

## 1 Bayes Optimal Classification (20 Points) (Yan)

1. We can write

$$\begin{aligned}
 \operatorname{argmin}_f \mathbb{E} \ell_{\alpha, \beta}(f(x), y) &= \operatorname{argmin}_f \mathbb{E}_{X, Y} [\alpha \mathbf{1}\{f(X) = 1, Y = 0\} + \beta \mathbf{1}\{f(X) = 0, Y = 1\}] \\
 &= \operatorname{argmin}_f \mathbb{E}_X [\mathbb{E}_{Y|X} [\alpha \mathbf{1}\{f(X) = 1, Y = 0\} + \beta \mathbf{1}\{f(X) = 0, Y = 1\}]] \\
 &= \operatorname{argmin}_f \mathbb{E}_X \left[ \int_y \alpha \mathbf{1}\{f(X) = 1, Y = 0\} + \beta \mathbf{1}\{f(X) = 0, Y = 1\} dP(y | x) \right] \\
 &= \operatorname{argmin}_f \int_x [\alpha \mathbf{1}\{f(x) = 1\} P(y = 0 | x) + \beta \mathbf{1}\{f(x) = 0\} P(y = 1 | x)] dP(x)
 \end{aligned}$$

We may minimize the integrand at each  $x$  by taking

$$f(x) = \begin{cases} 1 & \beta P(y = 1 | x) \geq \alpha P(y = 0 | x) \\ 0 & \alpha P(y = 0 | x) > \beta P(y = 1 | x) \end{cases}$$

2. Notice that

$$\begin{aligned}
 \mathbb{E} \ell_{\alpha, \beta}(f(x), y) &= \alpha P(f(x) = 1, y = 0) + \beta P(f(x) = 0, y = 1) \\
 &= \alpha P(f(x) = 1 | y = 0) P(y = 0) + \beta P(f(x) = 0 | y = 1) P(y = 1)
 \end{aligned}$$

which is same as the minimizer of the given risk  $R$  if  $\alpha = \frac{1}{P(y=0)}$  and  $\beta = \frac{1}{P(y=1)}$ .

3. Notice that since  $Y \sim \text{Ber}(\frac{1}{2})$ , we have  $P(Y = 1) = P(Y = 0) = \frac{1}{2}$ .

$$\begin{aligned}
 f^*(x) &= \operatorname{argmax}_y P(Y = y | X = x) = \operatorname{argmax}_y P(X = x | Y = y) P(Y = y) \\
 &= \operatorname{argmax}_y P(X = x | Y = y)
 \end{aligned}$$

Therefore,  $f^*(1) = 1$  since  $p = P(X = 1 | Y = 1) > P(X = 1 | Y = 0) = q$ , and  $f^*(0) = 0$  since  $1 - p = P(X = 0 | Y = 1) < P(X = 0 | Y = 0) = 1 - q$ . Hence,  $f^*(X) = X$ . The risk is  $R^* = P(f^*(X) \neq Y) = P(X \neq Y)$ .

$$R^* = P(Y = 1)P(X = 0 | Y = 1) + P(Y = 0)P(X = 1 | Y = 0) = \frac{1}{2} \cdot (1 - p) + \frac{1}{2} \cdot q$$

4. Figure 1 is a sample plot for this problem.

## 2 Regularized Linear Regression Using Lasso (20 Points) (Yan)

1. We can expand the equation as

$$\begin{aligned}
 &\frac{1}{2}(y^\top y - 2y^\top X\mathbf{w} + \mathbf{w}^\top X^\top X\mathbf{w}) + \lambda \|\mathbf{w}\|_1 \\
 &= \frac{1}{2}y^\top y + \sum_{i=1}^d -y^\top X_i \mathbf{w}_i + \frac{1}{2}\mathbf{w}_i^2 + \lambda |\mathbf{w}_i|
 \end{aligned}$$

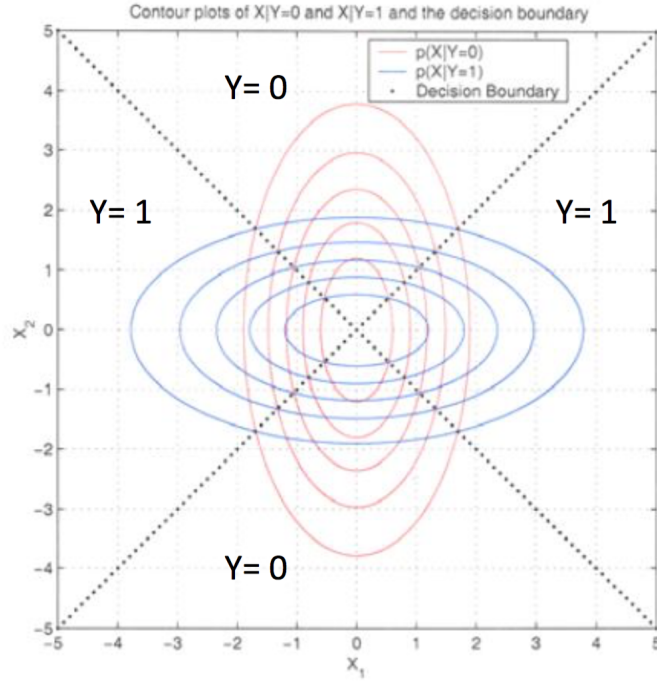


Figure 1: Sample plot for problem 1-4

2. If  $\mathbf{w}_i^* > 0$ , then we want to maximize

$$-y^\top X_{.i} \mathbf{w}_i + \frac{1}{2} \mathbf{w}_i^2 + \lambda \mathbf{w}_i$$

