

David J. Olive

# High Dimensional Statistics: an Asymptotic Viewpoint

July 31, 2024





# Preface

Many statistics departments offer a one semester graduate course in high dimensional statistics using texts such as Bühlmann and van de Geer (2011), Giraud (2022), Lederer (2022), or Wainwright (2019). Statistical learning texts are also used. See Hastie et al. (2009), Hastie et al. (2015), and James et al. (2021). Also see Fujikoshi, Ulyanov, and Shimizu (2010), Koch (2014), Olive (2023e), and Rish and Grabarnik (2015).

High dimensional statistics are used when  $n < 5p$  where  $n$  is the sample size and  $p$  is the number of predictors  $p$ . Consider the multiple linear regression model  $Y_i = \alpha + \mathbf{x}_i^T \boldsymbol{\beta} + e_i = \alpha + x_{i1}\beta_1 + \cdots + x_{ip}\beta_p + e_i$  for  $i = 1, \dots, n$ . Let the full model use all  $p$  predictors with  $\boldsymbol{\beta} = \boldsymbol{\beta}_F$ . In low dimensions where  $n \geq 10p$ , often  $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \stackrel{D}{\rightarrow} N_p(\mathbf{0}, \boldsymbol{\Sigma})$  where  $\boldsymbol{\Sigma}$  is estimated by  $\hat{\boldsymbol{\Sigma}} = \hat{\sigma}^2 \hat{\mathbf{C}}^{-1}$  where the errors  $e_i$  have variance  $V(E_i) = \sigma^2$  and where the inverse matrix  $\hat{\mathbf{C}}^{-1}$  does not exist if  $p > n$ . Much of the high dimensional literature seeks bounds on the Euclidean norm  $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|$ . However, if  $\hat{\boldsymbol{\beta}}$  is a  $\sqrt{n}$  consistent estimator of  $\boldsymbol{\beta}_F$ , then  $\hat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i$  is proportional to  $1/\sqrt{n}$ . Hence  $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|^2$  is proportional to  $p/n$  which tends to be large when  $p \gg n$ . Similar results hold for estimators  $\hat{\boldsymbol{\theta}}$  of  $\boldsymbol{\theta}$  for statistical models that depend on a  $p \times 1$  vector of parameters  $\boldsymbol{\theta}$ . Often the high dimensional literature imposes regularity conditions, **that are much too strong**, to force  $\|\hat{\boldsymbol{\beta}}_F - \boldsymbol{\beta}_F\|$  to be small as both  $n$  and  $p \rightarrow \infty$ .

This text uses large sample theory = asymptotic theory to justify many of the methods used in the text. Several dimension reduction techniques are used. One technique is to use data splitting and variable selection to choose a model  $I$  with  $k$  predictors where  $n \geq 10k$ , and then apply the standard low dimensional inference on the resulting model. This changes the high dimensional problem into a low dimensional problem. Sometimes we use the strong assumption that the cases  $(\mathbf{x}_i, Y_i)^T$  are independent and identically distributed (iid). Then variable selection methods often work because the conditional distribution  $Y|\mathbf{x}_I^T \boldsymbol{\beta}_I$  has much more information than the marginal distribution for  $Y$ .

A second technique is to use large sample theory such that  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{D} N_p(\mathbf{0}, \boldsymbol{\Sigma})$  where  $\boldsymbol{\Sigma}$  is estimated by  $\hat{\boldsymbol{\Sigma}} = \hat{\mathbf{C}}$  where the inverse matrix  $\hat{\mathbf{C}}^{-1}$  is not used. Then tests and confidence intervals for quantities that only use a few of the parameters, such as  $\theta_i$  or  $\theta_i - \theta_k$  can be derived. Hence low dimensional quantities are tested.

A third technique is to replace  $\boldsymbol{\theta}$  by the norm  $\|\boldsymbol{\theta}\|$  or  $\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2$  by the norm  $\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|$ , reducing the  $p$ -dimensional problem of testing  $H_0 : \boldsymbol{\theta} = \mathbf{0}$  or  $H_0 : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$  to the one-dimensional problem of testing  $H_0 : \|\boldsymbol{\theta}\| = 0$  or  $H_0 : \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| = 0$ .

The prerequisite for this text is a calculus based course in statistics at the level of Chihara and Hesterberg (2011), Hogg, Tanis, and Zimmerman (2020), Larsen and Marx (2011), Wackerly, Mendenhall and Scheaffer (2008) or Walpole, Myers, Myers and Ye (2016). Linear algebra and one computer programming class are essential. Knowledge of regression would be useful. See Olive (2017a) and Cook and Weisberg (1999). Knowledge of multivariate analysis would be useful. See Olive (2017b) and Johnson and Wichern (2007).

