss function where all misclassifications are charged a single unit. In essence, let \mathbf{L}_{ZO} denote the zero-one loss function. Then,

$$\mathbf{L}_{ZO}(k, l) = 1 - \delta_{kl} = \begin{cases} 1 & \text{if } k \neq l \\ 0 & \text{if } k = l \end{cases} \quad (2.47)$$

Proposition 2.7: Bayes Classifier

Let G denote a categorical variable and \hat{G} denote an estimate. We taken \mathcal{G} to be the set of all possible categorical values, defining $K := \text{card}(\mathcal{G})$. Hence, \mathcal{G}_j represents the j^{th} categorical value in the set. Let L denote the loss function for this categorical variable and estimate. The expected prediction error is defined as

$$EPE = E[L(G, \hat{G}(X))], \quad (2.48)$$

where expectation is taken with respect to the joint distribution $Pr(G, X)$. The minimizing solution for $\hat{G}(x)$ is given by

$$\hat{G}(x) = \mathcal{G}_k \text{ if } Pr(\mathcal{G}_k|X = x) = \max_{g \in \mathcal{G}} Pr(g|X = x). \quad (2.49)$$

This solution is known as the *Bayes classifier*, which says that we classify to the most probable class, using the conditional (discrete) distribution $Pr(G|X)$.

Proof. We begin by writing out the full loss function:

$$EPE = E_X \left[\sum_{k=1}^K L(\mathcal{G}_k, \hat{G}(X)) Pr(\mathcal{G}_k|X = x) \right]. \quad (2.50)$$

Subjecting this to pointwise minimization and substituting the relevant loss function, the minimizing solution

is given by

$$\hat{G}(x) = \operatorname{argmin}_{g \in \mathcal{G}} \sum_{k=1}^K L(\mathcal{G}_k, g) \Pr(\mathcal{G}_k | X = x) \quad (2.51)$$

$$= \operatorname{argmin}_{g \in \mathcal{G}} \sum_{k=1}^K (1 - \delta_{g\mathcal{G}_k}) \Pr(\mathcal{G}_k | X = x) \quad (2.52)$$

$$= \operatorname{argmin}_{g \in \mathcal{G}} [1 - \Pr(g | X = x)] \quad (2.53)$$

$$= \operatorname{argmax}_{g \in \mathcal{G}} \Pr(g | X = x) \quad (2.54)$$

Alternatively, this can be expressed by

$$\hat{G}(x) = \mathcal{G}_k \text{ if } \Pr(\mathcal{G}_k | X = x) = \max_{g \in \mathcal{G}} \Pr(g | X = x). \quad (2.55)$$

□

Definition 2.15: Bayes Rate

The error rate of the Bayes classifier is called the Bayes rate.

2.4 Local Methods in High Dimensions

Curse of Dimensionality

There are many manifestations of this problem, and we will examine a few here

Proposition 2.8: The Median Distance for N data points in p -dimensions

Consider N data points uniformly distributed in a p -dimensional unit ball centered at the origin. Suppose that we consider a nearest-neighbour estimate at the origin. The median distance from the origin to the closest data point is given by the expression

$$d(p, N) = \left(1 - \frac{1}{2}^{1/N}\right)^{1/p} \quad (2.56)$$

Proof. Let F denote the distribution. We first want to construct the density and distribution for these random variables. We note that the distribution for a point to lie in a concentric sphere of radius r is given by

$$F(r) = \int_{B_r(0)} d^p x f(x), \quad (2.57)$$

where $B_r(0)$ denotes the r -ball centered on the origin and $f(x)$ is the corresponding uniform density satisfying $f(x) = 1/\operatorname{Vol}(B_1(0))$. The volume of the unit ball in p -dimensions is given by

$$\operatorname{Vol}(B_1(0)) = \int_{B_1(0)} d^p x = \int d\Omega_{p-1} \int_0^1 dx x^{p-1} = \frac{\Omega_{p-1}}{p}, \quad (2.58)$$

where Ω_{p-1} denotes the surface area of the unit p -sphere. We can therefore compute the probability that a uniformly distributed variable in the p -dimensional unit ball will lie in some concentric sphere of radius r

by

$$Pr(X \leq r) = F(r) = \int d\Omega_{p-1} \int_0^r dx \left(\frac{p}{\Omega_{p-1}} \right) x^{p-1} = r^p. \quad (2.59)$$

We are interested in the median distance from the origin to the **closest** data point. To accommodate this, we require the notion of order statistics among N independent identically distributed random variables. In particular, suppose that we have N data points, then let's denote the corresponding ordered random variables by $X_{(1)}, X_{(2)}, \dots, X_{(N)}$, with $X_{(j)}$ denoting the j^{th} order statistic. By Proposition A.23, we have the closest statistic distribution

$$F_{(1)}(r) = Pr(X_{(1)} \leq r) = \sum_{k=1}^N \binom{N}{k} (1 - F(r))^{N-k} (F(r))^k = 1 - (1 - F(r))^N = 1 - (1 - r^p)^N, \quad (2.60)$$

where the second-last equality is obtained via binomial identity. We note that the median value is by definition the value d that satisfies

$$\frac{1}{2} = F_{(1)}(d). \quad (2.61)$$

Hence, solving for d in

$$\frac{1}{2} = 1 - (1 - d^p)^N, \quad (2.62)$$

we obtain

$$d(p, N) = \left(1 - \frac{1}{2}^{1/N} \right)^{1/p} \quad (2.63)$$

□

Lemma 2.1: Monotonicity of Exponentiation

Let $0 < \alpha < 1$. If $m < n$, then we have that

$$\alpha^{1/m} < \alpha^{1/n}. \quad (2.64)$$

The inequality is flipped if $\alpha > 1$.

Example 2.1: Curse of Dimensionality

By Proposition 2.8, we have a function for the closest median distance for N identical uniformly distributed random variables in the p -ball. Given that $0 < 1 - \frac{1}{2}^{1/N} < 1$ for any $N \in \mathbb{N}$, by Lemma ?? we can observe that for a fixed number of data points, they will tend to be closer to the boundary of the p -ball in higher dimensions. That is, if $p < q$, then $d(p, N) < d(q, N)$. This presents a problem as prediction is much more difficult near the edges of the training sample.

Another manifestation of the curse is that the sampling density is proportional to $N^{1/p}$, where p is the dimension of the input space and N is the sample size. Thus, in high dimensions all feasible training samples sparsely populate the input space.

Definition 2.16: Bias

Let $\hat{\nu}$ represent an estimator and ν be a parameter. Then, the bias of our estimator is defined by

$$\text{Bias}(\hat{\nu}) = E[\hat{\nu}] - \nu \quad (2.65)$$

This quantity is also referred to as the long-run average error of $\hat{\nu}$.

Proposition 2.9: Bias-Variance Decomposition

Let $\mathcal{T} = \{(x_j, y_j) : j \in \mathcal{I}\}$ be the set of all N data points, where $\mathcal{I} \simeq \mathbb{Z}_N$ is some index set^a and $x_j \in \mathbb{R}^p$, $y_j \in \mathbb{R}^d \forall j$ (for some $p, d \in \mathbb{N}$). When developing a learning algorithm, one is interested in partitioning the data set into a training data set as well as a test set to ensure that the algorithm is generalizing well. The way this is typically done is at random to ensure no selective bias enters into the training as this would hinder the algorithms ability to generalize. Suppose that we are interested in establishing a training set $D_s \subset \mathcal{T}$ of size $n < N$. Let X_i denote the random variable associated with the i^{th} data point in D_s ^b, drawn from the set \mathcal{T} by some joint probability density $Pr_D(X_1, \dots, X_n)$ where we have defined a random variable $D = (X_1, \dots, X_n)$ that generates our training data. In a supervised learning context, our prediction variable is taken to have the form

$$Y = f(X) + \epsilon, \quad (2.66)$$

where ϵ is a random variable satisfying $E_\epsilon[\epsilon] = 0$ and f is the ‘true’ underlying function that relates the x_j to y_j in the data. We want to measure the MSE of our model \hat{f}_D which aims to mimic f as closely as possible. We emphasize that \hat{f}_D is dependent on the training data D_s as this data is precisely what we use to inform us about \hat{f}_D ’s parameters. Hence, \hat{f}_D is a statistic as it is a function of the data. We now want to see how well our model \hat{f}_D performs on a fixed test variable $x \in \mathcal{T} \setminus D_s$, which can be decomposed into bias and variance components:

$$E_D[(Y - \hat{f}_D(x))^2] = \text{Var}_\epsilon[\epsilon] + [\text{Bias}_D(\hat{f}_D(x))]^2 + \text{Var}_D[\hat{f}_D(x)], \quad (2.67)$$

where E_D denotes the expectation with respect to the training data random variable D .

Proof. Although the proof is straightforward algebra, the underlying distributions that we take our expectation with respect to are paramount in reaching the desired result. With everything laid out, it should be clear how to perform this but I will nevertheless try to precisely demonstrate how the decomposition arises:

$$E_D[(Y - \hat{f}_D(x))^2] = E_D[(f(x) + \epsilon - \hat{f}_D(x))^2] \quad (2.68)$$

$$= E_D[(f(x) - \hat{f}_D(x))^2] + 2E_D[\epsilon(f(x) - \hat{f}_D(x))] + E_D[\epsilon^2] \quad (2.69)$$

$$= E_D[(f(x) - E_D[\hat{f}_D(x)] + E_D[\hat{f}_D(x)] - \hat{f}_D(x))^2] + \text{Var}_\epsilon[\epsilon] \quad (2.70)$$

$$= E_D[(f(x) - E_D[\hat{f}_D(x)])^2] + E_D[(E_D[\hat{f}_D(x)] - \hat{f}_D(x))^2] \quad (2.71)$$

$$+ 2E_D[(f(x) - E_D[\hat{f}_D(x)])(E_D[\hat{f}_D(x)] - \hat{f}_D(x))] + \text{Var}_\epsilon[\epsilon] \quad (2.72)$$

$$= (f(x) - E_D[\hat{f}_D(x)])^2 + \text{Var}_D[\hat{f}_D(x)] + \text{Var}_\epsilon[\epsilon] \quad (2.73)$$

$$+ 2(f(x) - E_D[\hat{f}_D(x)])E_D[E_D[\hat{f}_D(x)] - \hat{f}_D(x)] \quad (2.74)$$

$$= [\text{Bias}_D(\hat{f}_D(x))]^2 + \text{Var}_D[\hat{f}_D(x)] + \text{Var}_\epsilon[\epsilon] \quad (2.75)$$

To reach the end result, we used the fact that $f(x) - E_D[\hat{f}_D(x)]$ is a constant with respect to E_D (as both $f(x)$ and $E_D[\hat{f}_D(x)]$ are constants with respect to the expectation E_D). In addition, one can easily see that $E_D[E_D[\hat{f}_D(x)] - \hat{f}_D(x)] = 0$, thereby eliminating the last term from the second-last equality. \square

^aTypically, one takes $\mathcal{I} = \mathbb{Z}_N$ or $\mathcal{I} = \mathbb{Z}_N + 1$ unless there is some notation that would make an alternative index more convenient.

^bWith $X_i = (x_j, y_j)$ being the possible events in our probability space for some $j \in \mathcal{I}$.

2.5 Statistical Models, Supervised Learning and Function Approximation

Failure of Nearest-Neighbours

Our primary goal is to find a useful approximation $\hat{f}(x)$ to the function $f(x)$ that underlies the predictive relationship between the inputs and outputs. The class of nearest-neighbor methods can be viewed as direct estimates of the regression function seen in §2.3: $f(x) = E(Y|X = x)$. However, we have seen that they can fail in at least two ways:

- If dimension of input space is high, the nearest neighbors need not be close to the target point, and can result in large errors;
- If special structure is known to exist, this can be used to reduce both the bias and variance of the estimates.

2.5.1 A Statistical Model for the Joint Distribution $\Pr(X,Y)$

Proposition 2.10

Suppose that our data arose from a statistical model

$$Y = f(X) + \epsilon, \quad (2.76)$$

where the random error ϵ has $E(\epsilon) = 0$ and is independent of X . We note that for this model, using the squared error loss function gives us the solution $\hat{f}(x) = E[Y|X = x]$.

Proof. We compute the squared error loss:

$$EPE(f) = E[(Y - \hat{f}(X) - \epsilon)^2] = E[(Y - \hat{f}(X))^2 - 2\epsilon(Y - \hat{f}(X)) + \epsilon^2] \quad (2.77)$$

$$= E[(Y - \hat{f}(X))^2] - 2E[\epsilon]E[Y - \hat{f}(X)] + E[\epsilon^2] \quad (2.78)$$

$$= E[(Y - \hat{f}(X))^2] + Var[\epsilon] \quad (2.79)$$

The minimizing solution to this was shown in Proposition 2.3. We note that the presence of $Var[\epsilon]$ is irrelevant as it serves as background noise, that is not dependent on $f(x)$. Hence, we have our solution

$$\hat{f}(x) = E[Y|X = x] \quad (2.80)$$

□

This additive error model is a useful approximation to the truth. For most system, the input-output pairs (X,Y) will not have a deterministic relationship $Y = f(X)$. There will generally be other unmeasured variables that also contribute to Y , such as measurement error. The additive model therefore assumes that we can capture all these departures from a deterministic relationship through the error ϵ .

The assumption in (2.76) is that errors are independent and identically distributed. These assumptions are not strictly necessary but seems like a good model when we average squared errors uniformly in our EPE criterion.

Supervised Learning

We want to present the function-fitting paradigm from a machine learning point of view. Suppose for simplicity that the errors are additive and that the model $Y = f(X) + \epsilon$ is a reasonable assumption.

Supervised learning aims to learn f by example through a *teacher*. One assembles a *training* set of observations $\mathcal{T} = (x_i, y_i)$, $i = 1, \dots, N$ for the inputs and outputs in the system of study. There is a learning algorithm that is fed observed input values x_i which produces outputs $\hat{f}(x_i)$ in response to the inputs. The learning algorithm has the property that it can modify its input/output relationship \hat{f} in response to differences $y_i - \hat{f}(x_i)$ between original and generated outputs. This process is known as *learning by example*. Upon completion of the learning process, the hope is that artificial and real outputs are sufficiently close enough to be useful for all sets of inputs likely encountered in practice.

2.5.2 Function Approximation

Learning paradigm of previous section has been the motivation for research into supervised learning problems in the field of machine learning (with analogies to human reasoning) and neural networks (with biological analogies to the brain). One can consider the data pairs $\{x_i, y_i\}$ as points in $(p+1)$ -dimensional Euclidean space. The function $f(x)$ has domain equal to the p -dimensional input subspace, and is related to the data via a model such as $y_i = f(x_i) + \epsilon_i$. For convenience, we will assume the domain to be \mathbb{R}^p . The goal is to obtain a useful approximation to $f(x)$ for all x in some region of \mathbb{R}^p , given the representations in \mathcal{T} .

Instead of the aforementioned learning paradigm, we will treat supervised learning as a problem in function approximation. This encourages the geometrical concepts of Euclidean spaces and mathematical concepts of probabilistic inference to be applied to the problem. This will be the approach taken in this book.

Definition 2.17: Linear Basis Expansions

Many of the approximation that we will encounter have an associated set of parameters θ that can be modified to suit the data at hand. For instance, the linear model $f(x) = x^T \beta$ has $\theta = \beta$. Another class of useful approximators can be expressed as *linear basis expansions*

$$f_{\theta}(x) = \sum_{k=1}^K h_k(x) \theta_k, \quad (2.81)$$

where the h_k are a suitable set of functions or transformations of the input vector x . Some traditional examples for h_k are polynomials and trigonometric expressions. We'll also encounter nonlinear expansions, such as sigmoid transformation common to neural network models,

$$h_k(x) = \frac{1}{1 + \exp(-x^T \beta_k)} \quad (2.82)$$

Maximum Likelihood Estimation

Suppose that we had N random variables associated with N measurements that we denote by X_1, \dots, X_N . Then, we denote their joint density, parameterized by θ by

$$f_{\theta}(x_1, \dots, x_N) = f(x_1, \dots, x_N | \theta), \quad (2.83)$$

where we have observed values $X_1 = x_1, \dots, X_N = x_N$. We define the likelihood of θ as the function

$$lik(\theta) = f(x_1, \dots, x_N | \theta). \quad (2.84)$$

If the distribution is discrete, then f is the frequency distribution function. In essence, $lik(\theta)$ measures the probability that the observed data was generated by a distribution parameterized by θ . Hence, the principle of *maximum likelihood estimation* argues that one should fit the model by the θ that maximizes $lik(\theta)$. Formally, we want

$$\hat{\theta} = \underset{\theta \in S}{\operatorname{argmax}} \, lik(\theta), \quad (2.85)$$

where S denotes the set of all possible values for θ .

