

A Partial Solution Manual for: *The Elements of
Statistical Learning* by Jerome Friedman, Trevor
Hastie, and Robert Tibshirani

Wenhao Wu

wnhwu@ucdavis.edu

Dept. ECE, UC Davis

February 5, 2016

Contents

| | |
|---|----|
| Preface | 2 |
| Acknowledgment | 3 |
| 2 Overview of Supervised Learning | 4 |
| 3 Linear Methods for Regression | 5 |
| 4 Linear Methods for Classification | 6 |
| 5 Basis Expansions and Regularization | 7 |
| 6 Kernel Smoothing Methods | 8 |
| 7 Model Assessment and Selection | 9 |
| 8 Model Inference and Averaging | 10 |
| 9 Additive Models, Trees, and Related Methods | 11 |
| 10 Boosting and Additive Trees | 12 |
| 11 Neural Networks | 13 |
| 12 Support Vector Machines and Flexible Discriminants | 15 |
| 13 Prototype Methods and Nearest-Neighbors | 21 |
| 14 Unsupervised Learning | 22 |
| 15 Random Forests | 23 |
| 16 Ensemble Learning | 24 |
| 17 Undirected Graphical Models | 25 |
| 18 High-Dimensional Problems | 26 |

Preface

This work is expected to be used as a supplementary material for Weatherwax and Epstein's solution manual [?], which I found to be very helpful when self-studying this popular textbook. The numbering of chapters and problems are based on the 2nd edition (10th printing with corrections, Jan 2013) available online [?].

The author was not able to solve all the excercises. Even for the solutions included we expect many mistakes and shortcomings. It would be of great help if people could suggest possible solutions or help us find and correct the errors so this solution manual can be continuously improved to benefit more interested readers. We are also open to all comments and criticisms. Our contact information can be found at the website holding this draft [?].

Acknowledgment

Chapter 2

Overview of Supervised Learning

Chapter 3

Linear Methods for Regression

Chapter 4

Linear Methods for Classification

Chapter 5

Basis Expansions and Regularization

Chapter 6

Kernel Smoothing Methods

Chapter 7

Model Assessment and Selection

Chapter 8

Model Inference and Averaging

Chapter 9

Additive Models, Trees, and Related Methods

Chapter 10

Boosting and Additive Trees

Chapter 11

Neural Networks

Ex. 11.1

In (11.5), set $K = 1$, $g_1(T) = T$, we have

$$f_1(X) = \beta_{01} + \beta_1^T Z = \beta_{01} + \sum_{m=1}^M \beta_{m1} \sigma(\alpha_{0m} + \alpha_m^T X) \quad (11.1)$$

The correspondence between (11.1) and (11.5) becomes clearer, as enumerated in Table 11.1

Table 11.1: Correspondence between the project pursuit regression and the neural network

| (11.1) | (11.5) |
|--------------|---|
| ω_m | α_m |
| $g_m(\cdot)$ | $\beta_{01}, \beta_{m1} \sigma(\alpha_{0m} + \alpha_m^T X)$ |

Ex. 11.2

$$\frac{\partial f}{\partial X} = \sum_{m=1}^M \beta_m [\sigma(\cdot)(\sigma(\cdot) - 1)] \alpha_m \quad (11.2)$$

$$\frac{\partial^2 f}{\partial X \partial X^T} = \sum_{m=1}^M \beta_m [(2\sigma(\cdot) - 1)(\sigma(\cdot) - 1)\sigma(\cdot)] \alpha_m \alpha_m^T \quad (11.3)$$

Since $\sigma(\alpha_{0m} + \alpha_m^T X) \approx 1/2$ when $\alpha_{0m} \approx 0$ and $\alpha_m \approx 0$, therefore $\frac{\partial^2 f}{\partial X \partial X^T} \approx 0$, i.e. the resulting model is nearly linear.

Ex. 11.3

$$R(\theta) = - \sum_{i=1}^N R_i(\theta) = - \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log g_j(T) \quad (11.4)$$

Note that different from regression, each softmax function $g_j(T)$, $j = 1, \dots, K$ is a function

of all T_1, \dots, T_K .

$$\frac{\partial R_i}{\partial \beta_{km}} = - \sum_{j=1}^K \frac{y_{ij}}{g_j} \frac{\partial g_j}{\partial T_k} z_{mi} = \delta_{ki} z_{mi} \quad (11.5a)$$

$$\begin{aligned} \frac{\partial R_i}{\partial \alpha_{ml}} &= - \sum_{j=1}^K \frac{y_{ij}}{g_j} \sum_{k=1}^K \frac{\partial g_j}{\partial T_k} \beta_{km} \sigma'(\alpha_m^T x_i) x_{il} \\ &= \left[\sigma'(\alpha_m^T x_i) \sum_{k=1}^K \beta_{km} \delta_{ki} \right] x_{il} = s_{mi} x_{il} \end{aligned} \quad (11.5b)$$