Some highlights of this text follow.

- Prediction intervals are given that can be useful even if  $n < p$ .
- The response plot is useful for checking the model.
- The large sample theory for the elastic net, lasso, and ridge regression is greatly simplified.
- The large sample theory for some data splitting estimators, variable selection estimators, marginal maximum likelihood estimators, and one component partial least squares will be given. See Olive and Zhang (2023), Olive et al. (2024), and Rathnayake and Olive (2023).

**Downloading the book's R functions** *slpack.txt* and data files *sl-data.txt* into R: The commands

```
source("http://parker.ad.siu.edu/Olive/hdpack.txt")
source("http://parker.ad.siu.edu/Olive/hddata.txt")
```

The R software is used in this text. See R Core Team (2020). Some packages used in the text include `glmnet` Friedman et al. (2015), `leaps` Lumley (2009), `MASS` Venables and Ripley (2010), and `pls` Mevik et al. (2015).

### Acknowledgements

Teaching the material to Math 583 students at Southern Illinois University in 2023 was very useful. Trevor Hastie's website had a lot of useful information. Work by R. Dennis Cook and his coauthors was useful for figuring out OPLS.

# Contents

<b>1</b>	<b>Introduction</b> .....	1
1.1	Overview .....	1
1.2	Response Plots and Response Transformations .....	5
1.2.1	Response and Residual Plots .....	5
1.2.2	Response Transformations .....	8
1.3	The Multivariate Normal Distribution .....	13
1.4	Outlier Detection .....	16
1.4.1	The Location Model .....	17
1.4.2	Outlier Detection with Mahalanobis Distances .	18
1.4.3	Outlier Detection if $p > n$ .....	22
1.5	Large Sample Theory .....	28
1.5.1	The CLT and the Delta Method .....	29
1.5.2	Modes of Convergence and Consistency .....	32
1.5.3	Slutsky's Theorem and Related Results .....	39
1.5.4	Multivariate Limit Theorems .....	42
1.6	Mixture Distributions .....	47
1.7	A Review of Multiple Linear Regression .....	48
1.7.1	The ANOVA F Test .....	52
1.7.2	The Partial F Test .....	56
1.7.3	The Wald t Test .....	59
1.7.4	The OLS Criterion .....	60
1.7.5	The No Intercept MLR Model .....	63
1.8	Summary .....	64
1.9	Complements .....	68
1.10	Problems .....	68
<b>2</b>	<b>Multiple Linear Regression</b> .....	77
2.1	The MLR Model .....	77
2.2	Forward Selection .....	88
2.3	Principal Components Regression .....	91
2.4	Partial Least Squares .....	96

2.5	<b>Ridge Regression</b> .....	98
2.6	<b>Lasso</b> .....	106
2.7	<b>Lasso Variable Selection</b> .....	111
2.8	<b>The Elastic Net</b> .....	114
2.9	<b>OPLS</b> .....	117
2.10	<b>The MMLE</b> .....	119
2.11	<b><math>k</math>-Component Regression Estimators</b> .....	120
2.12	<b>Prediction Intervals</b> .....	122
2.13	<b>Cross Validation</b> .....	126
2.14	<b>Hypothesis Testing after Model Selection, <math>n/p</math> Large</b> .	131
2.15	<b>What if <math>n</math> is not <math>\gg p</math>?</b> .....	132
	2.15.1 <b>Sparse Models</b> .....	133
2.16	<b>Data Splitting</b> .....	134
2.17	<b>The Multitude of MLR Models</b> .....	136
2.18	<b>Summary</b> .....	136
2.19	<b>Complements</b> .....	141
2.20	<b>Problems</b> .....	146
<b>3</b>	<b>MLR with Heterogeneity</b> .....	155
	3.1 <b>OLS Large Sample Theory</b> .....	155
	3.2 <b>Bootstrap Methods and Sandwich Estimators</b> .....	156
	3.3 <b>Simulations</b> .....	158
	3.4 <b>OPLS in Low and High Dimensions</b> .....	160
	3.5 <b>Summary</b> .....	160
	3.6 <b>Complements</b> .....	160
	3.7 <b>Problems</b> .....	160
<b>4</b>	<b>Binary Regression</b> .....	161
	4.1 <b>Two Set Inference</b> .....	161
	4.2 <b>Summary</b> .....	161
	4.3 <b>Complements</b> .....	161
	4.4 <b>Problems</b> .....	161
<b>5</b>	<b>Poisson Regression</b> .....	163
	5.1 <b>Two Set Inference</b> .....	163
	5.2 <b>Summary</b> .....	163
	5.3 <b>Complements</b> .....	163
	5.4 <b>Problems</b> .....	163
<b>6</b>	<b>Other Regression Models</b> .....	165
	6.1 <b>Two Set Inference</b> .....	165
	6.2 <b>Summary</b> .....	165
	6.3 <b>Complements</b> .....	165
	6.4 <b>Problems</b> .....	165