Definition 2.18: Residual Sum of Squares, RSS

Suppose that we have a set of N data points (i.e pairing of observed input-outputs (x, y)). This set would be given by $\{(x_i, y_i) | i \in \mathbb{Z}_N\}$. Then, we can estimate the parameters θ in f_θ as we did for the linear model via minimizing residual sum-of-squares:

$$RSS(\theta) = \sum_{i=1}^N (y_i - f_\theta(x_i))^2 \quad (2.86)$$

Example 2.2: Max Likelihood Estimation: Least Squares

Suppose that we have a random sample $y_i, i = 1, \dots, N$ from a density $Pr_\theta(y)$ indexed by some parameters θ . Then, the likelihood function is given by

$$lik(\theta) = Pr(y_1, \dots, y_N | \theta) = \prod_{i=1}^N Pr(y_i; \theta), \quad (2.87)$$

where the second equality has assumed that y_i 's are independent. We note that although likelihood is often stated in a conditional formalism $Pr(Y|\theta)$, θ is not a random variable but an unknown parameter. Hence, we will typically write $Pr(Y; \theta)$. Then, the log-probability of the observed sample is given by

$$L(\theta) = \sum_{i=1}^N \log[Pr(y_i; \theta)]. \quad (2.88)$$

Since \log is a monotonically increasing function, it preserves the solution to MLE:

$$\hat{\theta} = \underset{\theta \in S}{\operatorname{argmax}} \prod_{i=1}^N Pr(y_i; \theta) = \underset{\theta \in S}{\operatorname{argmax}} \sum_{i=1}^N \log[Pr(y_i; \theta)]. \quad (2.89)$$

Then, suppose that we consider the additive error model $Y = f_\theta(X) + \epsilon$, with $\epsilon \sim N(0, \sigma^2)$. Then, suppose that the conditional likelihood was Gaussian

$$Pr(Y|X, \theta) = N(f_\theta(X), \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(Y - f_\theta(X))^2}{2\sigma^2}} \quad (2.90)$$

We therefore have

$$L(\theta) = \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - f_\theta(x_i))^2}{2\sigma^2}} \right] \quad (2.91)$$

$$= -\frac{N}{2} \log(2\pi) - N \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_\theta(x_i))^2. \quad (2.92)$$

Since only the last term contains θ , by MLE we have that

$$\hat{\theta} = \underset{\theta \in S}{\operatorname{argmax}} \left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_\theta(x_i))^2 \right) = \underset{\theta \in S}{\operatorname{argmin}} \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_\theta(x_i))^2 = \underset{\theta \in S}{\operatorname{argmin}} RSS(\theta), \quad (2.93)$$

which tells us that our solution to MLE is found by minimizing $RSS(\theta)$. We note that this is typically how loss functions can be derived. We choose some reasonable distribution that may have generated the data, invoke the principle of MLE and derive the function that requires minimization (i.e a loss function).

Example 2.3: Multinomial Likelihood

Consider the multinomial likelihood for the regression function $Pr(G|X)$ for a qualitative output G . Suppose that we have a model $Pr(G = \mathcal{G}_k|X = x) = p_{k,\theta}(x)$, $k = 1, \dots, K$ for the conditional probability of each class given X , indexed by parameter vector θ . Then, the log-likelihood (also referred to as cross-entropy) is

$$L(\theta) = \sum_{i=1}^N \log(p_{g_i, \theta}(x_i)), \quad (2.94)$$

and when maximized it delivers values of θ that best conforms with the data in the likelihood sense.

2.6 Structured Regression Models

We've seen that although nearest-neighbour and other local methods focus directly on estimating the function at a point, they face problems in high dimensions. However, they may also be inappropriate in low dimensions where more structured approaches can make more efficient use of the data. We'll introduce classes of such structured approaches in this section.

2.6.1 Difficulty of the Problem

Consider the RSS criterion for an arbitrary function f ,

$$RSS(f) = \sum_{i=1}^N (y_i - f(x_i))^2 \quad (2.95)$$

Note that minimizing (2.95) leads to infinitely many solutions, as any function \hat{f} passing through the training points (x_i, y_i) is a solution. One would also run into the problem of overfitting as any particular solution chosen might be a poor predictor at test points different from the training points. If there are multiple observation pairs (x_i, y_{il}) , $l = 1, \dots, N_i$ at each value of x_i , the risk becomes limited as the solutions would pass through average values of the y_{il} at each x_i .

In order to obtain useful results for finite N , we must restrict the eligible solutions to (2.95) to a smaller set of functions. Any restrictions imposed on f that lead to a unique solution to (2.95) does not really remove the ambiguity caused by the vast space of solutions. There are infinitely many possible restrictions, each leading to the unique solution.

In general, the imposed constraints by most learning methods can be described as *complexity* restrictions of some kind. Usually, this means some kind of regular behaviour in small neighbourhoods of the input space. That is, for all input points x sufficiently close to each other in some metric, \hat{f} exhibits some special structure such as nearly constant, linear or low-order polynomial behaviour. The estimator can then be obtained by averaging or polynomial fitting in that neighbourhood.

Methods such as splines, neural networks and basis-function methods implicitly define neighbourhoods of local behaviour. Any method that attempts to produce locally varying functions in small isotropic neighbourhoods will run into problems in high dimensions - curse of dimensionality. In addition, all method that overcome dimensionality problems have an associated - and often implicit or adaptive- metric for measuring neighbourhoods, which basically does not allow the neighbourhood to be simultaneously small in all directions.

2.7 Classes of Restricted Estimators

Here we will give a brief summary, since detailed descriptions are given in later chapters. Each of the classes has associated with it one or more parameters, sometimes appropriately called *smoothing* parameters.

2.7.1 Roughness Penalty and Bayesian Methods

Here, the class of functions is controlled by explicitly penalizing $RSS(f)$ with a roughness penalty

$$PRSS(f; \lambda) = RSS(f) + \lambda J(f). \quad (2.96)$$

The user-selected functional $J(f)$ will be large for functions f that vary too rapidly over small regions of input space.

Example 2.4: Cubic Smoothing Spline

For example, the popular *cubic smoothing spline* for one-dimensional inputs is the solution to the penalized least-squares criterion

$$PRSS(f; \lambda) = \sum_{i=1}^N (y_i - f(x_i))^2 + \lambda \int [f''(x)]^2 dx \quad (2.97)$$

The roughness penalty here controls large values of the second derivative of f , and the amount of penalty is dictated by $\lambda \geq 0$. For $\lambda = 0$ no penalty is imposed, and any interpolating function will do, while for $\lambda = \infty$ only functions linear in x are permitted. We can observe this fact by recognizing that the second term must go to zero if $\lambda \rightarrow \infty$. For this to occur, we require that $f''(x) = 0$, which are the class of functions $f(x) = \alpha x + b$ for some $\alpha, b \in \mathbb{R}$.

Example 2.5: Additive Functionals

Penalty functionals J can be constructed for functions in any dimension, and special versions can be created to impose special structure. For example, additive penalties $J(f) = \sum_{j=1}^p J(f_j)$ are used in conjunction with additive functions $f(X) = \sum_{j=1}^p f_j(X_j)$ to create additive models with smooth coordinate functions.

Example 2.6: Projection Pursuit Regression

Projection pursuit regression models have $f(X) = \sum_{m=1}^M g_m(\alpha_m^T X)$ for adaptively chosen directions α_m , and the functions g_m can each have an associated roughness penalty.

A penalty function, or *regularization methods*, express our prior belief that type of functions we seek exhibit a certain type of smooth behaviour, and can be cast in a Bayesian framework. The penalty J corresponds to a log-prior, and $PRSS(f; \lambda)$ the log-posterior distribution, and minimizing $PRSS(f; \lambda)$ amounts to finding the posterior mode.

2.7.2 Kernel Methods and Local Regression

These methods aim to explicitly provide estimates of the regression function or conditional expectation by specifying the nature of the local neighbourhood, and of class of regular functions fitted locally.

Definition 2.19: Kernel Density Estimation

In statistics, *kernel density estimation (KDE)* is a non-parametric way to estimate the probability density function of a random variable. Kernel density estimation is a fundamental data smoothing problem where inferences about the population are made, based on a finite data sample.

Definition 2.20: Kernel Function

The local neighbourhood is specified by a kernel function $K_\lambda(x_0, x)$ which assigns weights to points x in a region around x_0 . For example, the Gaussian kernel has a weight function based on the Gaussian density function

$$K_\lambda(x_0, x) = \frac{1}{\lambda} \exp\left[-\frac{\|x - x_0\|^2}{2\lambda}\right], \quad (2.98)$$

and assigns weights to points that die exponentially with their squared Euclidean distance from x_0 . The parameter λ corresponds to the variance of Gaussian density, and controls the width of the neighbourhood.

Definition 2.21: Nadaraya-Watson Weighted Average

The simplest form of kernel estimate is the Nadaraya-Watson weighted average

$$\hat{f}(x_0) = \frac{\sum_{i=1}^N K_\lambda(x_0, x_i) y_i}{\sum_{i=1}^N K_\lambda(x_0, x_i)} \quad (2.99)$$

Definition 2.22: Local Regression Estimate

In general, we can define a local regression estimate of $f(x_0)$ as $\hat{f}_\theta(x_0)$, where $\hat{\theta}$ minimizes

$$RSS(f_\theta, x_0) = \sum_{i=1}^N K_\lambda(x_0, x_i) (y_i - f_\theta(x_i))^2, \quad (2.100)$$

and f_θ is some parameterized functions, such as a low-order polynomials.

Example 2.7: Parameterized Functions for Local Regression Estimate

Some examples of parameterized functions for local regression estimate are

- $f_\theta(x) = \theta_0$, the constant function. This results in the Nadaraya-Watson estimate mentioned above in (2.99).
- $f_\theta(x) = \theta_0 + \theta_1 x$ gives the popular local linear regression model.

Nearest-neighbour methods can be thought of as kernel methods having a more data-dependent metric. Indeed, the metric for k -nearest neighbours is

$$K_k(x, x_0) = I(\|x - x_0\| \leq \|x_{(k)} - x_0\|), \quad (2.101)$$

where $x_{(k)}$ is the training observation ranked k^{th} in distance from x_0 , and $I(S)$ is the indicator of the set S .

2.7.3 Basis Functions and Dictionary Methods**Definition 2.23: Basis Functions**

This class of methods includes the familiar linear and polynomial expansions, but more importantly a wide

variety of more flexible models. The model for f is a linear expansion of basis functions

$$f_{\theta}(x) = \sum_{m=1}^M \theta_m h_m(x), \quad (2.102)$$

where each of the h_m is a function of the input x , and term linear here refers to the action of the parameters θ . This class covers wide variety of methods.

Definition 2.24: Radial Basis Functions

Radial basis functions are symmetric p -dimensional kernels located at particular centroids,

$$f_{\theta}(x) = \sum_{m=1}^M K_{\lambda_m}(\mu_m, x) \theta_m; \quad (2.103)$$

for instance, the Gaussian kernel $K_{\lambda}(\mu, x) = e^{-||x-\mu||^2/2\lambda}$ is popular. Radial basis functions have centroids μ_m and scales λ_m that have to be determined.

Definition 2.25: Dictionary Methods

A single-layer feed-forward neural network model with linear output weights can be thought of as an adaptive basis function method. The model has the form

$$f_{\theta}(x) = \sum_{m=1}^M \beta_m \sigma(\alpha_m^T x + b_m), \quad (2.104)$$

where $\sigma(x) = 1/(1 + e^{-x})$ is known as the activation function. The directions α_m and bias terms b_m have to be determined. These adaptively chosen basis function methods are also known as dictionary methods.

2.8 Model Selection and the Bias-Variance Tradeoff

All models described above and many others that will be discussed have a *smoothing* or *complexity* parameter that has to be determined:

- the multiplier of the penalty term;
- the width of the kernel;
- or the number of basis functions.

Proposition 2.11: Test / Generalization Error

Let $\hat{f}_k(x_0)$ denote the k -nearest neighbour regression fit. The consideration of the nearest neighbours usefully illustrates the competing forces that affect the predictive ability of such approximations. Suppose that data arises from a model $Y = f(X) + \epsilon$, with $E(\epsilon) = 0$ and $Var(\epsilon) = \sigma^2$. For simplicity, we assume that values of x_i in sample are fixed in advance (non-random). The expected prediction error at x_0 , also known as test

or generalization error, can be decomposed:

$$EPE_k(x_0) = E[(Y - \hat{f}_k(x_0))^2 | X = x_0] = \sigma^2 + \left[f(x_0) - \frac{1}{k} \sum_{l=1}^k f(x_{(l)}) \right]^2 + \frac{\sigma^2}{k} \quad (2.105)$$

Proof. We'll compute this by way of the property of conditional expectations [See Proposition A.24]

$$\begin{aligned} E[(Y - \hat{f}_k(x_0))^2 | X = x_0] &= E[(\epsilon + f(X) - \hat{f}_k(x_0))^2 | X = x_0] \\ &= E[\epsilon^2 + 2\epsilon(f(X) - \hat{f}_k(x_0)) + (f(X) - \hat{f}_k(x_0))^2 | X = x_0] \\ &= E[\epsilon^2 | X = x_0] + 2E[\epsilon(f(X) - \hat{f}_k(x_0)) | X = x_0] + E[(f(X) - \hat{f}_k(x_0))^2 | X = x_0] \\ &= E[\epsilon^2] + 2E[\epsilon]E[(f(X) - \hat{f}_k(x_0)) | X = x_0] + (f(x_0) - \hat{f}_k(x_0))^2 \\ &= \sigma^2 + [Bias^2(\hat{f}_k(x_0)) + Var_{\mathcal{T}}(\hat{f}_k(x_0))] \\ &= \sigma^2 + \left(f(x_0) - \frac{1}{k} \sum_{l=1}^k f(x_{(l)}) \right)^2 + \frac{\sigma^2}{k}, \end{aligned} \quad (2.106)$$

where we refer to the first term σ^2 as the *irreducible error*. This is the variance of the new test target-and is beyond our control, even if we know the true $f(x_0)$. The bias term is the squared difference between the true mean $f(x_0)$ and expected value of the estimate $[E_{\mathcal{T}}(\hat{f}_k(x_0)) - f(x_0)]^2$, where expectation averages randomness in training data. We expect that the second term will increase as k increases. We note that choosing larger number of nearest neighbours will influence the average approximation, which should begin to further deviate from $f(x_0)$ as k increases. \square

Lemma 2.2: Nearest Neighbours Bias

We consider the case under which the Bias increases as the number of nearest neighbours increases. We fix a point $x_0 \in \mathbb{R}^p$. Let (x_i, y_i) denote an input-output pair from a training sample of size m . We define $x_{(i)}$ ($1 \leq i \leq m$) as the ordered set of points that are closest to x_0 . Let N_k denote the k -nearest neighbours estimate

$$N_k = \frac{1}{k} \sum_{i=1}^k f(x_{(i)}), \quad (2.107)$$

where $f(x_{(i)}) = y_i$. The Bias term

$$Bias(k) = f(x_0) - N_k \quad (2.108)$$

increases as k increases if

$$N_k > \frac{1}{n-k} \sum_{i=k+1}^n f(x_{(i)}) \quad (2.109)$$

where $n > k$. The term on the right hand side of the inequality measures the average y_i value of the $k+1^{th}$ to n^{th} farthest away neighbour from x_0 . Hence, if the average of farther away neighbours tends to be less than k -nearest neighbours, the Bias will increase.

Proof. We want to observe the conditions under which

$$f(x_0) - N_k < f(x_0) - N_n \quad (2.110)$$

From this inequality, we have that

$$N_k > N_n \quad (2.111)$$

We can relate N_n to N_k via

$$N_n = \frac{k}{n}N_k + \frac{1}{n} \sum_{i=k+1}^n f(x_{(i)}). \quad (2.112)$$

Then, we have that

$$\frac{k}{n}N_k + \frac{1}{n} \sum_{i=k+1}^n f(x_{(i)}) < N_k \quad (2.113)$$

$$\rightarrow N_k > \frac{1}{n-k} \sum_{i=k+1}^n f(x_{(i)}) \quad (2.114)$$

□

Bias-Variance Tradeoff

The bias-variance tradeoff is a central problem in supervised learning. Ideally, one wants to choose a model that both accurately captures the regularities in its training data, but also generalizes well to unseen data. Unfortunately, it is typically impossible to do both simultaneously. High-variance learning methods may be able to represent their training set well but are at risk of overfitting to noisy or unrepresentative training data. In contrast, algorithms with high bias typically produce simpler models that don't tend to overfit but may underfit their training data, failing to capture important regularities. The variance term in (2.106) decreases as the inverse of k . Hence, as k varies, there is a bias-variance tradeoff.

