

Assignment 6

Alberto Mejia

RIN: 661514960

CSCI 4100 - Machine Learning from Data

October 14, 2019

1. (200) LFD Exercise 3.4

(a)

$$y = \mathbf{w}^{*T} \mathbf{x} + \epsilon$$

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{w}^* + \mathbf{H} \epsilon$$

$$\mathbf{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$$

$$\boldsymbol{\epsilon} = [\epsilon_1, \epsilon_2, \dots, \epsilon_N]^T$$

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{w}_{\text{lin}} \rightarrow \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\begin{aligned} \hat{\mathbf{y}} &= \mathbf{H} \mathbf{y} \rightarrow \mathbf{H}(\mathbf{w}^{*T} \mathbf{x} + \epsilon) \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{w}^{*T} \mathbf{x} + \epsilon) \\ &= \mathbf{X} \mathbf{w}^* + \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \\ &= \mathbf{X} \mathbf{w}^* + \mathbf{H} \epsilon \end{aligned}$$

$$(b) \hat{\mathbf{y}} - \mathbf{y} = (\mathbf{X} \mathbf{w}^* + \mathbf{H} \epsilon) - (\mathbf{X} \mathbf{w}^* + \epsilon) = (\mathbf{H} - \mathbf{I}) \epsilon$$

$$(c) E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2 \quad (3.3)$$

$$\begin{aligned} E_{\text{in}}(\mathbf{w}_{\text{lin}}) &= \frac{1}{N} \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \\ &= \frac{1}{N} \|(\mathbf{H} - \mathbf{I}) \epsilon\|^2 \\ &= \frac{1}{N} [(\mathbf{H} - \mathbf{I}) \epsilon]^T [(\mathbf{H} - \mathbf{I}) \epsilon] \\ &= \frac{1}{N} \epsilon^T [(\mathbf{H} - \mathbf{I})]^T [(\mathbf{H} - \mathbf{I}) \epsilon] \\ &= \frac{1}{N} \epsilon^T (\mathbf{H} - \mathbf{I})^2 \epsilon \end{aligned}$$

$\mathbf{H} - \mathbf{I}$ is symmetric, thus $(\mathbf{H} - \mathbf{I})^T = \mathbf{H} - \mathbf{I}$

Therefore, $(\mathbf{I} - \mathbf{H})^T = \mathbf{I} - \mathbf{H}$

$$\begin{aligned} &= \frac{1}{N} \epsilon^T (\mathbf{I} - \mathbf{H})^2 \epsilon \\ &= \frac{1}{N} \epsilon^T (\mathbf{I} - \mathbf{H}) \epsilon \end{aligned}$$

Exercise 3.3 part (c)

$$(d) \text{ Prove } E_{\mathbf{D}}[E_{\text{in}}(\mathbf{w}_{\text{lin}})] = \sigma^2 \left(1 - \frac{d+1}{N}\right)$$

$$\text{trace}(\mathbf{H}) = d + 1$$

$$E_{\mathbf{D}}[\epsilon] = \sigma^2$$

$$E_{\mathbf{D}}[\epsilon^T \epsilon] = N \sigma^2$$

$$E_{\mathbf{D}}\left[\frac{1}{N} \epsilon^T \epsilon\right] = \frac{1}{N} N \sigma^2 \rightarrow \frac{N \sigma^2}{N} \rightarrow \sigma^2$$

$$E_{\mathbf{D}}[\epsilon^T \mathbf{H} \epsilon] = \text{trace}(\mathbf{H}) \sigma^2$$

Exercise 3.3 part (d)

$E_{\mathbf{D}}[\epsilon^T \mathbf{H} \epsilon]$ is a diagonal matrix

$$\begin{aligned} E_{\mathbf{D}}[E_{\text{in}}(\mathbf{w}_{\text{lin}})] &= E_{\mathbf{D}}\left[\frac{1}{N} \epsilon^T (\mathbf{I} - \mathbf{H}) \epsilon\right] \\ &= E_{\mathbf{D}}\left[\frac{1}{N} \epsilon^T (\mathbf{I} \epsilon - \mathbf{H} \epsilon)\right] \\ &= E_{\mathbf{D}}\left[\frac{1}{N} \epsilon^T \epsilon - \frac{1}{N} \epsilon^T \mathbf{H} \epsilon\right] \end{aligned}$$