<b>7</b>	<b>One and Two Sample Tests</b> .....	167
	7.1 <b>Two Set Inference</b> .....	167
	7.2 <b>Summary</b> .....	167
	7.3 <b>Complements</b> .....	167
	7.4 <b>Problems</b> .....	167
<b>8</b>	<b>Classification</b> .....	169
	8.1 <b>Introduction</b> .....	169
	8.2 <b>LDA and QDA</b> .....	171
	8.2.1 <b>Regularized Estimators</b> .....	174
	8.3 <b>LR</b> .....	174
	8.4 <b>KNN</b> .....	176
	8.5 <b>Some Matrix Optimization Results</b> .....	178
	8.6 <b>FDA</b> .....	180
	8.7 <b>Estimating the Test Error</b> .....	186
	8.8 <b>Some Examples</b> .....	189
	8.9 <b>Classification Trees, Bagging, and Random Forests</b> ...	192
	8.9.1 <b>Pruning</b> .....	195
	8.9.2 <b>Bagging</b> .....	196
	8.9.3 <b>Random Forests</b> .....	197
	8.10 <b>Support Vector Machines</b> .....	197
	8.10.1 <b>Two Groups</b> .....	197
	8.10.2 <b>SVM With More Than Two Groups</b> .....	200
	8.11 <b>Summary</b> .....	200
	8.12 <b>Complements</b> .....	204
	8.13 <b>Problems</b> .....	205
<b>9</b>	<b>Multivariate Linear Regression</b> .....	213
	9.1 <b>Introduction</b> .....	213
	9.2 <b>Plots for the Multivariate Linear Regression Model</b> ..	217
	9.3 <b>Asymptotically Optimal Prediction Regions</b> .....	220
	9.4 <b>Testing Hypotheses</b> .....	225
	9.5 <b>An Example and Simulations</b> .....	235
	9.5.1 <b>Simulations for Testing</b> .....	240
	9.6 <b>The Robust <code>rmreg2</code> Estimator</b> .....	243
	9.7 <b>Bootstrap</b> .....	246
	9.7.1 <b>Parametric Bootstrap</b> .....	246
	9.7.2 <b>Residual Bootstrap</b> .....	246
	9.7.3 <b>Nonparametric Bootstrap</b> .....	247
	9.8 <b>Data Splitting</b> .....	247
	9.9 <b>Ridge Regression, PCR, and Other High Dimensional Methods</b> .....	247
	9.10 <b>Summary</b> .....	248
	9.11 <b>Complements</b> .....	254
	9.12 <b>Problems</b> .....	255

<b>10</b>	<b>Multivariate Analysis</b> .....	261
	10.1 <b>Two Set Inference</b> .....	261
	10.2 <b>Summary</b> .....	261
	10.3 <b>Complements</b> .....	261
	10.4 <b>Problems</b> .....	261
<b>11</b>	<b>Stuff for Students</b> .....	263
	11.1 <b>R</b> .....	263
	11.2 <b>Hints for Selected Problems</b> .....	266
	11.3 <b>Projects</b> .....	267
	11.4 <b>Tables</b> .....	270
	<b>Index</b> .....	281



# Chapter 1

## Introduction

This chapter provides a preview of the book, and some techniques useful for visualizing data in the background of the data are given in Section 1.2. Sections 1.3 and 1.7 review the multivariate normal distribution and multiple linear regression. Section 1.4 suggests methods for outlier detection. Some large sample theory is presented in Section 1.5, and Section 1.6 covers mixture distributions.

### 1.1 Overview

Statistical Learning could be defined as the statistical analysis of multivariate data. Machine learning, data mining, analytics, business analytics, data analytics, and predictive analytics are synonymous terms. The techniques are useful for Data Science and Statistics, the science of extracting information from data. The *R* software will be used. See R Core Team (2020).

Let  $\mathbf{z} = (z_1, \dots, z_k)^T$  where  $z_1, \dots, z_k$  are  $k$  random variables. Often  $\mathbf{z} = (Y, \mathbf{x}^T)^T$  where  $\mathbf{x}^T = (x_1, \dots, x_p)$  is the vector of predictors and  $Y$  is the variable of interest, called a response variable. Predictor variables are also called independent variables, covariates, or features. The response variable is also called the dependent variable. Usually context will be used to decide whether  $\mathbf{z}$  is a random vector or the observed random vector.