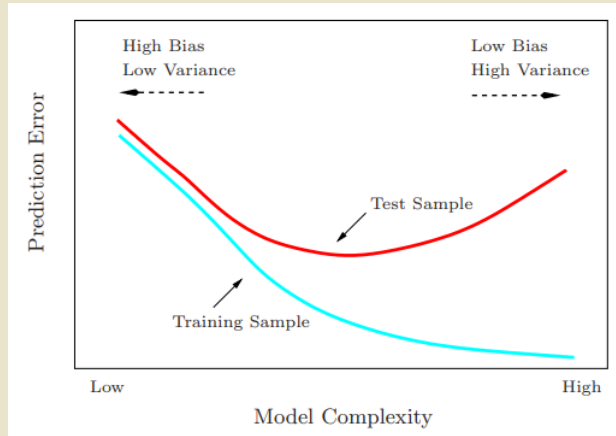


Figure 2.1: Test and training error as a function of model complexity.

Typically, we would like to choose our model complexity to trade bias off with variance in such a way as to minimize the test error. An obvious estimate of test error is the *training error* $\frac{1}{N} \sum_i (y_i - \hat{y}_i)^2$. However, training error is not a good estimate of test error, as it does not properly account for model complexity. Figure 2.1 shows typical behaviour of test and training error, as model complexity is varied. If we perform too much fitting, the model adapts itself too closely to the training data, and will not generalize well.

Chapter 3

Linear Methods for Regression

3.1 Introduction

A linear regression model assumes that the regression function $E(Y|X)$ is linear in the inputs X_1, \dots, X_p . For prediction purposes, linear models can sometimes outperform fancier nonlinear models, especially in situations with small numbers of training cases, low signal-to-noise ratio or sparse data. In this chapter, we describe linear methods for regression. The authors believe that an understanding of linear methods is essential for understanding nonlinear ones. Many nonlinear techniques can be considered direct generalizations of the linear methods that we'll discuss here.

3.2 Linear Regression Models and Least Squares

Linear Model

Just as in §2, we have an input vector $X^T = (X_1, X_2, \dots, X_p)$, and want to predict a real-valued output Y . The linear regression model has the form

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j \quad (3.1)$$

The linear model either assumes that the regression function $E(Y|X)$ is linear, or that the linear model is a reasonable approximation. Here the β_j 's are unknown parameters or coefficients, and the variables X_j can come from a variety of sources:

- quantitative inputs;
- transformations of quantitative inputs, such as log, square-root or square;
- basis expansions, such as $X_2 = X_1^2$, $X_3 = X_1^3$, leading to polynomial representation;
- interactions between variables, such as $X_3 = X_1 \cdot X_2$

Proposition 3.1

Let X be a $n \times m$ matrix. If X has full column rank, then $X^T X$ is positive definite.

Proof. Let x_1, \dots, x_n denote the column vectors of X . Then, we can express X in terms of these column

vectors as

$$X = \begin{bmatrix} | & | & \dots & | \\ | & | & \dots & | \\ x_1 & x_2 & \dots & x_n \\ | & | & \dots & | \\ | & | & \dots & | \end{bmatrix}, \quad (3.2)$$

If X has full column rank, then all of its columns are linearly independent. Let $v_i \in \mathbb{R} \forall i$, then the only solution to

$$\sum_{i=1}^n x_i v_i = \mathbf{0} \quad (3.3)$$

is if $v_i = 0 \forall i$. Let $v \in \mathbb{R}^n$. As a consequence of this, $Xv = \mathbf{0}$ iff $v = \mathbf{0}$. We will now show this. We can observe that

$$Xv = \sum_{i=1}^n x_i v_i, \quad (3.4)$$

which by its full column rank property means that $Xv \neq 0 \forall v \in \mathbb{R}^n \setminus \{\mathbf{0}\}$. Let $y = Xv$. Then, we consequently have

$$y^T y = v^T X^T X v = \sum_{i=1}^m y_i^2 > 0, \quad (3.5)$$

which by the definition of positive definite matrices, means that $X^T X$ is positive definite. \square

Corollary 3.1

Let A be a $n \times n$ matrix. If A is positive definite, then A is invertible.

Proof. Since A is positive definite, then $v^T A v > 0 \forall v \in \mathbb{R}^n \setminus \{\mathbf{0}\}$. Hence, $Av \neq \mathbf{0} \forall v \in \mathbb{R}^n \setminus \{\mathbf{0}\}$. Consequently, A must have full column rank. Since A has full column rank, A is invertible. \square

Least Squares Estimation: Linear Regression

Suppose that we have a set of training data $(x_1, y_1), \dots, (x_N, y_N)$ from which to estimate the parameters β . Each $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$ is a vector of feature measurements from the i^{th} case. Most popular estimation method is least squares, in which we pick coefficients $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$ to minimize residual sum of squares

$$RSS(\beta) = \sum_{i=1}^N (y_i - f(x_i))^2 = \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2. \quad (3.6)$$

From a statistical point of view the criterion is reasonable if the y_i 's are conditionally independent given the inputs x_i . We define \mathbf{X} as the $N \times (p+1)$ matrix with each row being an input vector. Similarly, let \mathbf{y} be the N -vector of outputs in the training set. Then, we can write the residual sum-of-squares as

$$RSS(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta). \quad (3.7)$$

Let $Y := \mathbf{y} - \mathbf{X}\beta$. Differentiating w.r.t β gets us

$$\frac{\partial RSS}{\partial \beta} = \frac{\partial Y^T}{\partial \beta} (\mathbf{y} - \mathbf{X}\beta) + (\mathbf{y} - \mathbf{X}\beta)^T \frac{\partial Y}{\partial \beta} \quad (3.8)$$

$$= -\mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) + (\mathbf{y} - \mathbf{X}\beta)^T \mathbf{X} \quad (3.9)$$

$$= -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) \quad (3.10)$$

$$\frac{\partial^2 RSS}{\partial \beta \partial \beta^T} = 2\mathbf{X}^T \mathbf{X} \quad (3.11)$$

Suppose that \mathbf{X} has full column rank, then by Proposition 3.1 and Corollary 3.1, $\mathbf{X}^T \mathbf{X}$ is positive definite and is invertible. Setting first derivative to zero:

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) = 0, \quad (3.12)$$

has unique solution

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.13)$$

The fitted values at the training inputs are

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (3.14)$$

where $\hat{y}_i = \hat{f}(x_i)$. The matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ appearing in equation (3.14) is sometimes referred to as the *hat matrix* because it puts the hat on \mathbf{y} .

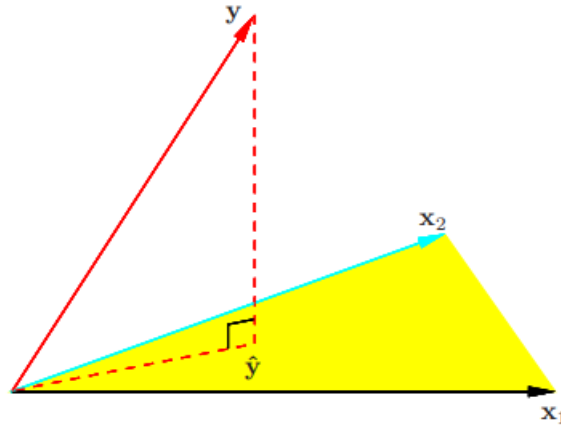


Figure 3.1: The N-dimensional geometry of least squares regression with two predictors. The outcome vector \mathbf{y} is orthogonally projected onto the hyperplane spanned by the input vectors \mathbf{x}_1 and \mathbf{x}_2 . The projection $\hat{\mathbf{y}}$ represents the vector of the least squares predictions.

Proposition 3.2: Geometrical Interpretation of Least Squares

Let $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_p$ be the column vectors of \mathbf{X} with $\mathbf{x}_0 \equiv 1$. Then, minimizing least squares amounts to choosing $\hat{\beta}$ so that the residual vector $\mathbf{y} - \hat{\mathbf{y}}$ is orthogonal to the subspace spanned by the column vectors $\{\mathbf{x}_i\}$ and the resulting estimate $\hat{\mathbf{y}}$ is the orthogonal projection of \mathbf{y} onto this subspace. Since the hat matrix \mathbf{H} computes the orthogonal projection, it is also known as a projection matrix.

Proof. We did most of the hard work in the “Least Squares Estimation” blurb. Consider (3.15) which we rewrite as

$$\mathbf{X}^T (\mathbf{y} - \hat{\mathbf{y}}) = 0. \quad (3.15)$$

Letting \mathbf{x}_i denote the column vectors of \mathbf{X} and $(\mathbf{x}_i)_j$ the j^{th} component of \mathbf{x}_i , we have that

$$[\mathbf{X}^T(\mathbf{y} - \hat{\mathbf{y}})]_k = \sum_{j=1}^N [\mathbf{X}^T]_{kj} (\mathbf{y}_j - \hat{\mathbf{y}}_j) \quad (3.16)$$

$$= \sum_{j=1}^N \mathbf{X}_{jk} (\mathbf{y}_j - \hat{\mathbf{y}}_j) \quad (3.17)$$

$$= \sum_{j=1}^N (\mathbf{x}_k)_j (\mathbf{y}_j - \hat{\mathbf{y}}_j) \quad (3.18)$$

$$= \mathbf{x}_k \cdot (\mathbf{y} - \hat{\mathbf{y}}) \quad (3.19)$$

$$= 0 \quad \forall k \quad (3.20)$$

Hence, we have that $(\mathbf{y} - \hat{\mathbf{y}})$ is orthogonal to the subspace spanned by the column vectors \mathbf{x}_k . More importantly, we see that $\hat{\mathbf{y}}$ is the orthogonal projection of \mathbf{y} onto this space by computing the inner product on \mathbb{R}^N :

$$\hat{\mathbf{y}}^T (\mathbf{y} - \hat{\mathbf{y}}) = \hat{\mathbf{y}}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X})^{-1})^T \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (3.21)$$

$$= \hat{\mathbf{y}}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \underbrace{\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}}_{\mathbf{I}} \mathbf{X}^T \mathbf{y} \quad (3.22)$$

$$= \hat{\mathbf{y}}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (3.23)$$

$$= \hat{\mathbf{y}}^T \mathbf{y} - \hat{\mathbf{y}}^T \mathbf{y} \quad (3.24)$$

$$= 0 \quad (3.25)$$

Hence, since $\hat{\mathbf{y}} \perp \mathbf{y} - \hat{\mathbf{y}}$ and $\mathbf{y} - \hat{\mathbf{y}} \perp \mathbf{v} \quad \forall \mathbf{v} \in \text{span}(\mathbf{x}_0, \dots, \mathbf{x}_p)$, we can conclude that $\hat{\mathbf{y}} \in \text{span}(\mathbf{x}_0, \dots, \mathbf{x}_p)$. \square

It can certainly be the case that the columns of \mathbf{X} are not linearly independent, so that \mathbf{X} is not of full rank. This can occur if for instance, two inputs were perfectly correlated (e.g. $\mathbf{x}_2 = 3\mathbf{x}_3$). This turns $\mathbf{X}^T \mathbf{X}$ singular and the least squares coefficients $\hat{\beta}$ are no longer unique. However, $\hat{\mathbf{y}} = \mathbf{X} \hat{\beta}$ is still a projection of \mathbf{y} onto the column space of \mathbf{X} , but there is more than one way of expressing that projection in terms of the column vectors of \mathbf{X} . Most regression software packages detect these redundancies and automatically implement some strategy for removing them.

Definition 3.1: Covariance Matrix of Random Vector

The covariance matrix of a random vector X is defined as

$$\text{Cov}(X) := E[(X - E[X])(X - E[X])^T], \quad (3.26)$$

which is sometimes denoted by $\text{Var}(X)$.

Proposition 3.3: Covariance Matrix of $\hat{\beta}$

Let $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, where y is the output vector, whose observations are assumed to be uncorrelated with constant variance σ^2 (i.e the covariance matrix has entries $\text{Cov}[\mathbf{y}]_{ij} = \sigma^2 \delta_{ij}$). Then, the covariance matrix of $\hat{\beta}$ is given by

$$\text{Cov}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2 \quad (3.27)$$

Definition 3.2: Chi-Squared (χ^2)-Distribution

Let Z_1, Z_2, \dots, Z_k be k independent standard-normal random variables. Then, the variable $Q = \sum_{i=1}^k Z_i^2$ is distributed according to the χ -squared distribution for k degrees of freedom. This is often denoted as

$$Q \sim \chi^2(k) \quad \text{or} \quad Q \sim \chi_k^2. \quad (3.28)$$

The probability density function of the χ^2 -distribution is given by

$$f(x; k) = \begin{cases} \frac{x^{k/2-1} e^{-x/2}}{2^{k/2} \Gamma(k/2)}, & x > 0 \\ 0, & \text{otherwise} \end{cases} \quad (3.29)$$

The expectation of the χ_k^2 distribution is given by

$$E[\chi_k^2] = k \quad (3.30)$$

Proposition 3.4: Unbiased Variance Estimator

Consider a training set of n observations $\{(\mathbf{x}_i, \mathbf{y}_i)\}$ where $\mathbf{x}_i \in \mathbb{R}^p$. We construct the $N \times p$ matrix \mathbf{X} with \mathbf{x}_i forming the i^{th} row. We consider the linear model $y = \mathbf{X}\beta + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$. We define the estimator

$$\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (3.31)$$

where $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$ and $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. Then, we have that

$$RSS \sim \sigma^2 \chi_{n-p}^2. \quad (3.32)$$

In particular, our estimator is unbiased as

$$E[\hat{\sigma}^2] = \sigma^2. \quad (3.33)$$

Proof. We take the linear model $\mathbf{y} = \mathbf{X}\beta + \epsilon$ and consider the vector of residuals $\hat{\epsilon}$

$$\hat{\epsilon} = \mathbf{y} - \mathbf{X}\hat{\beta} \quad (3.34)$$

$$= (\mathbb{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \mathbf{y} \quad (3.35)$$

$$:= Q\mathbf{y} \quad (3.36)$$

$$= Q(\mathbf{X}\beta + \epsilon) \quad (3.37)$$

$$= Q\epsilon, \quad (3.38)$$

where we have defined $Q := \mathbb{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$. Observe the following:

$$\text{tr}(Q) = \text{tr}(\mathbb{I}_{n \times n}) - \text{tr}(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \quad (3.39)$$

$$= n - \text{tr}(\mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1}) \quad (3.40)$$

$$= n - \text{tr}(\mathbb{I}_{p \times p}) \quad (3.41)$$

$$= n - p \quad (3.42)$$

However, we also observe the fact that Q is a normal matrix as $Q^2 = Q$. Since Q is idempotent, its only eigenvalues are either 0 or 1 [See B.3]. However, by the trace property of matrices, we have that

$$\text{tr}(Q) = \alpha_0(0) + \alpha_1(1) = n - p \quad (3.43)$$

where α_0, α_1 are the associated multiplicities for the 0 and 1 eigenvalues satisfying $\alpha_0 + \alpha_1 = n$. Hence, we must have that

$$\alpha_1 = n - p, \quad \alpha_0 = p. \quad (3.44)$$

We now use the fact that since Q is normal, then there exists a unitary matrix V that can diagonalize Q :

$$\Delta_{n-p}^p := V^T Q V = \text{diag}(\underbrace{1, 1, \dots, 1}_{n-p \text{ times}}, \underbrace{0, 0, \dots, 0}_{p \text{ times}}) \quad (3.45)$$

For convenience, we define an intermediate variable $K = V^T \hat{\epsilon}$. Then, we observe that

$$RSS := \sum_{i=1}^n \hat{\epsilon}_i^2 \quad (3.46)$$

$$= \hat{\epsilon}^T \hat{\epsilon} \quad (3.47)$$

$$= \hat{\epsilon}^T \underbrace{V V^T}_{\mathbb{I}} \hat{\epsilon} \quad (3.48)$$

$$= K^T K. \quad (3.49)$$

We also note that this can be cast in terms of ϵ :

$$K = V^T \hat{\epsilon} = V^T Q \epsilon \quad (3.50)$$

$$= V^T Q \underbrace{V V^T}_{\mathbb{I}} \epsilon \quad (3.51)$$

$$= \Delta_{n-p}^p V^T \epsilon \quad (3.52)$$

Hence,

$$RSS = K^T K = \epsilon^T V \Delta_{n-p}^p \Delta_{n-p}^p V^T \epsilon \quad (3.53)$$

$$= \epsilon^T V V^T \Delta_{n-p}^p \epsilon \quad (3.54)$$

$$= \epsilon^T \Delta_{n-p}^p \epsilon \quad (3.55)$$

$$= \sum_{i=1}^{n-p} \epsilon_i^2 \sim \sigma^2 \chi_{n-p}^2, \quad (3.56)$$

where in the second equality we have used the fact that any matrix A commutes with Δ_{n-p}^p : $\Delta_{n-p}^p A = A \Delta_{n-p}^p$. In the third equality, we used the fact that Δ_{n-p}^p is idempotent. The final equality is obtained from Δ_{n-p}^p 's property to select out the first $n - p$ entries and that $\epsilon_i \sim \mathcal{N}(0, \sigma^2) \forall i$. It therefore now follows that

$$E[RSS] = \sigma^2(n - p), \quad (3.57)$$

leading to the desired result

$$E\left[\frac{RSS}{n - p}\right] = E[\hat{\sigma}^2] = \sigma^2 \quad (3.58)$$

□

Definition 3.3: t-Distribution

Consider x_1, \dots, x_n as the n values observed in a sample from a continuously distributed population with expected value μ . Let \bar{x} and s denote the sample mean and sample variance respectively (s being the

Bessel-corrected sample mean with leading factor $1/(n-1)$). Then, the t -value is

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \quad (3.59)$$

The t -distribution with $n-1$ degrees of freedom is the sampling distribution of the t -value when the samples consist of independent identically distributed observations from a normally distributed population. The probability density function is given by

$$f(t; \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\nu/2)} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}, \quad (3.60)$$

where ν is the number of degrees of freedom.