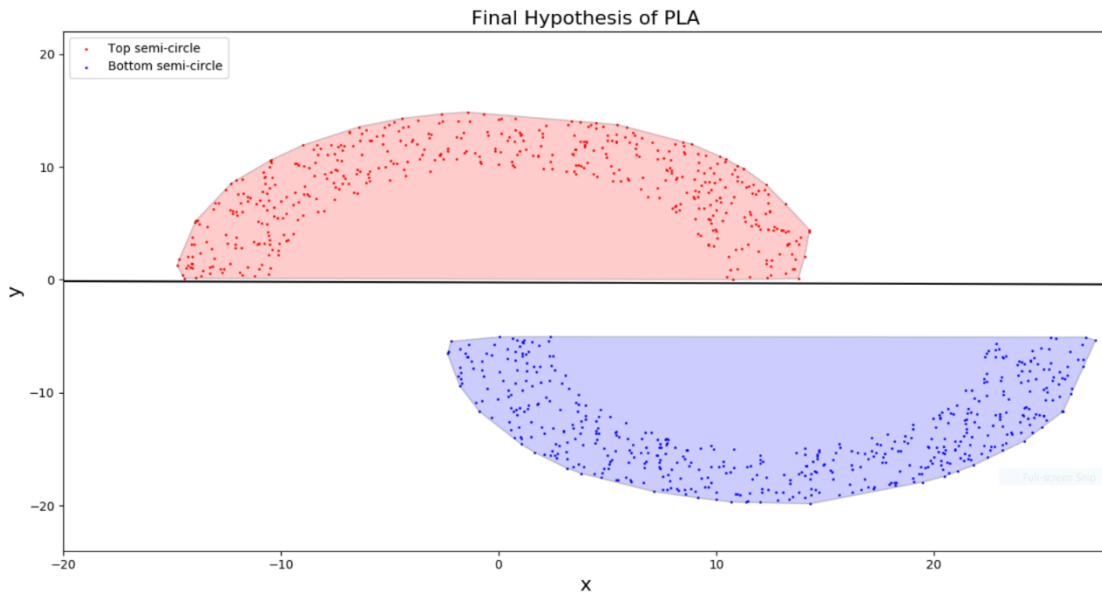
$$\begin{aligned}
&= \frac{1}{N} \mathbb{E}_D[\epsilon^T \epsilon] - \frac{1}{N} \mathbb{E}_D[\epsilon^T H \epsilon] \\
&= \frac{1}{N} (N \sigma^2 - \mathbb{E}_D[\epsilon^T H \epsilon]) \\
&= \frac{1}{N} (N \sigma^2 - \text{trace}(H) \sigma^2) \\
&= \frac{1}{N} (N \sigma^2 - (d+1) \sigma^2) \\
&= \sigma^2 - \frac{1}{N} \sigma^2 (d+1) \\
&= \sigma^2 - \sigma^2 \frac{(d+1)}{N} \\
&= \sigma^2 \left(1 - \frac{(d+1)}{N}\right)
\end{aligned}$$

(e) Prove $\mathbb{E}_{D, \epsilon'}[\mathbb{E}_{\text{test}}(w_{\text{lin}})] = \sigma^2 \left(1 + \frac{d+1}{N}\right)$