It is noted that

$$\frac{1}{g_j} \frac{\partial g_j}{\partial T_k} = \begin{cases} 1 - g_j & j = k \\ -g_k / \exp(T_j) & j \neq k \end{cases} \quad (11.6)$$

As a result, although $g_j(T)$ depends on all T_1, \dots, T_K , $(\partial g_j / \partial T_k) / g_j$ can still be locally evaluated and propagated downward over the link (T_k, g_j) . Consequently, the forward and backward propagation equations are pretty much the same as those for the square error loss function. In the forward pass for record x_i , $i = 1, \dots, N$, the weights β_{km} and α_{ml} are fixed and the predicted $\hat{g}_j(T_i)$ are evaluated. In the backward pass, $(y_{ij}/g_j)(\partial g_j / \partial T_k)$ are evaluated and propagated to T_k , where δ_{ki} is computed, and then back-propagated to give s_{mi} at Z_m . Then the gradients are evaluated as in Eq. (11.5). The gradient descent update is exactly the same as (11.13).

Ex. 11.4

If the network has no hidden layer, we have

$$g_j(x) = \frac{\exp(T_j)}{\sum_{k=1}^K \exp(T_k)} = \frac{\exp(\beta_j^T x)}{\sum_{k=1}^K \exp(\beta_k^T x)}, \quad (11.7)$$

exactly the same as the multinomial logistic model.

Ex. 11.5 (Program)

Ex. 11.6 (Program)

Ex. 11.7 (Program)

Chapter 12

Support Vector Machines and Flexible Discriminants

Ex. 12.1

Firstly, we prove that for (12.8), the optimal solution must satisfy $\hat{\xi}_i = [1 - y_i(x_i^T \hat{\beta} + \hat{\beta}_0)]_+$. To see this, from the constraints in (12.8), we have $\hat{\xi}_i \geq [1 - y_i(x_i^T \hat{\beta} + \hat{\beta}_0)]_+$. Assume for contradiction that $\exists i$ such that $\hat{\xi}_i > [1 - y_i(x_i^T \hat{\beta} + \hat{\beta}_0)]_+$, then setting $\hat{\xi}_i \leftarrow [1 - y_i(x_i^T \hat{\beta} + \hat{\beta}_0)]_+$ results in smaller objective in (12.8), which is in contradiction to the fact that $\hat{\xi}_i$ is from an optimal solution.

On the other hand, $\xi_i = [1 - y_i(x_i^T \beta + \beta_0)]_+ \Rightarrow \xi_i \geq 0, y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i$. Therefore, the solution to (12.8) is the same as

$$\min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i \quad (12.1)$$

$$\text{s.t. } \xi_i = [1 - y_i(x_i^T \beta + \beta_0)]_+, \forall i \quad (12.2)$$

which is exactly the same as (12.25).

Ex. 12.2

Define kernel $K(a, b) = \sum_{j=1}^p a_j b_j$, i.e. $\psi_j(x) = x_j, \gamma_j = 1$ for $j = 1, \dots, p$. Consequently, $g(x) = \sum_{j=1}^p \beta_j x_j \Leftrightarrow g(x) \in \mathcal{H}_K$. Consequently,

$$(12.25) \Leftrightarrow \min_{g, \beta_0} \sum_{i=1}^N [1 - y_i(g(x_i) + \beta_0)]_+ + \frac{\lambda}{2} \|g\|_{\mathcal{H}_K}^2 \quad (12.3)$$

Denote $L(y_i, g(x_i); \beta_0) = [1 - y_i(g(x_i) + \beta_0)]_+ = L_i(\beta_0)$, then

$$(12.25) \Leftrightarrow \min_{\beta_0} \left\{ \min_g \sum_{i=1}^N L_i(\beta_0) + \frac{\lambda}{2} \|g\|_{\mathcal{H}_K}^2 \right\}. \quad (12.4)$$

where the inner min must have a solution in the form of $g(x) = \sum_{i=1}^N \alpha_i K(x, x_i)$ as per