Definition 3.4: Z-score Hypothesis Testing

To test the hypothesis that a particular coefficient $\beta_j = 0$, we form the standardized coefficient or Z -score

$$z_j = \frac{\hat{\beta}_j}{\hat{\sigma}\sqrt{v_j}} \quad (3.61)$$

where v_j is the j^{th} diagonal element of $(\mathbf{X}^T \mathbf{X})^{-1}$. Under the null hypothesis that $\beta_j = 0$, z_j is distributed as t_{N-p-1} (a t distribution with $N-p-1$ degrees of freedom), and hence a large (absolute) value of z_j will lead to rejection of this null hypothesis.

3.3 Subset Selection

There are two main reasons why we are usually not satisfied with the least squares estimates given by Equation (3.13):

1. **Prediction Accuracy:** The least squares estimates often have low bias but large variance. Prediction accuracy can sometimes be improved by shrinking or setting some coefficients to zero. In doing so, we sacrifice a little bit of bias to reduce the variance of the predicted values and may thus improve the overall prediction accuracy.
2. **Interpretation:** With a large number of predictors, we often would like to determine a smaller subset that exhibits the strongest effects. In order to get the “big picture”, we are willing to sacrifice some of the small details.

3.3.1 Best-Subset Selection

Best subset regression finds for each $k \in \{0, 1, 2, \dots, p\}$ the subset of size k that gives smallest residual sum of squares. An efficient algorithm- the *leaps and bounds* procedure (Furnival and Wilson, 1974)- makes this feasible for p as large as 30 or 40. Note that the best subset of size 2, for example, need not include the variable that was in the best subset of size 1.

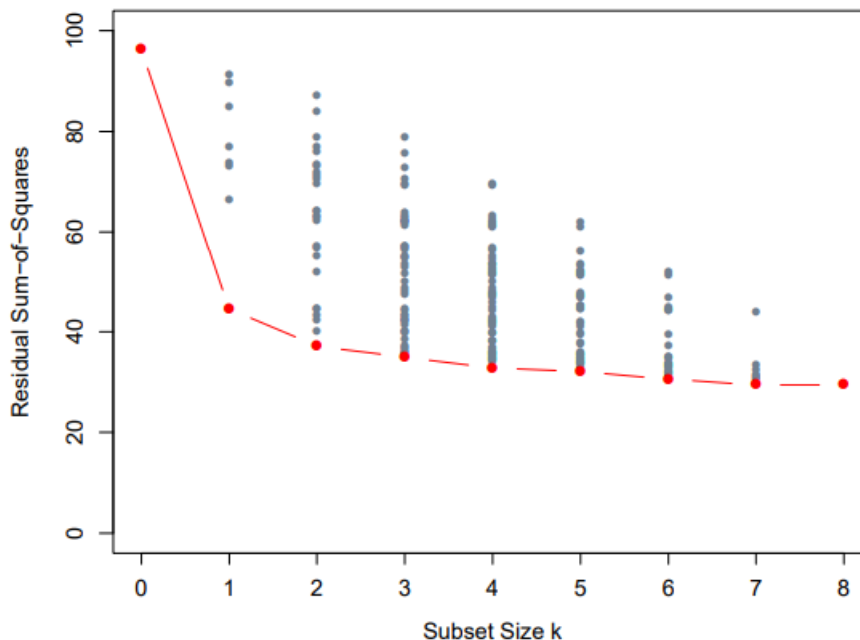


Figure 3.2: All possible subset models for the prostate cancer example. At each subset size is shown the residual sum-of-squares for each model of that size.

In the next section, we use cross-validation to estimate the prediction error and select k ; the AIC criterion is a popular alternative. We defer more detailed discussion of these and other approaches to §7.

3.3.2 Forward- and Backward-Stepwise Selection

Rather than searching through all possible subsets (which can become infeasible for p much larger than 40), we can seek a good path through them.

Forward-Stepwise Selection

Forward-stepwise selection starts with the intercept, and then sequentially adds into the model the predictor that most improves the fit. Though this may seem like a lot of computation, clever updating algorithms can exploit the QR decomposition for the current fit to rapidly establish the next candidate.

Forward-stepwise selection is a *greedy algorithm*, producing a nest sequence of models. While it may seem sub-optimal compared to best-subset selection, there are several reasons why it can be preferred:

1. *Computational*: For large p we cannot compute the best subset sequence, but we can always compute the forward stepwise sequence (even when $p \gg N$).
2. *Statistical*: A price is paid in variance for selecting the best subset of each size; forward stepwise is a more constrained search, and will have lower variance, but perhaps more bias.

Backward-Stepwise Selection

Backward-stepwise selection starts with the full model, and sequentially deletes the predictor that has the least impact on the fit. The candidate for dropping is the variable with the smallest Z -score. Backward selection can only be used when $N > p$, while forward stepwise can always be used.

3.3.3 Forward-Stagewise Regression

Forward-stagewise regression (FS) begins like forward-stepwise regression, with an intercept equal to \bar{y} , and centred predictors with coefficients initially all 0. At each step, the algorithm identifies the variable most

correlated with the current residual. It then computes the simple linear regression coefficient of the residual on this chosen variable, and then adds it to the current coefficient for that variable. This procedure continues until none of the variables have correlation with the residuals-i.e the least squares fit when $N > p$.

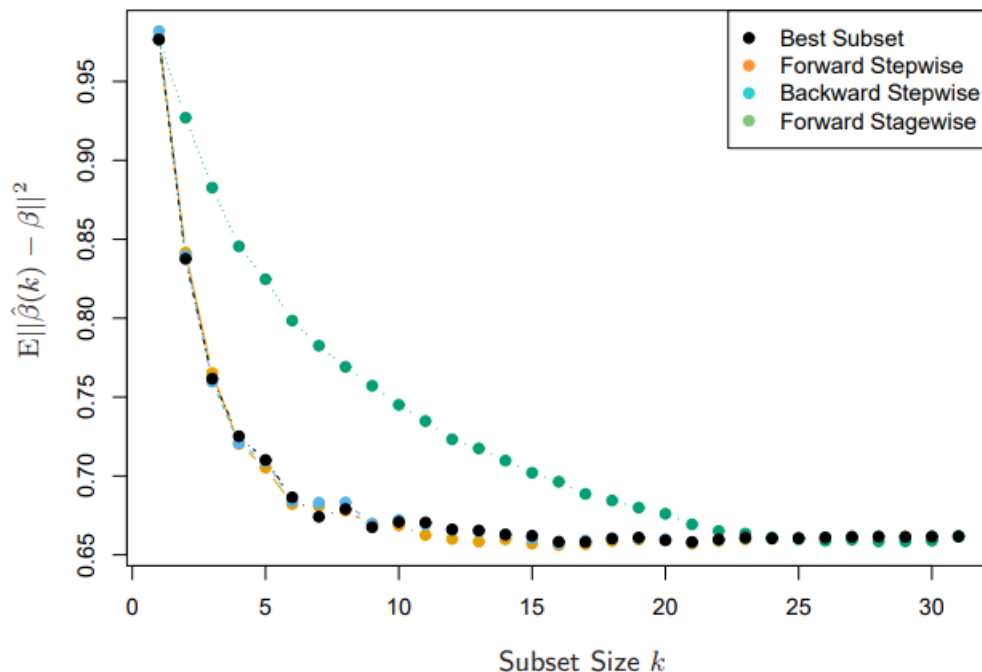


Figure 3.3: Comparison of four subset-selection techniques on a simulated linear regression problem $Y = X^T\beta + \epsilon$. There are $N = 300$ observations on $p = 31$ standard Gaussian variables, with pairwise correlations all equal to 0.85. For 10 of the variables, the coefficients are drawn at random from a $\mathcal{N}(0, 0.4)$ distribution; the rest are zero. The noise $\epsilon \sim \mathcal{N}(0, 6.25)$, resulting in a signal-to-noise ratio of 0.64. Results are averaged over 50 simulations. Shown is the mean-squared error of the estimated coefficient $\hat{\beta}(k)$ at each step from the true β .

3.4 Shrinkage Methods

3.4.1 Ridge Regression

Definition 3.5: Ridge Regression

Ridge regression shrinks the regression coefficients by imposing a penalty on their size. The ridge coefficients minimize a penalized residual sum of squares by the minimization problem

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}. \quad (3.62)$$

Here, $\lambda \geq 0$ is a complexity parameter that controls the amount of shrinkage: the larger the value of λ , the greater the amount of shrinkage. The coefficients are shrunk toward zero (and each other). The idea of penalizing by the sum-of-squares of the parameters is also used in neural networks, where it is known as

weight decay. An equivalent way to write the ridge problem is

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2. \quad (3.63)$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq t, \quad (3.64)$$

which makes explicit the size constraint on the parameters. There is a one-to-one correspondence between the parameters λ and t .

We notice that the intercept β_0 has been left out of the penalty term. Penalization of the intercept would make the procedure depend on the origin chosen for Y ; that is, adding a constant c to each of the targets y_i would not simply result in a shift of the predictions by the same amount c .

Proposition 3.5: Ridge Regression Solution

Consider the Ridge Regression problem, given by

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \quad (3.65)$$

$$= \underset{\beta}{\operatorname{argmin}} \{ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta \}. \quad (3.66)$$

We first standardize our inputs by the approximation $\beta_0 = \frac{1}{n} \sum_{i=1}^n y_i$ and $x_{ij} \mapsto x_{ij} - \bar{x}_j$. Then, the unique solution is given by

$$\hat{\beta}^{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.67)$$

We note that the addition of a positive constant λ to the diagonals of $\mathbf{X}^T \mathbf{X}$ ensures invertibility, guaranteeing that $(\mathbf{X}^T \mathbf{X} + \lambda \mathbb{I})^{-1}$ is well-defined [See Proposition 3.7].

Proposition 3.6: Ridge Regression Breaks Scaling Symmetry

Consider the regular RSS solution to Equation 3.6

$$\hat{\beta} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (3.68)$$

which is well-defined provided that $\mathbf{X}^T \mathbf{X}$ is invertible. Then, if we scale all the inputs $x_i \mapsto \alpha x_i$, we have the scaling $\mathbf{X} \mapsto \alpha \mathbf{X}$. Hence, our RSS solution scales as

$$\hat{\beta} \mapsto \frac{1}{\alpha} \hat{\beta}. \quad (3.69)$$

However, the presence of a non-zero λ term in ridge regression breaks this scaling symmetry.

Lemma 3.1: Inner Product Adjoint Property

Let A be an $n \times m$ real matrix. Then its matrix adjoint is simply A^T . It follows by the definition of an adjoint that

$$\langle Au, v \rangle = \langle u, A^T v \rangle \quad (3.70)$$

for all $u \in \mathbb{R}^p$ and $v \in \mathbb{R}^n$ [See Definition B.3].

Proposition 3.7: Invertibility of Ridge Regression

Let X be a $n \times p$ matrix with real entries. If $\lambda > 0$, then the matrix

$$X^T X + \lambda \mathbb{I} \quad (3.71)$$

is invertible.

Proof. Let $\langle \cdot, \cdot \rangle$ denote the inner product on \mathbb{R}^p . Then, we have that

$$\langle (X^T X + \lambda \mathbb{I})u, u \rangle = \lambda \langle u, u \rangle + \langle X^T X u, u \rangle \quad (3.72)$$

$$= \lambda \langle u, u \rangle + \langle Xu, Xu \rangle_n \quad (3.73)$$

$$\geq \lambda \langle u, u \rangle \quad (3.74)$$

$$= \lambda \|u\|^2 \quad (3.75)$$

$$> 0, \quad (3.76)$$

where $\langle \cdot, \cdot \rangle_n$ denotes the inner product on \mathbb{R}^n . We have also used Lemma 3.1 to have the equivalence $\langle Xu, Xu \rangle_n = \langle X^T X u, u \rangle$. Hence, we must have that $\forall u \in \mathbb{R}^p \setminus \{\mathbf{0}\}$, $(X^T X + \lambda \mathbb{I})u \neq \mathbf{0}$. This tells us that $\ker(X^T X + \lambda) := \{v : (X^T X + \lambda)v = \mathbf{0}\} = \{\mathbf{0}\}$. Hence, by the rank-nullity theorem, $X^T X + \lambda$ must have full rank and is therefore invertible. \square

Proposition 3.8: Ridge Estimates for Orthonormal Inputs

Suppose that our inputs were orthonormal. That is, if \tilde{x}_i denotes the i^{th} column of \mathbf{X} , then $\tilde{x}_i \cdot \tilde{x}_j = \delta_{ij}$. It can easily be shown that $\mathbf{X}^T \mathbf{X} = \mathbb{I}$ as a consequence of this. This gives us the ridge estimates as just a scaled version of the least squares estimates:

$$\hat{\beta}^{\text{ridge}} = \frac{\hat{\beta}}{1 + \lambda}. \quad (3.77)$$

Lemma 3.2: Decomposition Equivalence

Let M be a $p \times p$ matrix that can be decomposed as $M = \mathbf{U} \Sigma \mathbf{U}^T$ where \mathbf{U} is a $p \times p$ matrix and Σ is a $p \times p$ diagonal matrix. Let $\sigma_i = \Sigma_{ii}$ be the diagonal entries and \mathbf{u}_j denote the $p \times 1$ j^{th} column vector of \mathbf{U} . Then, we have the equivalence

$$M = \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{u}_i^T \quad (3.78)$$

Proposition 3.9: Singular Value Decomposition of Input Matrix

Consider the centered $n \times p$ input matrix \mathbf{X} . Its singular value decomposition is taken to be

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T, \quad (3.79)$$

where \mathbf{U} and \mathbf{V} are $n \times p$ and $p \times p$ orthogonal matrices respectively. The columns of \mathbf{U} span the column space of \mathbf{X} and the columns of \mathbf{V} span the row space. \mathbf{D} is a $p \times p$ diagonal matrix, with diagonal entries $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ called the singular values of \mathbf{X} . If one or more values $d_j = 0$, then \mathbf{X} is singular.

Using the singular value decomposition, we can write the least squares fitted vector as

$$\mathbf{X} \hat{\beta}^{\text{ls}} = \mathbf{U} \mathbf{U}^T \mathbf{y}. \quad (3.80)$$

The ridge solutions are given by

$$\mathbf{X}\hat{\beta}^{\text{ridge}} = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j \mathbf{u}_j^T \mathbf{y}, \quad (3.81)$$

where \mathbf{u}_j is the j^{th} column of \mathbf{U} .

Definition 3.6: Sample Covariance Matrix

Let $\mathbf{x}_i \in \mathbb{R}^p$ denote the i^{th} input data vector with $1 \leq i \leq n$ and x_{ij} denote the j^{th} component of \mathbf{x}_i . Let $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$. Then the sample covariance matrix, denoted by \mathbf{Q} has matrix elements

$$Q_{ij} = \frac{1}{n-1} \sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j). \quad (3.82)$$

If we centre our data inputs for \mathbf{X} so that $\mathbf{X}_{ij} = x_{ij} - \bar{x}_j$, then the matrix $\mathbf{X}^T \mathbf{X}$ is the sample covariance matrix for our data.

3.4.2 The Lasso

Definition 3.7: Lasso Estimate

The lasso is a shrinkage method like ridge, with subtle but important differences. The lasso estimate is defined by

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \quad (3.83)$$

$$\text{subject to } \sum_{j=1}^p |\beta_j| \leq t. \quad (3.84)$$

Just as in ridge regression, we can write the lasso problem in the equivalent Lagrangian form

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\} \quad (3.85)$$

One can notice the similarity to the ridge regression problem, the L_2 ridge penalty $\sum_{j=1}^p \beta_j^2$ is replaced by the L_1 lasso penalty $\sum_{j=1}^p |\beta_j|$.