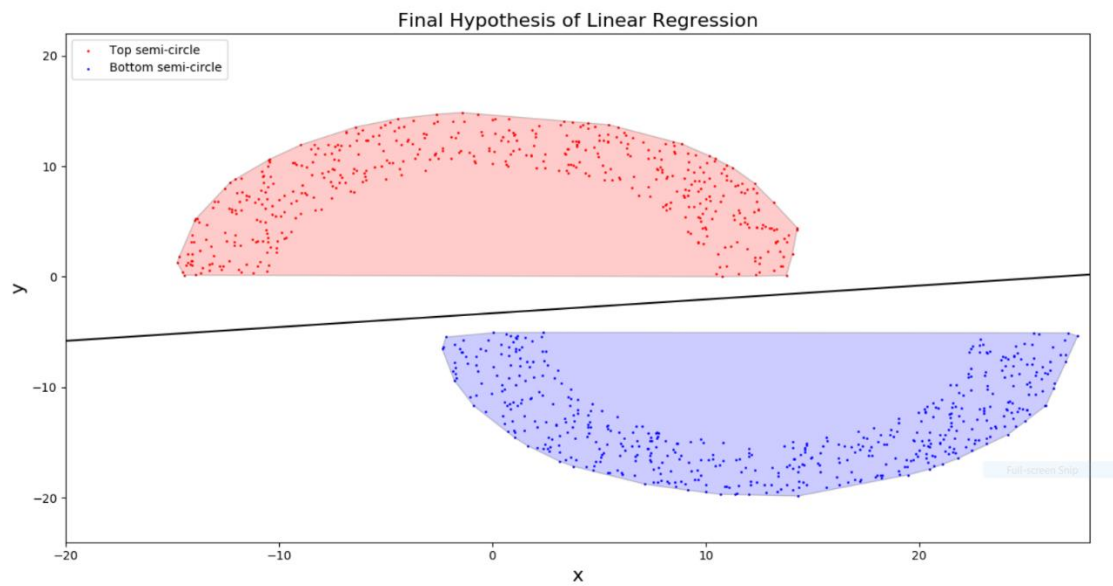
$$\begin{aligned}
\mathbb{E}_{D, \epsilon'}[\mathbb{E}_{\text{test}}(w_{\text{lin}})] &= \mathbb{E}_{D, \epsilon'}\left[\frac{1}{N} \|Xw - y\|^2\right] \\
&= \frac{1}{N} \mathbb{E}_{D, \epsilon'}[\|(H\epsilon - \epsilon')\|^2] \\
&= \frac{1}{N} \mathbb{E}_{D, \epsilon'}[\|(H\epsilon - \epsilon')^T (H\epsilon - \epsilon')\|^2] \\
&= \frac{1}{N} \mathbb{E}_{D, \epsilon'}[\|(H^T \epsilon^T - \epsilon'^T)(H\epsilon - \epsilon')\|^2] \\
&= \frac{1}{N} \mathbb{E}_{D, \epsilon'}[\|H^T \epsilon^T H \epsilon - H^T \epsilon^T \epsilon' - \epsilon'^T H \epsilon + \epsilon'^T \epsilon'\|^2] \\
&= \frac{1}{N} \mathbb{E}_{D, \epsilon'}[\|H^T \epsilon^T H \epsilon - 2\epsilon'^T H^T \epsilon' + \epsilon'^T \epsilon'\|^2] \\
&= \frac{1}{N} \mathbb{E}_{D, \epsilon'}[\|H^T \epsilon^T H \epsilon - 2\epsilon'^T H^T \epsilon' + \epsilon'^T \epsilon'\|^2] \\
&= \mathbb{E}_{D, \epsilon'}\left[\frac{1}{N} \mathbb{E}_{\epsilon'}[\|H^T \epsilon^T H \epsilon - 2\epsilon'^T H^T \epsilon' + \epsilon'^T \epsilon'\|^2]\right] \\
&= \mathbb{E}_{D, \epsilon'}\left[\frac{1}{N} \mathbb{E}_{\epsilon'}[\|H^T \epsilon^T H \epsilon - 2\epsilon'^T H^T \epsilon' + \epsilon'^T \epsilon'\|^2]\right] \\
&= \frac{1}{N} (\sigma^2 (d+1) + N \sigma^2) \\
&= \sigma^2 \left(1 + \frac{d+1}{N}\right)
\end{aligned}$$

2. (200) LFD Problem 3.1

(a)