(5.50)(5.51), and we have $\|g\|_{\mathcal{H}_K}^2 = \alpha^T K \alpha$. Therefore

$$(12.25) \Leftrightarrow \min_{\beta_0} \left\{ \min_{\alpha} \sum_{i=1}^N [1 - y_i (\sum_{j=1}^N \alpha_j K(x_j, x_i) + \beta_0)]_+ + \frac{\lambda}{2} \alpha^T K \alpha \right\} \quad (12.5)$$

$$\Leftrightarrow \min_{\beta_0, \alpha} \sum_{i=1}^N [1 - y_i (\sum_{j=1}^N \alpha_j K(x_j, x_i) + \beta_0)]_+ + \frac{\lambda}{2} \alpha^T K \alpha \quad (12.6)$$

Ex. 12.3

Similar to Ex. (12.2). Denote $g(x) = \sum_{m=1}^M \beta_m h_m(x)$. Without penalting the constant term, we have

$$H(\beta, \beta_0) = \sum_{i=1}^N V(y_i - \beta_0 - g(x_i)) + \frac{\lambda}{2} \sum_{m=1}^M \beta_m^2 \quad (12.7)$$

Again we break the minimization problem into 2 steps:

$$\min_{\beta_0, \beta} H(\beta, \beta_0) = \min_{\beta_0} \left\{ \min_{\beta | \beta_0} H(\beta, \beta_0) \right\} \quad (12.8)$$

Consider saquare error loss $V(r) = r^2$, the inner min problem is in the form of

$$\min_{\beta} \sum_{i=1}^N (y_i - \beta_0 - H\beta)^2 + \frac{\lambda}{2} \beta^T \beta \quad (12.9)$$

$$\Leftrightarrow \min_{\alpha} \|\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{K}\alpha\|_F^2 + \frac{\lambda}{2} \alpha^T \mathbf{K}\alpha \quad (12.10)$$

whose solution is $\hat{\alpha} = (\mathbf{K} + \lambda \mathbf{I}/2)^{-1} \mathbf{y}_{\beta_0}$, $\mathbf{y}_{\beta_0} = \mathbf{y} - \beta_0 \mathbf{1}$. Consequently, the outer min problem w.r.t β_0 is in the form of

$$\min_{\beta_0} \mathbf{y}_{\beta_0}^T [\mathbf{I} - (\mathbf{K} + \lambda \mathbf{I}/2)^{-1} \mathbf{K}] \mathbf{y}_{\beta_0} \quad (12.11)$$

which is a quadratic problem.

Ex. 12.4

(a)

$$\text{Left} = (x - \bar{x}_k)^T U U^T (x - \bar{x}_k) - (x - \bar{x}_{k'})^T U U^T (x - \bar{x}_{k'}) \quad (12.12)$$

where $U = W^{-1/2} V^*$, the L columns of V^* are the eigen vectors of $B^* = (W^{-1/2})^T B W^{-1/2}$,

where B is the between-class covariance.

$$\text{Right} = (x - \bar{x}_k)^T W^{-1} (x - \bar{x}_k) - (x - \bar{x}_{k'})^T W^{-1} (x - \bar{x}_{k'}) \quad (12.13)$$

Consequently,

$$\begin{aligned} & \text{Left} - \text{Right} \\ &= 2(\bar{x}_k - \bar{x}_{k'})^T (W^{-1} - UU^T)x + (\bar{x}_k - \bar{x}_{k'})^T [W^{-1} - UU^T](\bar{x}_k + \bar{x}_{k'}) \end{aligned} \quad (12.14)$$

$$\begin{aligned} &= 2(\bar{x}_k - \bar{x}_{k'})^T W^{-1/2} (I - V^*(V^*)^T) (W^{-1/2})^T x \\ &+ (\bar{x}_k - \bar{x}_{k'})^T W^{-1/2} (I - V^*(V^*)^T) (W^{-1/2})^T (\bar{x}_k + \bar{x}_{k'}) \end{aligned} \quad (12.15)$$

Since $(\bar{x}_k - \bar{x}_{k'})^T \in R(M)$ (row space), $(\bar{x}_k - \bar{x}_{k'})^T W^{-1/2} \in R(M^*)$, therefore $(W^{-1/2})^T (\bar{x}_k - \bar{x}_{k'}) \in C(V^*)$ (column space). Therefore

$$(\bar{x}_k - \bar{x}_{k'})^T W^{-1/2} (I - V^*(V^*)^T) = 0 \quad (12.16)$$

thus Left = Right.

(b) ???