3.5 Principal Component Analysis

While ESLR didn't have a focused section on principle component analysis (at least in Chapter 3), they did introduce terminology and some definitions associated with PCA but didn't have a well structured layout of the algorithm and reasoning. I have therefore used external resources to write this section for myself as PCA is something I deem to be quite important.

Lemma 3.3

Let X be a random vector of size p . Let α be a vector of size p . Then

$$\text{Var}[\alpha^T X] = \alpha^T \Sigma \alpha, \quad (3.86)$$

where $\Sigma = \text{Cov}(X, X)$ is the covariance matrix for X given by $\Sigma = E[XX^T] - E[X]E[X^T]$.

Lemma 3.4

Let X be a random vector of size p , α be a vector of size p and β be a vector of size p . Then $\alpha^T X$ and $\beta^T X$ are uncorrelated^a if and only if

$$\alpha^T \Sigma \beta = 0, \quad (3.87)$$

where $\Sigma = \text{Cov}(X, X)$ is the covariance matrix for X given by $\Sigma = E[XX^T] - E[X]E[X^T]$.

^aWe say that two random variables X and Y are uncorrelated if $\text{Cov}(X, Y) = 0$.

Definition 3.8: Principal Component Analysis

Suppose that a random vector X of size p constitutes a basis for the feature space of our data. It's quite probable that this isn't the most relevant set of coordinates for our data and we instead seek out "optimal coordinates". In essence, we want to construct a new set of linearly transformed coordinates (Z_1, \dots, Z_k) with $Z_i = \alpha_i^T X$ for all $1 \leq i \leq k$, where α_i satisfy some desirable properties.

We define the first principal component by $Z_1 = \alpha_1^T X$ as where α_1 is a linear transformation that maximizes the variance:

$$\alpha_1 := \underset{\alpha \in \mathbb{R}^p}{\text{argmax}} \text{Var}[\alpha^T X] \quad (3.88)$$

$$\text{subject to } \alpha_1^T \alpha_1 = 1 \quad (3.89)$$

We want to maximize independence among our new random variables and so we impose a new condition on the j^{th} transformation α_j to be uncorrelated with all transformed components α_i for all $1 \leq i \leq j-1$. The j^{th} principal component is therefore defined by $Z_j = \alpha_j^T X$ where α_j satisfies

$$\alpha_j := \underset{\alpha \in \mathbb{R}^p}{\text{argmax}} \text{Var}[\alpha^T X] \quad (3.90)$$

$$\text{subject to } \alpha_j^T \alpha_j = 1, \quad (3.91)$$

$$\text{and } \text{Cov}(\alpha_j^T X, \alpha_i^T X) = 0 \ \forall \ i = 1, \dots, j-1 \quad (3.92)$$

Theorem 3.1: Principal Components

Let X be a random vector of size p , Σ denote the covariance matrix of X and α_i be the principal transformation vectors that ensure $Z_i = \alpha_i^T X$ are the principal components. If Σ 's k largest eigenvalues are non-zero, then α_i are all orthonormal eigenvectors of Σ with descending eigenvalues.

$$\Sigma \alpha_1 = \lambda_1 \alpha_1, \quad (3.93)$$

$$\vdots \quad (3.94)$$

$$\Sigma \alpha_k = \lambda_k \alpha_k, \quad (3.95)$$

where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$. Furthermore, α_k are unique up to a sign.

Proof. We begin with the first principal component. Since we have a constraint $\alpha_1^T \alpha_1 = 1$, we can introduce a Lagrange multiplier λ_1 to then solve for

$$\alpha_1 = \operatorname{argmax}_{\alpha \in \mathbb{R}^p} \operatorname{Var}[\alpha^T X] - \lambda_1(\alpha^T \alpha - 1) \quad (3.96)$$

$$= \operatorname{argmax}_{\alpha \in \mathbb{R}^p} \alpha^T \Sigma \alpha - \lambda_1(\alpha^T \alpha - 1) \quad (3.97)$$

Letting $f_1(\alpha) = \alpha^T \Sigma \alpha - \lambda_1(\alpha^T \alpha - 1)$, our critical points satisfy

$$0 = \nabla_{\alpha} f_1 = 2\Sigma \alpha - 2\lambda_1 \alpha. \quad (3.98)$$

Hence, we have the requirement that

$$\Sigma \alpha_1 = \lambda_1 \alpha_1, \quad (3.99)$$

that is; α_1 is an eigenvector of Σ with eigenvalue λ_1 . Since, we are maximizing the variance $\operatorname{Var}[\alpha^T X]$, we have that

$$\operatorname{Var}[\alpha_1^T X] = \alpha_1^T \Sigma \alpha_1 = \lambda_1. \quad (3.100)$$

Hence, α_1 must be an eigenvector of Σ associated with the largest eigenvalue λ_1 . We will now show that the rest of the solutions are eigenvectors via induction.

Base Case: $j=2$

For the second principal component, we have the two constraints $\alpha_2^T \alpha_2 = 1$ and $\alpha_2^T \Sigma \alpha_1 = \lambda_1 \alpha_2^T \alpha_1 = 0$. Since α_1 is an eigenvector of Σ , then demanding that $\alpha_2^T X$ and $\alpha_1^T X$ to be uncorrelated is equivalent to α_2 being orthogonal to α_1 (provided that $\lambda_1 \neq 0$). Just as for the first principal component, we can express this optimization-constraint problem by Lagrange multipliers. We have that

$$\alpha_2 = \operatorname{argmax}_{\alpha \in \mathbb{R}^p} \alpha^T \Sigma \alpha - \lambda_2(\alpha^T \alpha - 1) - \phi \alpha^T \alpha_1 \quad (3.101)$$

Letting $f_2(\alpha) = \alpha^T \Sigma \alpha - \lambda_2(\alpha^T \alpha - 1) - \phi \alpha^T \alpha_1$, our critical points satisfy

$$0 = \nabla_{\alpha} f_2 = 2\Sigma \alpha - 2\lambda_2 \alpha - \phi \alpha_1 \quad (3.102)$$

Hence, we must have that $\Sigma \alpha_2 = \lambda_2 \alpha_2 + \frac{1}{2} \phi \alpha_1$. Left-multiplying α_1 onto this expression, one obtains

$$\alpha_1^T \Sigma \alpha_2 = \lambda_2 \alpha_1^T \alpha_2 + \phi \alpha_1^T \alpha_1. \quad (3.103)$$

However, we have that $\alpha_1^T \Sigma \alpha_2 = 0$, $\alpha_1^T \alpha_2 = 0$ and $\alpha_1^T \alpha_1 = 1$ so we require $\phi = 0$. Hence, we have that

$$\Sigma \alpha_2 = \lambda_2 \alpha_2. \quad (3.104)$$

Since we are trying to maximize the variance $\operatorname{Var}[\alpha_2^T X] = \lambda_2$, then λ_2 must be the second largest eigenvalue of Σ with eigenvector α_2 .

Inductive Step: Assume that it holds for $j = 2, \dots, k-1$

As per the inductive hypothesis, we assume that the principal component transformations α_j are eigenvectors of Σ with associated eigenvalues λ_j for all $1 \leq j \leq k-1$. We will show that this implies that α_k is also an eigenvector of Σ for some eigenvalue λ_k . The PCA optimization problem imposes j constraints on α_j and we therefore introduce $j-1$ Lagrange multipliers ϕ_i^j (j is an upper index here) for the correlation conditions and λ_j for the normalization condition. Hence, the solution to the k^{th} principal component transformation is given by

$$\alpha_k = \operatorname{argmax}_{\alpha \in \mathbb{R}^p} \alpha^T \Sigma \alpha - \lambda_k(\alpha^T \alpha - 1) - \sum_{i=1}^{k-1} \phi_i^k \alpha^T \Sigma \alpha_i. \quad (3.105)$$

Letting $f_k(\alpha) = \alpha^T \Sigma \alpha - \lambda_2(\alpha^T \alpha - 1) - \sum_{i=1}^{k-1} \phi_i^k \alpha^T \Sigma \alpha_i$, the critical points satisfy

$$0 = \nabla_{\alpha} f_k = 2\Sigma \alpha - 2\lambda_k \alpha - \sum_{i=1}^{k-1} \phi_i^k \Sigma \alpha_i. \quad (3.106)$$

Hence, we have the requirement that

$$\Sigma \alpha_k = \lambda_k \alpha_k + \frac{1}{2} \sum_{i=1}^{k-1} \phi_i^k \Sigma \alpha_i. \quad (3.107)$$

Consider left-multiplying α_j for $1 \leq j \leq k-1$, so we obtain

$$\alpha_j^T \Sigma \alpha_k = \lambda_k \alpha_j^T \alpha_k + \frac{1}{2} \sum_{i=1}^{k-1} \phi_i^k \alpha_j^T \Sigma \alpha_i \quad (3.108)$$

By the correlation conditions, we have that $\alpha_j^T \Sigma \alpha_k = 0 \forall 1 \leq j \leq k-1$. Similarly, by the inductive hypothesis, we have that $\alpha_j^T \Sigma \alpha_k = \alpha_k^T \Sigma \alpha_j = \lambda_j \alpha_k^T \alpha_j = 0 \forall 1 \leq j \leq k-1$. By the assumption that $\lambda_j \neq 0$, this implies the orthogonality condition among α_k, α_j for $k \neq j$. Hence, (3.108) becomes

$$0 = \frac{1}{2} \sum_{i=1}^{k-1} \phi_i^k \lambda_i \alpha_j^T \alpha_i = \frac{1}{2} \sum_{i=1}^{k-1} \phi_i^k \lambda_i \delta_{ij} = \frac{1}{2} \phi_j^k \lambda_j \quad \forall 1 \leq j \leq k-1 \quad (3.109)$$

Provided that $\lambda_j \neq 0$, we must have $\phi_j^k = 0 \forall 1 \leq j \leq k-1$. Hence, it follows that α_k satisfies

$$\Sigma \alpha_k = \lambda_k \alpha_k, \quad (3.110)$$

with the variance $\text{Var}[\alpha_k^T X] = \lambda_k$ being the k^{th} largest eigenvalue of Σ with eigenvector α_k . By the induction hypothesis, this demonstrates that the linear transformations α_j are eigenvectors of Σ . In addition, orthonormality emerges very naturally in this proof.

Eigenvector Class of Solutions

Finally, we are interested in whether α_j are unique or must be a member of some class. It's quite obvious that the map $\alpha_j \mapsto -\alpha_j$ leaves the optimization problem directly intact and so if α_j is a solution, then so is $-\alpha_j$. We'll show that this family pair is indeed the only solutions to PCA for a given eigenvector. Suppose that for a given eigenvector α_k satisfying all the PCA properties, there exists another eigenvector β_k similarly satisfying all the PCA properties. Since $\{\alpha_1, \dots, \alpha_p\}$ form an orthogonal basis for \mathbb{R}^p , then one can write any vector $v \in \mathbb{R}^p$ as

$$v = \sum_{i=1}^p \langle v, \alpha_i \rangle \alpha_i, \quad (3.111)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product on \mathbb{R}^p . Since β_k also satisfy all the PCA properties, then it is orthogonal to every vector in $\text{span}(\alpha_1, \dots, \hat{\alpha}_k, \dots, \alpha_p)$, with $\hat{\alpha}_k$ indicating that α_k isn't present. Hence, one can write

$$\beta_k = \langle \beta_k, \alpha_k \rangle \alpha_k. \quad (3.112)$$

In essence, $\beta_k \in \text{span}(\alpha_k)$. We therefore have $\beta_k = c \alpha_k$ for some $c \in \mathbb{R}$. By the normalization condition assumed on β_k and α_k , we have

$$1 = \beta_k^T \beta_k = c^2 \alpha_k^T \alpha_k = c^2, \quad (3.113)$$

whose only solutions are given by $c = \pm 1$, concluding the proof. \square

Chapter 4

Linear Methods for Classification

4.1 Introduction

Definition 4.1: Hyperplane

Let $a_1, a_2, \dots, a_n \in \mathbb{R} \setminus \{0\}$ be scalars. Let x_i denote the i^{th} component of a vector $x \in \mathbb{R}^n$. The set

$$\{x \in \mathbb{R}^n : \sum_{i=1}^n a_i x_i = c\} \quad (4.1)$$

for some constant $c \in \mathbb{R}$ is a subspace of \mathbb{R}^n called a *hyperplane*.

Definition 4.2: Linear Decision Boundaries

Suppose that there are K classes, for convenience labeled $1, 2, \dots, K$ and the fitted linear model for the k^{th} indicator response variable is $\hat{f}_k(x) = \hat{\beta}_{k0} + \hat{\beta}_k^T x$ where $x \in \mathbb{R}^n$. The decision boundary (which is linear for our model) between class k and l is the set of points for which $\hat{f}_k(x) = \hat{f}_l(x)$, that is, the set

$$\{x \in \mathbb{R}^n : (\hat{\beta}_{k0} - \hat{\beta}_{l0}) + (\hat{\beta}_k - \hat{\beta}_l)^T x = 0\}, \quad (4.2)$$

which is a hyperplane of \mathbb{R}^n .

Definition 4.3: Odds

In statistics, the odds of an event occurring give some measure of likelihood of the event occurring. The odds are defined as the ratio between the probability that the event occurs and probability that the event doesn't occur. If p denotes the probability of an event occurring, then its odds is given by

$$o = \frac{p}{1-p}, \quad p \in [0, 1) \quad (4.3)$$

Definition 4.4: Posterior Probabilities for 2 Classes: Logit Transformation

A popular model for modelling the posterior class probability $\Pr(G = k|X = x)$ for 2 classes is given by

$$\Pr(G = 1|X = x) = \frac{\exp(\beta_0 + \beta^T x)}{1 + \exp(\beta_0 + \beta^T x)}, \quad (4.4)$$

$$\Pr(G = 2|X = x) = \frac{1}{1 + \exp(\beta_0 + \beta^T x)}. \quad (4.5)$$

The logit transformation $\log[p/(1 - p)]$ is a monotone transformation that satisfies

$$\log\left(\frac{\Pr(G = 1|X = x)}{\Pr(G = 2|X = x)}\right) = \beta_0 + \beta^T x. \quad (4.6)$$

One can therefore see that the decision boundary is the set of points for which the log-odds are zero, and this hyperplane is defined by $\{x|\beta_0 + \beta^T x = 0\}$.

4.2 Linear Regression of an Indicator Matrix

Definition 4.5: Indicator Response Matrix

Suppose that we have a set of training data $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ where $x^{(i)} \in \mathbb{R}^p \forall i$. Suppose that we choose to use a linear model so that $\hat{y}_i = \tilde{x}^{(i)T} \beta$ where $\tilde{x}^{(i)} = (x^{(i)}, 1) \in \mathbb{R}^{p+1}$. Since we are doing classification, suppose that \mathcal{G} has K classes. Then, there will be K indicators $Y_k, k = 1, \dots, K$, with $Y_k = 1$ if $G = k$ and 0 otherwise. We collect these into a vector $Y = (Y_1, \dots, Y_K)$. The N training instances of these form an $N \times K$ indicator response matrix \mathbf{Y} . Hence, \mathbf{Y} is a matrix of 0's and 1's with each row having a single 1.

Proposition 4.1

Suppose that we have a set of training data $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$ where $x^{(i)} \in \mathbb{R}^p \forall i$. Let \mathbf{Y} denote the indicator response matrix. Following the procedure of Chapter 3, we want to find the best set of parameters β for the linear model that minimizes the residual sum of squares (RSS). Let \mathbf{X} denote the $N \times (p + 1)$ matrix with each row corresponding to our input vectors. In essence

$$\mathbf{X}_{ij} = \begin{cases} x_j^{(i)} & \text{if } 1 \leq j \leq p \\ 1 & \text{if } j = p + 1 \end{cases} \quad (4.7)$$

Regressing linearly, one obtains the best fit to the linear regression model as

$$\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (4.8)$$

Hence, the $(p + 1) \times K$ coefficient matrix is given by $\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$.