**Definition 1.1.** A **case** or **observation** consists of  $k$  random variables measured for one person or thing. The  $i$ th case  $\mathbf{z}_i = (z_{i1}, \dots, z_{ik})^T$ . The **training data** consists of  $\mathbf{z}_1, \dots, \mathbf{z}_n$ . A statistical model or method is fit (trained) on the training data. The **test data** consists of  $\mathbf{z}_{n+1}, \dots, \mathbf{z}_{n+m}$ , and the test data is often used to evaluate the quality of the fitted model.

Following James et al. (2013, p. 30), the previously unseen test data is not used to train the Statistical Learning method, but interest is in how well the

method performs on the test data. If the training data is  $(\mathbf{x}_1, Y_1), \dots, (\mathbf{x}_n, Y_n)$ , and the previously unseen test data is  $(\mathbf{x}_f, Y_f)$ , then particular interest is in the accuracy of the estimator  $\hat{Y}_f$  of  $Y_f$  obtained when the Statistical Learning method is applied to the predictor  $\mathbf{x}_f$ . The two Pelawa Watagoda and Olive (2021b) prediction intervals, developed in Section 2.2, will be tools for evaluating Statistical Learning methods for the additive error regression model  $Y_i = m(\mathbf{x}_i) + e_i = E(Y_i|\mathbf{x}_i) + e_i$  for  $i = 1, \dots, n$  where  $E(W)$  is the expected value of the random variable  $W$ . The multiple linear regression (MLR) model,  $Y_i = \beta_1 + x_2\beta_2 + \dots + x_p\beta_p + e = \mathbf{x}^T\boldsymbol{\beta} + e$ , is an important special case. Olive, Rathnayake, and Haile (2022) give prediction intervals for parametric regression models such as generalized linear models (GLMs), generalized additive models (GAMs), and some survival regression models.

The estimator  $\hat{Y}_f$  is a *prediction* if the response variable  $Y_f$  is continuous, as occurs in regression models. If  $Y_f$  is categorical, then  $\hat{Y}_f$  is a *classification*. For example, if  $Y_f$  can be 0 or 1, then  $\mathbf{x}_f$  is classified to belong to group  $i$  if  $\hat{Y}_f = i$  for  $i = 0$  or  $1$ .

Following Marden (2006, pp. 5,6), the focus of *supervised learning* is predicting a future value of the response variable  $Y_f$  given  $\mathbf{x}_f$  and the training data  $(Y_1, \mathbf{x}_1), \dots, (Y_n, \mathbf{x}_1)$ . Hence the focus is not on hypothesis testing, confidence intervals, parameter estimation, or which model fits best, although these four inference topics can be useful for better prediction. The focus of *unsupervised learning* is to group  $\mathbf{x}_1, \dots, \mathbf{x}_n$  into clusters. *Data mining* is looking for relationships in large data sets.

**Notation:** Typically lower case boldface letters such as  $\mathbf{x}$  denote column vectors, while upper case boldface letters such as  $\mathbf{S}$  or  $\mathbf{Y}$  are used for matrices or column vectors. If context is not enough to determine whether  $\mathbf{y}$  is a random vector or an observed random vector, then  $\mathbf{Y} = (Y_1, \dots, Y_p)^T$  may be used for the random vector, and  $\mathbf{y} = (y_1, \dots, y_p)^T$  for the observed value of the random vector. An upper case letter such as  $Y$  will usually be a random variable. A lower case letter such as  $x_1$  will also often be a random variable. An exception to this notation is the generic multivariate location and dispersion estimator  $(T, \mathbf{C})$  where the location estimator  $T$  is a  $p \times 1$  vector such as  $T = \bar{\mathbf{x}}$ .  $\mathbf{C}$  is a  $p \times p$  dispersion estimator and conforms to the above notation.

The main focus of the first three chapters is developing tools to analyze the multiple linear regression (MLR) model  $Y_i = \mathbf{x}_i^T\boldsymbol{\beta} + e_i$  for  $i = 1, \dots, n$ . Classical regression techniques use (ordinary) least squares (OLS) and assume  $n \gg p$ , but Statistical Learning methods often give useful results if  $p \gg n$ . OLS forward selection, lasso, ridge regression, marginal maximum likelihood (MMLE), one component partial least squares (OPLS), the elastic net, partial least squares (PLS), and principal component regression (PCR) will be some of the techniques examined. See Chapter 3.