Ex. 12.5 (Program)

Ex. 12.6

(a) The i -th row of $\mathbf{Y}\theta$ is

$$(\mathbf{Y}\theta)_i = \sum_{j=1}^K 1(Y_{ij} = 1)\theta_j = \theta(g_i) \quad (12.17)$$

(since there are exactly one j where $Y_{ij} = 1$ for each i). The i -th row of $\mathbf{H}\beta$ is

$$(\mathbf{H}\beta)_i = \sum_{j=1}^K \beta_j h_j(x_i) \quad (12.18)$$

therefore

$$\sum_{i=1}^N (\theta(g_i) - \beta^T h(x_i))^2 = \|\mathbf{Y}\theta - \mathbf{H}\beta\|^2 \quad (12.19)$$

(b) According to the definition, $(\mathbf{D}_\pi)_{kk}$ is the empirical frequency of class k , and θ_k is the score for class k . $\theta^T \mathbf{D}_\pi \mathbf{1} = 0$ implies that the average score over the N records is 0; $\theta^T \mathbf{D}_\pi \theta = 1$ means the variance of the over the N records is 1.

(c) Fixing θ the optimal β is

$$\hat{\beta} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{Y} \theta \quad (12.20)$$

therefore (12.65) can be rewritten as

$$\min_{\theta} \|(\mathbf{I} - \mathbf{S}) \mathbf{Y} \theta\|^2 \Leftrightarrow \min_{\theta} \theta^T \mathbf{Y}^T \mathbf{Y} \theta - \theta^T \mathbf{Y}^T \mathbf{S} \mathbf{Y} \theta \quad (12.21)$$

where $\mathbf{S} = \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$. Since $\theta^T \mathbf{Y}^T \mathbf{Y} \theta = N$, this minimization is equivalent to

$$\max_{\theta} \theta^T \mathbf{Y}^T \mathbf{S} \mathbf{Y} \theta \quad (12.22)$$

(d) Suppose that the SVD of $\mathbf{H} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, then $\mathbf{S} = \mathbf{U} \mathbf{U}^T$ where \mathbf{U} is a N -by- l orthonormal matrix. Therefore \mathbf{S} has L eigenvalues of 1 and $N - L$ eigenvalues of 0. Since constant function is included in h_j , $\mathbf{H} \neq 0$, therefore $L > 0$, so the largest eigenvalue is 1.

(e) (12.53) can be rewritten as

$$ASR = \frac{1}{N} \|\mathbf{Y} \mathbf{\Theta} - \mathbf{H} \mathbf{B}\|_F^2 \quad (12.23)$$

Similar to (c) the solution is the same as

$$\max_{\mathbf{\Theta}} \text{tr}\{\mathbf{\Theta}^T \mathbf{Y}^T \mathbf{S} \mathbf{Y} \mathbf{\Theta}\} \quad (12.24)$$

$$\text{s.t. } \mathbf{\Theta}^T \mathbf{Y}^T \mathbf{Y} \mathbf{\Theta} = \mathbf{I} \quad (12.25)$$

Therefore $\mathbf{Y} \mathbf{\Theta}$ are the K largest eigenvectors of \mathbf{S} .

Ex. 12.7

The penalized optimal scoring problem is in the form of

$$\min_{\mathbf{\Theta}, \mathbf{B}} \|\mathbf{Y} \mathbf{\Theta} - \mathbf{H} \mathbf{B}\|_F^2 + \lambda \text{tr}(\mathbf{B}^T \mathbf{\Omega} \mathbf{B}) \quad (12.26)$$

Given $\mathbf{\Theta}$, the optimal \mathbf{B} is

$$\hat{\mathbf{B}} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{\Omega})^{-1} \mathbf{H}^T \mathbf{Y} \mathbf{\Theta} \quad (12.27)$$

Substitute into Eq. (12.26), we have

$$\min_{\mathbf{\Theta}} \text{tr}(\mathbf{\Theta}^T \mathbf{Y}^T \mathbf{Y} \mathbf{\Theta} - \mathbf{\Theta}^T \mathbf{Y}^T \mathbf{S} \mathbf{Y} \mathbf{\Theta}) \quad (12.28)$$

$$\text{s.t. } \mathbf{\Theta}^T \mathbf{D}_{\pi} \mathbf{\Theta} = \mathbf{I} \quad (12.29)$$

where $\mathbf{S} = \mathbf{H}(\mathbf{H}^T \mathbf{H} + \lambda \mathbf{\Omega})^{-1} \mathbf{H}^T$. Therefore $\mathbf{Y}\mathbf{\Theta}$ are still the eigenvectors of \mathbf{S} .