Algorithm 1: Simple Linear Classification

Data: We assume a linear model has already been solved for an RSS loss function with coefficient matrix $\hat{\mathbf{B}}$.

input : Data point $x \in \mathbb{R}^p$

output: Classification \hat{y} for x

- 1 Compute the fitted output $\hat{f}(x)^T = (1, x^T)\hat{\mathbf{B}}$, a K vector;
- 2 $\hat{y} \leftarrow \operatorname{argmax}_{k \in \mathcal{G}} \hat{f}_k(x)$;

4.3 Linear Discriminant Analysis

Definition 4.6: Bayes Theorem for Posterior Probabilities

Suppose that $f_k(x)$ is the class-conditional density of X in class $G = k$, and let π_k be the prior probability of class k , with $\sum_{k=1}^K \pi_k = 1$. A simple application of Bayes Theorem gives us

$$\Pr(G = k | X = x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l} \quad (4.9)$$

Hence, one can see that in terms of the ability to classify, having $f_k(x)$ is almost equivalent to having the quantity $\Pr(G = k | X = x)$. Many techniques are based on models for the class densities:

- Linear and Quadratic discriminant analysis use Gaussian densities;
- Flexible mixtures of Gaussians allow for nonlinear decision boundaries;
- Nonparametric density estimates for each class density allow most flexibility;
- *Naive Bayes* models are a variant of the previous case, and assume that each of the class densities are products of marginal densities (They assume that inputs are conditionally independent in each class).

Proposition 4.2: Emergence of LDA for Multivariate Gaussian

Suppose that we model each class density as multivariate Gaussian

$$f_k(x) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)}. \quad (4.10)$$

Appendix A

Probability Theory

A.1 Constructing a Probability Space

The first thing we must do is establish the fundamental structures comprising a probability space. We will define a set Ω such that its points ω are associated with possible outcomes of a measurement. We also denote \mathcal{A} to be a nonempty collection of subsets of Ω which will represent collection of *events* that will be assigned probabilities.

Definition A.1: Sample Space, Ω

A set Ω with outcomes s_1, s_2, \dots, s_n (i.e. $\Omega = \{s_1, s_2, \dots, s_n\}$) must meet some conditions in order to be a sample space:

- The outcomes must be mutually exclusive, i.e. if s_j takes place, then no other s_i will take place, $\forall i, j \in \{1, 2, \dots, n\} \quad i \neq j$.
- The outcomes must be collectively exhaustive, i.e., on every experiment (or random trial) there will always take place some outcome $s_i \in \Omega$ for $i \in \{1, 2, \dots, n\}$.
- The sample space Ω must have the right granularity depending on what we are interested in. We must remove irrelevant information from the sample space. In other words, we must choose the right abstraction (forget some irrelevant information).

Definition A.2: σ -algebra / σ -field \mathcal{A} [Event Space]

A non-empty collection of subsets \mathcal{A} of set Ω is called a σ -field of subsets of Ω provided that the following two properties hold:

- (i) If A is in \mathcal{A} , then A^c is also in \mathcal{A} .
- (ii) If A_n is in \mathcal{A} , $n = 1, 2, \dots$, then $\bigcup_{n=1}^{\infty} A_n$ and $\bigcap_{n=1}^{\infty} A_n$ are both in \mathcal{A} .

Definition A.3: Event

Given a σ -field \mathcal{A} that corresponds to some sample space Ω . We say that if $A \in \mathcal{A}$, then A is an *event*.

The statement 'the event A occurs' means that the outcome of our experiment is represented by some point $\omega \in A$. For an event A , if we let $P(A)$ denote the probability of the event, then we have $0 \leq P(A) \leq 1$.

Definition A.4: Probability Measure

A probability measure P on a σ -field of subsets \mathcal{A} of a set Ω is a real valued function having domain \mathcal{A} satisfying the following properties:

- (i) $P(\Omega) = 1$
- (ii) $P(A) \geq 0 \quad \forall A \in \mathcal{A}$
- (iii) If $A_n, n = 1, 2, 3, \dots$ are mutually disjoint sets in \mathcal{A} , then $P(\cup_{n=1}^{\infty} A_n) = \sum_{n=1}^{\infty} P(A_n)$.

Definition A.5: Probability Space

A probability space, denoted by (Ω, \mathcal{A}, P) is a set Ω , a σ -field of subsets \mathcal{A} , and probability measure P defined on \mathcal{A} .

A.2 Standard Definitions and Properties

Definition A.6: Conditional Probability

Let A and B be two events such that $P(A) > 0$. Then the conditional probability of B given A , written $P(B|A)$, is defined to be

$$P(B|A) = \frac{P(B \cap A)}{P(A)} \quad (\text{A.1})$$

If $P(A) = 0$, then the conditional probability of B given A is undefined.

Proposition A.1

Let A be an event and A^c be its complement, defined as $A^c = \Omega - A$. It follows from the properties of disjoint probability sets that

$$P(A^c) = 1 - P(A) \quad (\text{A.2})$$

Definition A.7: Independent Events

Two events A and B are independent if and only if

$$P(A \cap B) = P(A)P(B) \quad (\text{A.3})$$

This definition emerges as a consequence of wanting to construct a notion of an event's occurrence having no influence on the the occurrence of the other event. Through the conditional probabilistic lens, this would mean $P(B|A) = P(B)$ (i.e. Given that A has occurred, this does not affect the probability that B will occur). Therefore, it follows that $P(A \cap B) = P(A)P(B)$.

Definition A.8: Mutual Exclusivity

Events A and B are said to be two mutually exclusive events if both cannot occur. In essence, their intersection is disjoint $A \cap B = \emptyset$ so that they have the following properties:

$$P(A \cap B) = 0 \quad (\text{A.4})$$

$$P(A \cup B) = P(A) + P(B) \quad (\text{A.5})$$

Definition A.9: Discrete Random Variable

A discrete real-valued random variable X on a probability space (Ω, \mathcal{A}, P) is a function X with domain Ω and range that is a finite or countably infinite subset $\{x_1, x_2, \dots\}$ of the real numbers \mathbb{R} such that $\{\omega : X(\omega) = x_i\}$ is event for all i .

Hence, $\{\omega : X(\omega) = x_i\}$ is an event and we usually will write $\{X = x_i\}$ for brevity and denote the probability of this event as $P(X = x_i)$.

Definition A.10: Discrete Density Function

The real-valued function f defined on \mathbb{R} by $f(x) = P(X = x)$ is called the discrete density function of X . A number x is called a possible value of X if $f(x) > 0$.

We note that a real-valued function f defined on \mathbb{R} is called a discrete density function provided that it satisfies the following properties:

- (i) $f(x) \geq 0, x \in \mathbb{R}$.
- (ii) $\{x : f(x) \neq 0\}$ is a finite or countably infinite subset of \mathbb{R} . Let $\{x_1, x_2, \dots\}$ denote this set. Then
- (iii) $\sum_i f(x_i) = 1$.

We can compute the probability of X taking on value in some set A via

$$P(X \in A) = \sum_{x \in A} f(x) \quad (\text{A.6})$$

Definition A.11: Discrete r -dimensional Random Vector

We let \mathbb{R}^r denote the collection of all r -tuples of real numbers. A point $\mathbf{x} = (x_1, x_2, \dots, x_r)$ of \mathbb{R}^r is usually called an r -dimensional vector. Thus for each $\omega \in \Omega$, the r values $X_1(\omega), \dots, X_r(\omega)$ define a point

$$\mathbf{X}(\omega) = (X_1(\omega), X_2(\omega), \dots, X_r(\omega)) \quad (\text{A.7})$$

of \mathbb{R}^r . This defines an r -dimensional vector-valued function on Ω , $\mathbf{X} : \Omega \rightarrow \mathbb{R}^r$, which is usually written as $\mathbf{X} = (X_1, X_2, \dots, X_r)$.

A discrete r -dimensional random vector \mathbf{X} is a function \mathbf{X} from Ω to \mathbb{R}^r taking on a finite or countably infinite number of values $\mathbf{x}_1, \mathbf{x}_2, \dots$ such that

$$\{\omega : \mathbf{X}(\omega) = \mathbf{x}_i\} \quad (\text{A.8})$$

is an event for all i .

Definition A.12: Discrete Density Function for Random Vector

The discrete density function f for the random vector \mathbf{X} is defined by

$$f(x_1, \dots, x_r) = P(X_1 = x_1, \dots, X_r = x_r) \quad (\text{A.9})$$

or equivalently

$$f(\mathbf{x}) = P(\mathbf{X} = \mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^r \quad (\text{A.10})$$

The probability that \mathbf{X} belongs to the subset A of \mathbb{R}^r can be found by using the analog of (A.6), namely

$$P(\mathbf{X} \in A) = \sum_{\mathbf{x} \in A} f(\mathbf{x}) \quad (\text{A.11})$$

Definition A.13: Mutually Independent Random Variables

Let X_1, X_2, \dots, X_r be r discrete random variables having densities f_1, f_2, \dots, f_r respectively. These random variables are said to be *mutually independent* if their joint density function f is given by

$$f(x_1, x_2, \dots, x_r) = f_1(x_1)f_2(x_2) \cdots f_r(x_r) \quad (\text{A.12})$$

Consider two independent discrete random variables having densities f_X and f_Y , respectively. Then for any two subsets A and B of R , we have

$$P(X \in A, Y \in B) = P(X \in A)P(Y \in B) \quad (\text{A.13})$$

Definition A.14: Probability Generating Function

Let X be a non-negative integer-valued random variable. The probability generating function Φ_X of X is defined as

$$\Phi_X(t) = \sum_{x=0}^{\infty} P(X = x)t^x = \sum_{x=0}^{\infty} f_X(x)t^x, \quad -1 \leq t \leq 1 \quad (\text{A.14})$$

Definition A.15: Random Variable

A random variable X on a probability space (Ω, \mathcal{A}, P) is a real-valued function $X(\omega)$, $\omega \in \Omega$, such that for $-\infty < x < \infty$, $\{\omega | X(\omega) \leq x\}$ is an event.

Definition A.16: Continuous Random Variable

A random variable X is called a *continuous random variable* if

$$P(X = x) = 0, \quad -\infty < x < \infty \quad (\text{A.15})$$

We can observe that X is a continuous random variable if and only if its distribution function F is continuous at every x , that is, F is a continuous function.

Definition A.17: Symmetric Random Variable

A random variable X is said to be *symmetric* if X and $-X$ have the same distribution function.

Definition A.18: Median

For any probability distribution on the real line \mathbb{R} with cumulative distribution function F , regardless of whether it is any kind of continuous probability distribution, in particular an absolutely continuous distri-

bution, or a discrete probability distribution, a median is by definition any real number m that satisfies the inequalities

$$P(X \leq m) = \frac{1}{2}, \quad P(X \geq m) = \frac{1}{2} \quad (\text{A.16})$$

A.3 Distributions and Densities

Let X and Y be two discrete random variables. For any real numbers x and y , the set $\{\omega | X(\omega) = x \text{ and } Y(\omega) = y\}$ is an event that we will usually denote by $\{X = x, Y = y\}$.

Definition A.19: Joint Density and Marginal Density

Let $\mathbf{X} = (X_1, X_2, \dots, X_r)$ be an r -dimensional random vector with density f . Then the function f is usually called the *joint density* of the random variables X_1, X_2, \dots, X_r . The density function of the random variable X_i is then called the i^{th} *marginal density* of \mathbf{X} or of f .

Definition A.20: (Cumulative) Distribution Function [Discrete]

The function $F(t)$, $-\infty < t < \infty$, defined by

$$F(t) = P(X \leq t) = \sum_{x \leq t} f(x), \quad -\infty < t < \infty \quad (\text{A.17})$$

is called the *distribution function* of the random variable X or of the density f . One immediate consequence of this is that it satisfies:

$$P(a < X \leq b) = P(X \leq b) - P(X \leq a) = F(b) - F(a) \quad (\text{A.18})$$

Proposition A.2

Let X and Y be independent, non-negative integer-valued random variables. Then

$$\Phi_{X+Y}(t) = \Phi_X(t)\Phi_Y(t) \quad (\text{A.19})$$

Definition A.21: (Cumulative) Distribution Function [Continuous]

The distribution function F of a random variable X is the function

$$F(x) = P(X \leq x), \quad -\infty < x < \infty \quad (\text{A.20})$$

Proposition A.3: Properties of Distribution Functions

Not all functions can arise as distribution functions, for the latter must satisfy certain conditions. Let X be a random variable and let F be its distribution function. Then

- (i) $0 \leq F(x) \leq 1$ for all x .
- (ii) F is a non-decreasing function of x .
- (iii) $F(-\infty) = 0$ and $F(+\infty) = 1$.

(iv) $F(x+) = F(x)$ for all x . (F is a right-continuous function)

We note that a distribution function is any function F satisfying properties (i)-(iv).

Definition A.22: Probability Density Function (PDF) / Density

A density function / PDF (with respect to integration) is a non-negative function f such that

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad (\text{A.21})$$

Note that if f is density function, then the function F defined by

$$F(x) = \int_{-\infty}^x f(y) dy, \quad -\infty < x < \infty \quad (\text{A.22})$$

is a continuous function satisfying properties (i)-(iv) in **Prop A.3**.

Definition A.23: Uniform Density

Let Ω be a sample space with finite measure $\text{Vol}(\Omega) < \infty$. Then, a uniform density is a constant function f , such that

$$1 = \int_{\Omega} f dV \quad (\text{A.23})$$

Hence, $f = 1/\text{Vol}(\Omega)$.

Example A.1: Uniform Density / Distribution on a Real Line Interval

Let a and b be constants with $a < b$. The uniform density on the interval (a, b) is the density f defined by

$$f(x) = \begin{cases} (b-a)^{-1} & \text{for } a < x < b, \\ 0 & \text{elsewhere} \end{cases} \quad (\text{A.24})$$

The distribution function corresponding to (A.24) is given by

$$F(x) = \begin{cases} 0 & x < a, \\ (x-a)/(b-a), & a \leq x \leq b, \\ 1, & x > b. \end{cases} \quad (\text{A.25})$$

Definition A.24: Binomial Density

Let $0 < p < 1$. Then, the real valued function f defined on \mathbb{R} by

$$f(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x}, & x = 0, 1, 2, \dots, n, \\ 0, & \text{elsewhere} \end{cases} \quad (\text{A.26})$$

is called the *binomial density* with parameters n and p .

Definition A.25: Geometric Density

Let $0 < p < 1$. Then the real valued function f defined on \mathbb{R} by

$$f(x) = \begin{cases} p(1-p)^x, & x = 0, 1, 2, \dots, \\ 0, & \text{elsewhere} \end{cases} \quad (\text{A.27})$$

is a discrete density function called the *geometric density* with parameter p .

Definition A.26: Poisson Density

Let $0 < p < 1$ and let λ be a positive number. Then, the real valued function f defined on \mathbb{R} by

$$f(x) = \begin{cases} \frac{\lambda^x e^{-\lambda}}{x!}, & x = 0, 1, 2, \dots, \\ 0, & \text{elsewhere.} \end{cases} \quad (\text{A.28})$$

is called the *Poisson density* with parameter λ .

Proposition A.4: Binomial Theorem

Let $0 < p < 1$ and $x < n \in \mathbb{Z}$. Then, we have that

$$1 = \sum_{x=0}^n \binom{n}{x} p^x (1-p)^{n-x} \quad (\text{A.29})$$

Which follows from the binomial theorem

$$(a+b)^n = \sum_{x=0}^n \binom{n}{x} a^x b^{n-x} \quad (\text{A.30})$$

Proposition A.5

Let ϕ be a differentiable strictly increasing or strictly decreasing function on an interval I , and let $\phi(I)$ denote the range of ϕ and ϕ^{-1} the inverse function to ϕ . Let X be a continuous random variable having density f such that $f(x) = 0$ for $x \notin I$. Then $Y = \phi(X)$ has density g given by $g(y) = 0$ for $y \notin \phi(I)$ and

$$g(y) = f(\phi^{-1}(y)) \left| \frac{d}{dy} \phi^{-1}(y) \right|, \quad y \in \phi(I) \quad (\text{A.31})$$

It is a bit more suggestive to write this in the following form:

$$g(y) = f(x) \left| \frac{dx}{dy} \right|, \quad y \in \phi(I) \quad \text{and} \quad x = \phi^{-1}(y) \quad (\text{A.32})$$

Definition A.27: Cauchy Density

The following function f , is a density known as the *Cauchy Density*.

$$f(x) = \frac{1}{\pi(1+x^2)}, \quad -\infty < x < \infty \quad (\text{A.33})$$

The corresponding distribution function is given by

$$F(x) = \frac{1}{2} + \frac{1}{\pi} \arctan(x), \quad -\infty < x < \infty \quad (\text{A.34})$$

Definition A.28: Symmetric Density

A density function f is called *symmetric* if $f(-x) = f(x)$ for all x . The Cauchy density and the uniform density on $(-a, a)$ are both symmetric.

Proposition A.6

Let X be a random variable that has a density. Then f has a symmetric density if and only if X is a symmetric random variable.

Definition A.29: Standard Normal Density

The following density, ϕ

$$\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad -\infty < x < \infty \quad (\text{A.35})$$

The standard normal density is clearly symmetric.