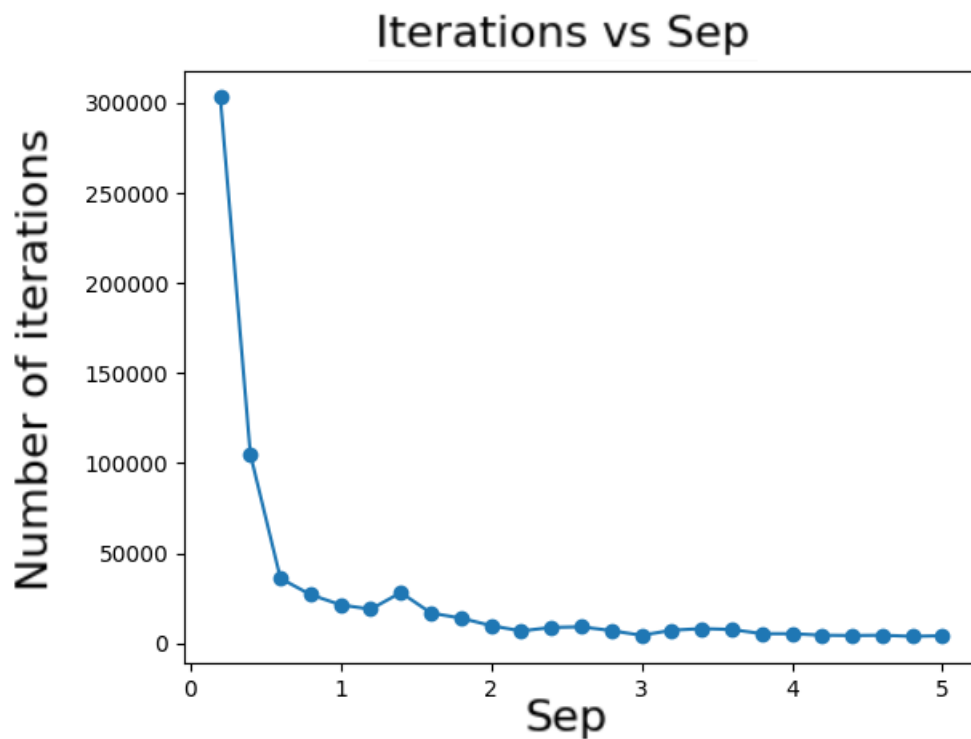


(b)



As we can observe, using linear regression we get a smaller in sample error as it produces a more precise result.

3. (200) LFD Problem 3.2



We can observe a couple of things from the graph. One of them is that as the Sep increases, the number of iterations decrease, AKA the number of steps for convergence. Thus, the points also shrink in distance as Sep increases.

4. (200) **LFD Problem 3.8**

For Linear Regression: $E_{\text{out}}(h) = E[(h(x) - y)^2]$

$$y = h^*(x) + \epsilon(x)$$

$$\epsilon(x) = y - h^*(x)$$

$$(a+b)^2 = a^2 + 2ab + b^2$$

(1) Algebraic formula

Show that $E_{\text{out}}(h)$ is minimized by: $h^*(x) = E[y|x]$ and show $E[\epsilon(x)] = 0$

$$\begin{aligned} E_{\text{out}}(h) &= E[(h(x) - y)^2] \\ &= E[(h(x) - (h^*(x) + \epsilon(x)))^2] \\ &= E[(h(x) - h^*(x) + y - h^*(x))^2] \\ &= E[((h(x) - h^*(x)) - y + h^*(x))^2] \\ &= E[((h(x) - h^*(x)) + (h^*(x) - y))^2] \\ &= E[(h(x) - h^*(x))^2 + (h^*(x) - y)^2] \\ &= E[(h(x) - h^*(x))^2 + 2(h(x) - h^*(x))(h^*(x) - y) + (h^*(x) - y)^2] \quad \text{Apply (1)} \\ &= E[(h(x) - h^*(x))^2] + 2E[(h(x) - h^*(x))(h^*(x) - y)] + E[(h^*(x) - y)^2] \end{aligned}$$