Ex. 12.8

I found the proof to this problem on [?]. I am trying to follow it the best I can and here is my interpretation. Assuming that $\bar{x} = 0$. We first perform the generalized SVD:

$$(\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \mathbf{U} \mathbf{D} \mathbf{V}^T \quad (12.30)$$

$$\text{s.t. } \mathbf{U}^T \mathbf{Y}^T \mathbf{Y} \mathbf{U} = \mathbf{I}, \mathbf{V}^T \mathbf{X}^T \mathbf{X} \mathbf{V} = \mathbf{I} \quad (12.31)$$

Later we will show that both β_l and v_l are proportional to the columns of \mathbf{V} . From the GSVD, \mathbf{U} and \mathbf{V} satisfy the following 2 equations:

$$\mathbf{U}^T \mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \mathbf{U} = \mathbf{D}^2 \quad (12.32a)$$

$$\mathbf{V}^T \mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \mathbf{V} = \mathbf{D}^2 \quad (12.32b)$$

of which the proof is trivial. First we show that the LDA's discriminant directions v_l are parallel to the columns of \mathbf{V} :

Proposition 12.1. *For the LDA problem (Fisher)*

$$\max_{\mathbf{A}} \text{tr}(\mathbf{A}^T \mathbf{B} \mathbf{A}), \text{ s.t. } \mathbf{A}^T \mathbf{W} \mathbf{A} = \mathbf{I} \quad (12.33)$$

where

$$\mathbf{B} = \mathbf{X}^T \mathbf{Y} (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{X} \quad (12.34a)$$

$$\mathbf{W} = \mathbf{T} - \mathbf{B} \quad (12.34b)$$

$$\mathbf{T} = \mathbf{X}^T \mathbf{X} \quad (12.34c)$$

are the between-class, within-class and total variance (up to normalization), the solution is

$$\hat{\mathbf{A}} = \mathbf{V}(\mathbf{I} - \mathbf{D}^2)^{-1/2} \quad (12.35)$$

Proof. From Eq. (12.32b) and the second constraint of the GSVD, it is easy to see $\hat{\mathbf{A}}^T \mathbf{W} \hat{\mathbf{A}} = \mathbf{I}$. On the other hand, $\hat{\mathbf{A}}$ diagonalizes \mathbf{B} by $\hat{\mathbf{A}}^T \mathbf{B} \hat{\mathbf{A}} = (\mathbf{I} - \mathbf{D}^2)^{-1} \mathbf{D}^2$. \square

Next we show that the β_l from optimal scoring are also parallel to the columns of \mathbf{V}

Proposition 12.2. *The optimal scoring problem as in Eq. (12.24) has solution $\hat{\mathbf{\Theta}} = \mathbf{U}$.*

Proof. From the first constraint of GSVD, obviously $\mathbf{U}^T \mathbf{Y}^T \mathbf{Y} \mathbf{U} = \mathbf{I}$. from Eq. (12.32a), \mathbf{U} diagonalizes $\mathbf{Y}^T \mathbf{S} \mathbf{Y}$ by $\mathbf{U}^T \mathbf{Y}^T \mathbf{S} \mathbf{Y} \mathbf{U} = \mathbf{D}^2$. \square

Consequently, we have $\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \mathbf{U} = \mathbf{V} \mathbf{D}$. We can see that v_l (columns of $\hat{\mathbf{A}}$) and β_l (columns of $\hat{\mathbf{B}}$) differ by only a diagonal matrix $(\mathbf{I} - \mathbf{D}^2)^{-1} \mathbf{D}$.

Ex. 12.9

The reduced features are simply

$$\mathbf{X}^* = \mathbf{X} \hat{\mathbf{B}} = \mathbf{S} \mathbf{Y} \quad (12.36)$$

therefore the optimal scoring can be computed by

$$\max_{\boldsymbol{\Theta}} \boldsymbol{\Theta}^T \mathbf{Y}^T [\mathbf{S} \mathbf{Y} (\mathbf{Y}^T \mathbf{S} \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{S}^T] \mathbf{Y} \boldsymbol{\Theta} \quad (12.37)$$

$$\text{s.t. } \boldsymbol{\Theta}^T \mathbf{Y}^T \mathbf{Y} \boldsymbol{\Theta} = N \mathbf{I} \quad (12.38)$$

with trivial manipulations one can see that the objective function is exactly the same as optimal scoring on original features.

Ex. 12.10

Chapter 13

Prototype Methods and Nearest-Neighbors

Chapter 14

Unsupervised Learning

Chapter 15

Random Forests

Chapter 16

Ensemble Learning

Chapter 17

Undirected Graphical Models

Chapter 18

High-Dimensional Problems