The normal density with mean μ and variance σ^2 is often denoted by $n(\mu, \sigma^2)$ or $n(y; \mu, \sigma^2)$, $-\infty < y < \infty$. Thus,

$$n(y; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(y-\mu)^2/(2\sigma^2)}, \quad -\infty < y < \infty \quad (\text{A.36})$$

Definition A.30: Exponential Density

The exponential density with parameter λ is given by

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (\text{A.37})$$

The corresponding distribution function is

$$F(x) = \begin{cases} 1 - e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (\text{A.38})$$

Proposition A.7

Let X be a random variable such that the following holds:

$$P(X > a + b) = P(X > a)P(X > b), \quad a \geq 0 \quad \text{and} \quad b \geq 0 \quad (\text{A.39})$$

Then either $P(X > 0) = 0$ or X is exponentially distributed.

Proposition A.8: Sum of Random Variables

Let X, Y be continuous random variables with densities f_X and f_Y respectively. Then, the random variable $Z = X + Y$ has density f_Z , given by

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(z - k) f_Y(k) dk \quad (\text{A.40})$$

A.4 Expectations

Notation: Let \mathbf{X} be a discrete r -dimensional random vector having possible values $\mathbf{x}_1, \mathbf{x}_2, \dots$ and density f , and let ϕ be a real-valued function defined on \mathbb{R}^r . Then $\sum_{\mathbf{x}} \phi(\mathbf{x}) f(\mathbf{x})$ is defined as

$$\sum_{\mathbf{x}} \phi(\mathbf{x}) f(\mathbf{x}) = \sum_j \phi(\mathbf{x}_j) f(\mathbf{x}_j) \quad (\text{A.41})$$

Definition A.31: Expectation Value

Let X be any discrete random variable that assumes a finite number of values x_1, \dots, x_r . Then the expected value of X , denoted by EX , $E[X]$ or μ , is the number

$$E[X] = \sum_{i=1}^r x_i f(x_i) \quad (\text{A.42})$$

The expected value $E[X]$ is also called the mean of X .

Definition A.32: Finite / Undefined Expectation

Let X be a discrete random variable having density f . If $\sum_j |x_j| f(x_j) < \infty$, then we say that X has finite expectation and we define its expectation by (A.42). On the other hand if $\sum_j |x_j| f(x_j) = \infty$, then we say X does not have finite expectation and $E[X]$ is undefined.

Proposition A.9

Let \mathbf{X} be a discrete random vector having density f , and let ϕ be a real-valued function defined on \mathbb{R}^r . Then the random variable $Z = \phi(\mathbf{X})$ has finite expectation if and only if

$$\sum_{\mathbf{x}} |\phi(\mathbf{x})| f(\mathbf{x}) < \infty \quad (\text{A.43})$$

and, when (A.43) holds,

$$E[Z] = \sum_{\mathbf{x}} \phi(\mathbf{x}) f(\mathbf{x}) \quad (\text{A.44})$$

Proposition A.10: Properties of Expectation Operator

Let X and Y be two random variables having finite expectation.

(i) If c is a constant and $P(X = c) = 1$, then $E[X] = c$.

(ii) Linearity:

- a) If c is a constant, then cX has finite expectation and $E[cX] = cE[X]$.
 b) $X + Y$ has finite expectation and^a

$$E[X + Y] = E[X] + E[Y] \quad (\text{A.46})$$

- (iii) Suppose that $P(X \geq Y) = 1$. Then $E[X] \geq E[Y]$; moreover, $E[X] = E[Y]$ if and only if $P(X = Y) = 1$.
 (iv) $|E[X]| \leq E[|X|]$.

^aMore explicitly, we note that these are expectations w.r.t different densities:

$$E_{X+Y}[X + Y] = E_X[X] + E_Y[Y] \quad (\text{A.45})$$

Proposition A.11

Let X be a random variable such that for some constant M , $P(|X| \leq M) = 1$. Then X has finite expectation and $|E[X]| \leq M$.

Proposition A.12

Let X and Y be two independent random variables having finite expectations. Then XY has finite expectation and

$$E[XY] = E[X]E[Y] \quad (\text{A.47})$$

Proposition A.13

Let X be a non-negative integer-valued random variable. Then X has finite expectation if and only if the series $\sum_{x=1}^{\infty} P(X \geq x)$ converges. If the series does converge, then

$$E[X] = \sum_{x=1}^{\infty} P(X \geq x) \quad (\text{A.48})$$

Definition A.33: Moments / Central Moments

Let X be a discrete random variable, and let $r \geq 0$ be an integer. We say that X has a *moment* of order r if X^r has finite expectation. In that case we define the r^{th} of X as $E[X^r]$.

If X has a moment of order r then the r^{th} moment of $X - \mu$, where μ is the mean of X , is called the *central moment* (or the r^{th} about the mean) of X .

Proposition A.14

If the random variables X and Y have moments of order r , then $X + Y$ also has a moment of order r .

Definition A.34: Variance

Let X be a random variable having a finite second moment. Then the *variance* of X , denoted by $\text{Var}[X]$ or $V[X]$, is defined by

$$\text{Var}[X] = E[(X - E[X])^2] \quad (\text{A.49})$$

Through expanding, this works out to the following:

$$\text{Var}[X] = E[X^2] - (E[X])^2 \quad (\text{A.50})$$

Definition A.35: Standard Deviation

We often denote $\text{Var } X$ by σ^2 . The non-negative number $\sigma = \sqrt{\text{Var } X}$ is called the *standard deviation* of X or of f_X .

Definition A.36: Covariance

Let X and Y be two random variables each having finite second moment. We define a quantity called the *covariance* of X and Y written as $\text{Cov}(X, Y)$. Thus we have the formula

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y] \quad (\text{A.51})$$

We note that $X + Y$ has a finite second moment and finite variance. We therefore have an important formula:

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2 \text{Cov}(X, Y) \quad (\text{A.52})$$

Definition A.37: Correlation Coefficient

Let X and Y be two random variables having finite nonzero variances. One measure of the degree of dependence between the two random variables is the *correlation coefficient* $\rho(X, Y)$ defined by

$$\rho = \rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{(\text{Var}[X])(\text{Var}[Y])}} \quad (\text{A.53})$$

These random variables are said to be *uncorrelated* if $\rho = 0$. We can automatically see that independent random variables are uncorrelated. However, it is possible for dependent random variables to be uncorrelated as well. We observe that the correlation coefficient ρ is always between -1 and 1, and that $\rho = 1$ if and only if $P(X = aY) = 1$ for some constant a .

Definition A.38: Cross-Correlation Matrix

Let \mathbf{X}, \mathbf{Y} be random vectors. Then, we define the *cross-correlation matrix* by $E[\mathbf{X}\mathbf{Y}^T]$, where the matrix elements in the standard basis are given by $[E[\mathbf{X}\mathbf{Y}^T]]_{ij} = E[x_i y_j]$.

Theorem A.1: The Schwartz Inequality

Let X and Y have finite second moments. Then

$$[E[XY]]^2 \leq (E[X^2])(E[Y^2]) \quad (\text{A.54})$$

Furthermore, equality holds in (A.54) if and only if either $P(Y=0) = 1$ or $P(X = aY) = 1$ for some constant a .

Proposition A.15: Chebyshev's Inequality

Let X be a random variable with mean μ and finite variance σ^2 . Then for any real number $t > 0$

$$P(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2} \quad (\text{A.55})$$

Theorem A.2: Weak Law of Large Numbers

Let X_1, X_2, \dots, X_n be independent random variables having a common distribution with finite mean μ and set $S_n = X_1 + \dots + X_n$. Then for any $\delta > 0$

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{S_n}{n} - \mu\right| \geq \delta\right) = 0 \quad (\text{A.56})$$

A.5 Jointly Distributed Random Variables

Definition A.39: Joint Distribution Function

Let X and Y be two random variables defined on the same probability space. Their joint distribution function F is defined by

$$F(x, y) = P(X \leq x, Y \leq y), \quad -\infty < x, y < \infty \quad (\text{A.57})$$

Definition A.40: Marginal Distribution Functions

The one-dimensional distribution functions F_X and F_Y defined by

$$F_X(x) = P(X \leq x) \quad \text{and} \quad F_Y(y) = P(Y \leq y) \quad (\text{A.58})$$

are called the marginal distribution functions of X and Y . They are related to the joint distribution function F by

$$F_X(x) = F(x, \infty) = \lim_{y \rightarrow \infty} F(x, y) \quad (\text{A.59})$$

$$F_Y(y) = F(\infty, y) = \lim_{x \rightarrow \infty} F(x, y) \quad (\text{A.60})$$

Definition A.41: Joint Density Function

If there is a nonnegative function f such that

$$F(x, y) = \int_{-\infty}^x \left(\int_{-\infty}^y f(u, v) dv \right) du, \quad -\infty < x, y < \infty, \quad (\text{A.61})$$

then f is called a joint density function (with respect to integration) for the distribution function F or the pair of random variables X, Y .

$$P((X, Y) \in A) = \int_A \int f(x, y) dx dy \quad (\text{A.62})$$

By letting A be the entire plane we obtain

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \, dx \, dy = 1 \quad (\text{A.63})$$

Definition A.42: Marginal Density

Let F be the distribution function for a pair of random variables X, Y . Then, the marginal distribution F_X has marginal density f_X given by

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) \, dy \quad (\text{A.64})$$

Similarly, F_Y has marginal density f_Y given by

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) \, dx \quad (\text{A.65})$$

We note that it satisfies

$$F_X(x) = \int_{-\infty}^x f_X(u) \, du \quad (\text{A.66})$$

We can observe that

$$\frac{\partial^2}{\partial x \partial y} F(x, y) = f(x, y) \quad (\text{A.67})$$

Definition A.43: Independent Random Variables

The variables X and Y are called independent random variables if whenever $a \leq b$ and $c \leq d$, then

$$P(a < X \leq b, c < Y \leq d) = P(a < X \leq b)P(c < Y \leq d) \quad (\text{A.68})$$

By letting $a = c = -\infty$, $b = x$, and $d = y$, it follows that if X and Y are independent, then

$$F(x, y) = F_X(x)F_Y(y), \quad -\infty < x, y < \infty \quad (\text{A.69})$$

Proposition A.16

If X and Y are independent and A and B are unions of a finite or countably infinite number of intervals, then

$$P(X \in A, Y \in B) = P(X \in A)P(Y \in B) \quad (\text{A.70})$$

In other words, the events

$$\{\omega | X(\omega) \in A\} \quad \text{and} \quad \{\omega | X(\omega) \in B\} \quad (\text{A.71})$$

are independent events.

Proposition A.17

Let X and Y be random variables having marginal densities f_X and f_Y . Then X and Y are independent if and only if the function f defined by

$$f(x, y) = f_X(x)f_Y(y), \quad -\infty < x, y < \infty \quad (\text{A.72})$$

is a joint density for X and Y .

Definition A.44: Bivariate Density Function

A two-dimensional (or bivariate) density function f is a non-negative function on \mathbb{R}^2 such that

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \, dx \, dy = 1 \quad (\text{A.73})$$

Definition A.45: Standard Bivariate Normal Density

The density given below by f is referred to as the standard bivariate normal density.

$$f(x, y) = \frac{1}{2\pi} e^{-(x^2+y^2)/2}, \quad -\infty < x, y < \infty \quad (\text{A.74})$$

Proposition A.18

Let X and Y be random variables having joint density f . In many contexts, we will have a random variable Z defined in terms of X and Y and we wish to calculate the density of Z . Let $Z = \phi(X, Y)$, where ϕ is a real-valued function whose domains contains the range of X and Y . For fixed z the event $\{Z \leq z\}$ is equivalent to the event $\{(X, Y) \in A_z\}$ where A_z is the subset of \mathbb{R}^2 defined by

$$A_z = \{(x, y) | \phi(x, y) \leq z\} \quad (\text{A.75})$$

Thus,

$$F_Z(z) = P(Z \leq z) \quad (\text{A.76})$$

$$= P((X, Y) \in A_z) \quad (\text{A.77})$$

$$= \int_{A_z} \int f(x, y) \, dx \, dy \quad (\text{A.78})$$

If we can find a non-negative function g such that

$$\int_{A_z} \int f(x, y) \, dx \, dy = \int_{-\infty}^z g(v) \, dv, \quad -\infty < z < \infty \quad (\text{A.79})$$

then g is necessarily a density of Z .

Proposition A.19

Let X and Y be independent random variables having the respective normal densities $n(\mu_1, \sigma_1^2)$ and $n(\mu_2, \sigma_2^2)$. Then $X + Y$ has the normal density

$$n(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2) \quad (\text{A.80})$$

Definition A.46: Conditional Density (Discrete)

Let X and Y be discrete random variables having joint density f . If x is a possible value of X , then

$$P(Y = y | X = x) = \frac{P(X = x, Y = y)}{P(X = x)} = \frac{f(x, y)}{f_X(x)} \quad (\text{A.81})$$

The function $f_{Y|X}$ defined by

$$f_{Y|X}(y|x) = \begin{cases} \frac{f(x,y)}{f_X(x)}, & f_X(x) \neq 0 \\ 0, & f_X(x) = 0 \end{cases} \quad (\text{A.82})$$

is called the conditional density of Y given x .

Definition A.47: Conditional Density (Continuous)

Let X and Y be continuous random variables having joint density f . The conditional density $f_{Y|X}$ is defined by

$$f_{Y|X}(y|x) = \begin{cases} \frac{f(x,y)}{f_X(x)}, & 0 < f_X(x) < \infty, \\ 0 & \text{elsewhere.} \end{cases} \quad (\text{A.83})$$

If f_X is continuous and $f_X(x) \neq 0$, we have

$$P(a \leq Y \leq b | X = x) = \frac{\int_a^b f(x, y) dy}{f_X(x)} \quad (\text{A.84})$$

Proposition A.20: Bayes Rule

Let X and Y be random variables with marginal densities f_X and f_Y respectively and conditional densities $f_{X|Y}$ and $f_{Y|X}$. We have the continuous analog to Bayes' rule given below:

$$f_{X|Y}(x|y) = \frac{f_X(x)f_{Y|X}(y|x)}{\int_{-\infty}^{\infty} f_X(x)f_{Y|X}(y|x) dx} \quad (\text{A.85})$$

Definition A.48: Joint Distribution Function (Multivariate)

Let X_1, \dots, X_n be n random variables defined on a common probability space. Their joint distribution function F is defined by

$$F(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n), \quad -\infty < x_1, \dots, x_n < \infty \quad (\text{A.86})$$

Definition A.49: Marginal Distribution Function (Multivariate)

The marginal distribution functions $F_{X_m}, m = 1, \dots, n$ are defined by

$$F_{X_m}(x_m) = P(X_m \leq x_m), \quad -\infty < x_m < \infty \quad (\text{A.87})$$

The value of $F_{X_m}(x_m)$ can be obtained from F by letting $x_1, \dots, x_{m-1}, x_{m+1}, \dots, x_n$ all approach $+\infty$.

Definition A.50: Joint Density Function (Multivariate)

A non-negative function f is called a joint density function (with respect to integration) for the joint distribution function F , or for the random variables X_1, \dots, X_n if

$$F(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} f(u_1, \dots, u_n) du_1 \cdots du_n, \quad -\infty < x_1, \dots, x_n < \infty \quad (\text{A.88})$$

We also note that

$$f(x_1, \dots, x_n) = \frac{\partial^n}{\partial x_1 \cdots \partial x_n} F(x_1, \dots, x_n) \quad (\text{A.89})$$

is valid at the continuity points of F .

Definition A.51: Marginal Density Function (Multivariate)

The random variable X_m has the marginal density f_{X_m} obtained by integrating f over the remaining $n - 1$ variables. For example,

$$f_{X_2}(x_2) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_1 dx_3 \cdots dx_n \quad (\text{A.90})$$

Definition A.52: Independent Random Variables (Multivariate)

In general, the random variables X_1, \dots, X_n are called independent whenever $a_m \leq b_m$ for $m = 1, \dots, n$, then

$$P(a_1 < X_1 \leq b_1, \dots, a_n < X_n \leq b_n) = P(a_1 < X_1 \leq b_1) \cdots P(a_n < X_n \leq b_n) \quad (\text{A.91})$$

Proposition A.21

A necessary and sufficient condition for independence is that

$$F(x_1, \dots, x_n) = F_{x_1}(x_1) \cdots F_{x_n}(x_n), \quad -\infty < x_1, \dots, x_n < \infty \quad (\text{A.92})$$

If F has a density f , then X_1, \dots, X_n are independent if and only if f can be chosen so that

$$f(x_1, \dots, x_n) = f_{x_1}(x_1) \cdots f_{x_n}(x_n), \quad -\infty < x_1, \dots, x_n < \infty \quad (\text{A.93})$$

If X_1, \dots, X_n are random variables whose joint density is given by (A.93) then X_1, \dots, X_n are independent and X_m has the marginal density f_m .