Assume minimization of $E_{\text{out}}(h)$ by $h^*(x) = E[y|x]$

$$E[(h(x) - h^*(x))(h^*(x) - y)]$$

$$\nabla E(h^*) = E[(h^*(x) - y)] = 0$$

$$E[y|x] = E[h^*(x) + \epsilon(x)|x]$$

$$= E[h^*(x)|x] + E[\epsilon(x)|x]$$

$$E[\epsilon(x)|x] = 0$$

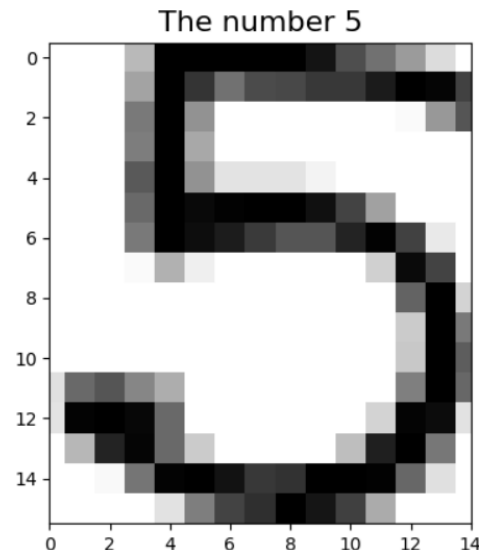
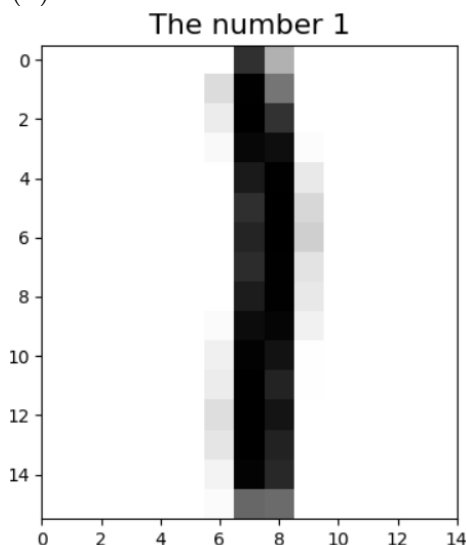
$$E[(h(x) - h^*(x))(h^*(x) - y)] = 0$$

$$\therefore E_{\text{out}}(h) \text{ is minimized when } h^*(x) = E[y|x]$$

5. (200) **LFD Problem 3.6 (6xxx Level Only)**

6. (200) Handwritten Digits Data - Obtaining Features

(a)



(b) We choose intensity and symmetry as our two features to measure the properties of the images that will help us distinguish between numbers like 1 and 5. These features are chosen not only because Professor Malik told us too, but also because the data is preprocessed already using Computer Vision to fit all numbers into the same aspect ratio(16×16) and then extra noise is filtered out (background is strictly one color and the numbers are strictly grayscale).

Let D be a matrix for a digit.

For intensity, we just do the following for each digit: Take the sum of all the grayscale values to determine the intensity of black it has per pixel. Thus, we get the following formula

$$Intensity = \sum_{i=0}^{15} (D[0][i] + D[1][j] + \dots + D[n][i])$$

For Symmetry, the logic is as follows: Calculate the horizontal symmetry of the image, which is defined as the difference between the image and its opposite-axis inverted version. Mathematically, we'll define it as

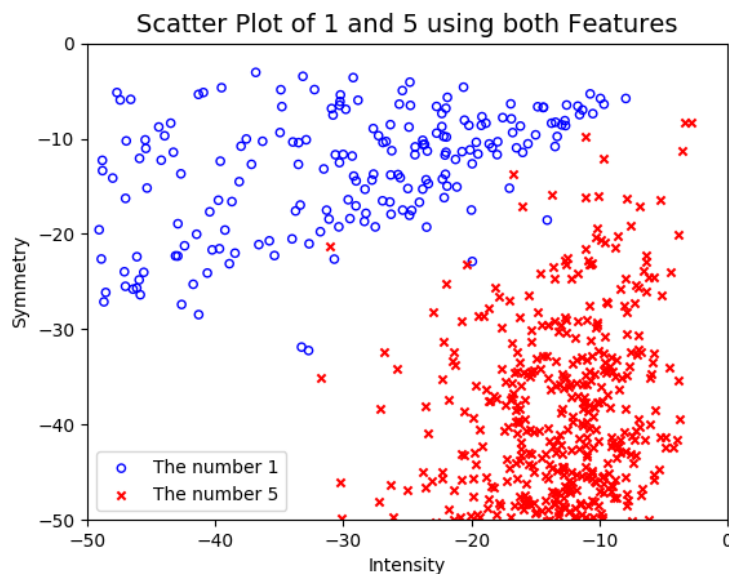
$$(D[i][j] - D[k][l]) = (r) \quad \text{if } (D[i][j] - D[k][l] > 0)$$

$$(D[i][j] - D[k][l]) = -1*(r) \quad \text{if } (D[i][j] - D[k][l] < 0)$$

We can see that this is the piecewise function for absolute value.

$$Symmetry = \sum_{i=0}^{15} \sum_0^7 |D[i][j] - D[i][15-j]|$$

(c)



As we can see, an increase in symmetry shows an increase of ones where as fives are more situated in the places with less symmetry. Fives and ones also have enough intensity difference to show, fives have more black pixels and thus more 'average'

intensity. Therefore they will be found more along increasing intensity. Being that we are just comparing pixels for symmetry, the fives will start to dip into the symmetry side of the graph but not as much as the ones. Less symmetrical ones will also dip into the intensity side since some ones will have “five” like attributes (like those fancy ones with an angled top and wide bottom)