Proposition A.22

Let X_1, \dots, X_n be independent random variables. Let Y be a random variable defined in terms of X_1, \dots, X_m and let Z be a random variable defined in terms X_{m+1}, \dots, X_n (where $1 < m < n$). Then Y and Z are independent.

Definition A.53: Conditional Density (Multivariate)

If X_1, \dots, X_n has a joint density f , then any subcollection of these random variables has a joint density which can be found by integrating over the remaining variables. For example, if $1 \leq m < n$,

$$f_{X_1, \dots, X_m}(x_1, \dots, x_m) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, \dots, x_n) dx_{m+1} \cdots dx_n \quad (\text{A.94})$$

The conditional density of a subcollection of X_1, \dots, X_n given the remaining variables can also be defined in an obvious manner. Thus the conditional density of X_{m+1}, \dots, X_n given X_1, \dots, X_m is defined by

$$f_{X_{m+1}, \dots, X_n | X_1, \dots, X_m}(x_{m+1}, \dots, x_n | x_1, \dots, x_m) = \frac{f(x_1, \dots, x_n)}{f_{X_1, \dots, X_m}(x_1, \dots, x_m)} \quad (\text{A.95})$$

Where f is the joint density of X_1, \dots, X_n .

Definition A.54: Order Statistics

Let U_1, \dots, U_n be independent continuous random variables, each having distribution F and density function f . Let X_1, \dots, X_n be random variables obtained by letting $X_1(\omega), \dots, X_n(\omega)$ be the set $U_1(\omega), \dots, U_n(\omega)$ permuted so as to be in increasing order. In particular, we define X_1 and X_n to be the functions

$$X_1(\omega) = \min\{U_1(\omega), \dots, U_n(\omega)\} \quad (\text{A.96})$$

and

$$X_n(\omega) = \max\{U_1(\omega), \dots, U_n(\omega)\} \quad (\text{A.97})$$

The random variable X_k is called the k^{th} order statistic. Another related variable of interest is the range R , defined by

$$R(\omega) = X_n(\omega) - X_1(\omega) \quad (\text{A.98})$$

$$= \max\{U_1(\omega), \dots, U_n(\omega)\} - \min\{U_1(\omega), \dots, U_n(\omega)\} \quad (\text{A.99})$$

Proposition A.23: Distributions for Order Statistics

Let X_1, X_2, \dots, X_n be identically distributed and independent random variables. Let their common CDF be denoted by F . We define $(X_{(1)}, X_{(2)}, \dots, X_{(n)})$ as the vector of order statistics of X_1, X_2, \dots, X_n . Then, the distribution for $X_{(k)}$ in a sample of size n is given by

$$F_{(k,n)}(x) = P(X_{(k)} \leq x) = \sum_{j=k}^n \binom{n}{j} [F(x)]^j [1 - F(x)]^{n-j} \quad (\text{A.100})$$

Proof. We will break the event $(X_{(k)} \leq x)$ into disjoint sub-events given by

$$(X_{(k)} \leq x) = (X_{(n)} \leq x) \cup (X_{(n)} > x, X_{(n-1)} \leq x) \cup \dots \cup (X_{(n)} > x, \dots, X_{(k+1)} > x, X_{(k)} \leq x). \quad (\text{A.101})$$

Recall from the property of probability measures that if A_i 's are all mutually disjoint sets, then $P(\cup_{i=1}^l A_i) = \sum_{i=1}^l P(A_i)$. Hence, it amounts to identify the probability of the events contained within each of these terms. Consider the event $(X_{(n)} > x, \dots, X_{(j+1)} > x, X_{(j)} \leq x)$. This event tells us that the first j ordered variables have a value less than x while the rest have a value lying above x . Since the CDF, $F(x)$ to each one tells us the probability of them occupying a value less than x , one can view this through the lens of fail / successes among n independent variables. In essence, we use the multiplicative property of independent events and combinatorics to establish the number of combinations one can arrange such an ordering of variables. We therefore have

$$P(X_{(n)} > x, \dots, X_{(j+1)} > x, X_{(j)} \leq x) = \binom{n}{j} [F(x)]^j [1 - F(x)]^{n-j}. \quad (\text{A.102})$$

Hence, it therefore follows that

$$F_{(k,n)}(x) = P(X_{(k)} \leq x) = \sum_{j=k}^n \binom{n}{j} [F(x)]^j [1 - F(x)]^{n-j} \quad (\text{A.103})$$

□

Theorem A.3: Change of Variables

Let X_1, \dots, X_n be continuous random variables having joint density f and let random variables Y_1, \dots, Y_n be defined by

$$Y_i = \sum_{j=1}^n a_{ij} X_j, \quad i = 1, \dots, n, \quad (\text{A.104})$$

where the matrix $A = [a_{ij}]$ has nonzero determinant $\det A$. Then Y_1, \dots, Y_n have joint density f_{Y_1, \dots, Y_n} given by

$$f_{Y_1, \dots, Y_n}(y_1, \dots, y_n) = \frac{1}{|\det A|} f(x_1, \dots, x_n), \quad (\text{A.105})$$

where the x 's are defined in terms of y 's as the unique solution to the equations $y_i = \sum_{j=1}^n a_{ij} x_j$.

A.6 Expectations and the Central Limit Theorem**Definition A.55: Expectation (Continuous)**

Let X be a continuous random variable having density f . We say that X has finite expectation if

$$\int_{-\infty}^{\infty} |x| f(x) dx < \infty, \quad (\text{A.106})$$

and in that case we define its expectation by

$$E[X] = \int_{-\infty}^{\infty} x f(x) dx \quad (\text{A.107})$$

Theorem A.4

Let X_1, \dots, X_n be continuous random variables having joint density f and let Z be a random variable defined in terms of X_1, \dots, X_n be $Z = \phi(X_1, \dots, X_n)$. Then Z has finite expectation if and only if

$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} |\phi(x_1, \dots, x_n)| f(x_1, \dots, x_n) dx_1 \cdots dx_n < \infty \quad (\text{A.108})$$

in which case

$$E[Z] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \phi(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \cdots dx_n < \infty \quad (\text{A.109})$$

Definition A.56: Moments (Continuous)

Let X be a continuous random variable having density f and mean μ . If X has finite m^{th} moment, then we have

$$E[X^m] = \int_{-\infty}^{\infty} x^m f(x) dx \quad (\text{A.110})$$

If X has finite second moment, its variance σ^2 is given by

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx \quad (\text{A.111})$$

Definition A.57: Conditional Expectation

Let X and Y be continuous random variables having joint density f and suppose that Y has finite expectation. Recall that we defined the conditional density of Y given $X = x$ by

$$f_{Y|X}(y|x) = \begin{cases} \frac{f(x,y)}{f_X(x)}, & 0 < f_X(x) < \infty, \\ 0 & \text{elsewhere.} \end{cases} \quad (\text{A.112})$$

For each x such that $0 < f_X(x) < \infty$ the function $f_{Y|X}(y|x)$, $-\infty < y < \infty$, is a density function with respect to **Def A.22**. Thus we can talk about various moments of this density. Its mean is called the *conditional expectation* of Y given $X = x$ and is denoted by $E[Y|X = x]$ or $E[Y|x]$. Thus

$$E[Y|X = x] = \int_{-\infty}^{\infty} y f(y|x) dy \quad (\text{A.113})$$

$$= \frac{\int_{-\infty}^{\infty} y f(x, y) dy}{f_X(x)} \quad (\text{A.114})$$

when $0 < f_X(x) < \infty$. We define $E[Y|X = x] = 0$ elsewhere.

Proposition A.24: Properties of Conditional Expectation

Let X, Y, Z be random variables and $a, b \in \mathbb{R}$, and $g : \mathbb{R} \rightarrow \mathbb{R}$. Assuming all the following expectations exist, we have that

- (i) $E[a|Y] = a$
- (ii) $E[aX + bY|Z] = aE[X|Z] + bE[Y|Z]$
- (iii) $E[X|Y] \geq 0$ if $X \geq 0$.
- (iv) $E[X|Y] = E[X]$ if X and Y are independent.
- (v) $E[E[X|Y]] = E[X]$
- (vi) $E[Xg(Y)|Y] = g(Y)E[X|Y]$. In particular, $E[g(Y)|Y] = g(Y)$.
- (vii) $E[X|Y, g(Y)] = E[X|Y]$
- (viii) $E[E[X|Y, Z]|Y] = E[X|Y]$

Definition A.58: Regression Function

In statistics, the function m defined by $m(x) = E[Y|X = x]$ is called the *regression function* of Y on X .

Lemma A.1

Let X be a random variable with density f_X . Then, $\frac{X-\alpha}{\beta}$ is a random variable with density

$$f_{(X-\alpha)/\beta}(z) = \beta f_X(\beta z + \alpha) \quad (\text{A.115})$$

Theorem A.5: Central Limit Theorem

Let X_1, X_2, \dots be independent, identically distributed random variables having mean μ and finite nonzero variance σ^2 . Set $S_n = X_1 + \dots + X_n$. Then

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq x\right) = \Phi(x), \quad -\infty < x < \infty \quad (\text{A.116})$$

Where we recall that Φ is the CDF for the normal density:

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du \quad (\text{A.117})$$

Proof. Let X_1, X_2, \dots be independent, identically distributed random variables having mean μ and finite nonzero variance σ^2 . Let their density be denoted by f . We define the random variable $S_n = X_n + S_{n-1}$, noting that its density is given by

$$f_{S_n}(x_n) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} dx_1 dx_2 \dots dx_{n-1} \left(\prod_{i=1}^{n-1} f(x_{i+1} - x_i) \right) f(x_1) \quad (\text{A.118})$$

We define a new random variable $G_n := \frac{S_n - n\mu}{\sigma\sqrt{n}}$. Its density is given by

$$f_{G_n}(x_n) = \sigma\sqrt{n} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} dx_1 dx_2 \dots dx_{n-1} f(\sigma\sqrt{n}x_n + n\mu - x_{n-1}) \left(\prod_{i=1}^{n-2} f(x_{i+1} - x_i) \right) f(x_1) \quad (\text{A.119})$$

□

Appendix B

Linear Algebra

Proposition B.1

Let X be a $n \times p$ matrix. Then $(X^T X)^T = X^T X$ and if $X^T X$ is invertible, we then have $((X^T X)^{-1})^T = (X^T X)^{-1}$.

Proof. The first statement is trivial. For the second statement, we observe that

$$\begin{aligned} (X^T X)(X^T X)^{-1} &= 1 \\ ((X^T X)(X^T X)^{-1})^T &= 1 \\ ((X^T X)^{-1})^T (X^T X) &= 1 \\ ((X^T X)^{-1})^T &= (X^T X)^{-1} \end{aligned} \tag{B.1}$$

□

Proposition B.2: Recasting Weighted Sum Into a Matrix Equation

Let $x^{(i)} \in \mathbb{R}^{1 \times n_x}$ be a vector and let $v_i \in \mathbb{R}$ where $1 \leq i \leq m$. Then, we define the $n_x \times m$ matrix X such that $x^{(i)}$ are stacked beside each other in columns^a:

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \tag{B.2}$$

Hence, the matrix elements are given by $X_{ij} = x_i^{(j)}$ with $x_i^{(j)}$ indicating the i^{th} entry of the vector $x^{(j)}$. Similarly, we define a $1 \times m$ vector V by

$$V = [v_1 \quad v_2 \quad \dots \quad v_m] \tag{B.3}$$

We now establish the following identity:

$$B = \sum_{i=1}^m v_i x^{(i)} = X V^T \tag{B.4}$$

Proof. We consider both these objects equivalent if their underlying elements are the same. We observe that

B is a $m \times 1$ vector. Hence, we have that

$$B_j = \sum_{i=1}^m v_i x_j^{(i)}. \quad (\text{B.5})$$

Similarly, we observe that XV^T is a $m \times 1$ vector. We note that its elements are given by

$$\left[XV^T \right]_j = \sum_{i=1}^m X_{ji} V_i = \sum_{i=1}^m x_j^{(i)} v_i = B_j. \quad (\text{B.6})$$

We therefore have the identity:

$$\sum_{i=1}^m v_i x^{(i)} = XV^T. \quad (\text{B.7})$$

□

^aIf $x^{(i)}$ correspond to the input elements of a training set, then X is the training matrix, defined by Def ??.

Definition B.1: Positive / Negative Definite Matrix

Let M be an $n \times n$ symmetric real matrix. Then, we say that M is positive definite iff the scalar

$$x^T M x > 0 \quad \forall x \in \mathbb{R}^n \setminus \mathbf{0}. \quad (\text{B.8})$$

Similarly, we say that M is negative definite iff the scalar

$$x^T M x < 0 \quad \forall x \in \mathbb{R}^n \setminus \mathbf{0} \quad (\text{B.9})$$

Definition B.2: Partial Identity Matrix Properties

Let $\Delta_m^n := \text{diag}(\underbrace{1, 1, \dots, 1}_{n \text{ times}}, \underbrace{0, 0, \dots, 0}_{m \text{ times}})$ be an $(m+n) \times (m+n)$ matrix. In essence, the upper-left $m \times m$ quadrant is an identity matrix and the rest of the block quadrant are zero matrices. Let v be an $(m+n) \times 1$ vector with entries denoted as v_j . Then, we have that

$$\|\Delta_m^n v\|^2 := v^T (\Delta_m^n)^T \Delta_m^n v = \sum_{i=1}^n v_i^2 \quad (\text{B.10})$$

Proposition B.3: Idempotent Matrix Eigenvalues

Let A be an idempotent matrix (i.e it satisfies $A^2 = A$). Then, the only two eigenvalues for A are given by $\lambda = 0, 1$.

Proof. Let $v \in \mathbb{R}^N \setminus \{0\}$ be an eigenvector of A satisfying

$$Av = \lambda v, \quad (\text{B.11})$$

then, we have that

$$Av = A^2 v = \lambda^2 v = \lambda v, \quad (\text{B.12})$$

so that λ must satisfy

$$\lambda^2 - \lambda = 0, \quad (\text{B.13})$$

whose only solutions are given by $\lambda = 0, 1$. □

Theorem B.1: Invertible Matrix Theorem

Let A be a $n \times n$ matrix. A is invertible if and only if any (and hence, all) of the following hold:

1. A is row-equivalent to the $n \times n$ identity matrix \mathbb{I}_n .
2. A has n pivot positions.
3. The equation $A\mathbf{x} = \mathbf{0}$ has only the trivial solution $\mathbf{x} = \mathbf{0}$.
4. The columns of A form a linearly independent set.
5. The linear transformation $x \mapsto Ax$ is one-to-one.
6. For every column vector $b \in \mathbb{R}^n$, the equation $Ax = b$ has a unique solution.
7. The columns of A span \mathbb{R}^n .
8. The linear transformation $x \mapsto Ax$ is a surjection.
9. There is an $n \times n$ matrix C such that $CA = \mathbb{I}_n$.
10. There is an $n \times n$ matrix D such that $AD = \mathbb{I}_n$.
11. The transpose matrix A^T is invertible.
12. The columns of A form a basis for \mathbb{R}^n .
13. The rank of A is n .
14. The null space of A is $\{\mathbf{0}\}$.
15. 0 fails to be an eigenvalue of A .
16. The determinant of A is not zero.
17. The matrix A has n non-zero singular values.

Definition B.3: Adjoint, T^*

Let V and W be two vector spaces and suppose $T \in \mathcal{L}(V, W)$. The **adjoint** of T is the function $T^* : W \rightarrow V$ such that

$$\langle Tv, w \rangle = \langle v, T^*w \rangle \quad (\text{B.14})$$

for every $v \in V$ and every $w \in W$.

Definition B.4: Singular Value Decomposition (SVD)

Let M be an $n \times m$ matrix. Then, a singular value decomposition (SVD) for M takes the form

$$M = U\Sigma V^*, \quad (\text{B.15})$$

where U is a $n \times n$ unitary matrix, V is a $m \times m$ unitary matrix with V^* denoting its conjugate transpose and Σ is a $m \times n$ rectangular diagonal matrix with diagonal entries $\Sigma_{ii} = \sigma_i$. The diagonal entries σ_i are referred to as the *singular values* of M . Typically, one constructs this so that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p$ where

$p := \min\{n, m\}$.

If M contains real entries, then its SVD is given by $M = U\Sigma V^T$ where U and V are orthogonal matrices satisfying

$$U^T U = \mathbb{I}_{m \times m} \tag{B.16}$$

$$V^T V = \mathbb{I}_{n \times n}. \tag{B.17}$$

We note that the singular value decomposition is not unique but is guaranteed to exist for any matrix. We also note that some authors take the SVD as $M = UDV^T$ where U is $n \times m$, D is $m \times m$ and V is $m \times m$.