Take the derivative and equate to 0, we have:

$$\mathbf{w}_i^* = y^\top X_{.i} - \lambda$$

3. If  $\mathbf{w}_i^* < 0$ , then we want to maximize

$$-y^\top X_{.i} \mathbf{w}_i + \frac{1}{2} \mathbf{w}_i^2 - \lambda \mathbf{w}_i$$

Take the derivative and equate to 0, we have:

$$\mathbf{w}_i^* = y^\top X_{.i} + \lambda$$

4. From the previous questions, we know  $\mathbf{w}_i^* = 0$  if none of the above conditions hold, that is

$$y^\top X_{.i} - \lambda \leq 0 \quad y^\top X_{.i} + \lambda \geq 0$$

Combining them, we get

$$-\lambda \leq y^\top X_{.i} \leq \lambda$$

5. If the lasso is replaced by  $\frac{1}{2} \lambda \|\mathbf{w}\|_2^2$ , the optimization problem regarding  $\mathbf{w}_i$  is given by

$$-y^\top X_{.i} \mathbf{w}_i + \frac{1}{2} \mathbf{w}_i^2 + \frac{1}{2} \lambda \mathbf{w}_i^2$$

Take the derivative and equate to 0, we have:

$$\mathbf{w}_i^* = \frac{y^\top X_{.i}}{1 + \lambda}$$

It is equal to 0 if  $y^\top X_{.i} = 0$  or  $\lambda$  goes to infinity. In contrast,  $\mathbf{w}_i^* = 0$  when  $|y^\top X_{.i}| < \lambda$  in Lasso regression. This is why the L1 norm regularization encourages sparsity.

### 3 Multinomial Logistic Regression (20 Points) (Yan)

The solution for this question absorbs the intercept term  $w_{c0}$  into the vector of  $\mathbf{w}_c$ .

1. When  $C = 2$ ,

$$\begin{aligned} p(y = 1 \mid \mathbf{x}, W) &= \frac{\exp(w_{10} + \mathbf{w}_1^\top \mathbf{x})}{\exp(w_{10} + \mathbf{w}_1^\top \mathbf{x}) + \exp(w_{20} + \mathbf{w}_2^\top \mathbf{x})} = \frac{1}{1 + \exp(w_{20} - w_{10} + (\mathbf{w}_2 - \mathbf{w}_1)^\top \mathbf{x})} \\ p(y = 2 \mid \mathbf{x}, W) &= \frac{\exp(w_{20} + \mathbf{w}_2^\top \mathbf{x})}{\exp(w_{10} + \mathbf{w}_1^\top \mathbf{x}) + \exp(w_{20} + \mathbf{w}_2^\top \mathbf{x})} = \frac{\exp(w_{20} - w_{10} + (\mathbf{w}_2 - \mathbf{w}_1)^\top \mathbf{x})}{1 + \exp(w_{20} - w_{10} + (\mathbf{w}_2 - \mathbf{w}_1)^\top \mathbf{x})} \end{aligned}$$

This is equivalent with logistic regression that has weights  $(w_{20} - w_{10}, \mathbf{w}_2 - \mathbf{w}_1)$ .

2. Let  $\mu_{ic} = P(y_i = c \mid \mathbf{x}_i, W)$ ,  $y_{ic} = \mathbf{1}\{y_i = c\}$ .

(a)

$$\ell(W) = \log \prod_{i=1}^n \prod_{c=1}^C \mu_{ic}^{y_{ic}} = \sum_{i=1}^n \sum_{c=1}^C y_{ic} \log \mu_{ic} = \sum_{i=1}^n \left( \sum_{c=1}^C y_{ic} \mathbf{w}_c^\top \mathbf{x}_i - \log \sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i) \right)$$

(b)

$$\begin{aligned} g_c(W) &= \frac{\partial}{\partial \mathbf{w}_c} \sum_{i=1}^n \left( \sum_{c=1}^C y_{ic} \mathbf{w}_c^\top \mathbf{x}_i - \log \sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i) \right) \\ &= \sum_{i=1}^n \left( \frac{\partial}{\partial \mathbf{w}_c} \sum_{c=1}^C y_{ic} \mathbf{w}_c^\top \mathbf{x}_i - \frac{\partial}{\partial \mathbf{w}_c} \log \sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i) \right) \\ &= \sum_{i=1}^n \left( y_{ic} \mathbf{x}_i - \frac{\frac{\partial}{\partial \mathbf{w}_c} \sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)} \right) \\ &= \sum_{i=1}^n \left( y_{ic} \mathbf{x}_i - \frac{\exp(\mathbf{w}_c^\top \mathbf{x}_i) \mathbf{x}_i}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)} \right) \\ &= \sum_{i=1}^n (y_{ic} - \mu_{ic}) \mathbf{x}_i \end{aligned}$$

(c)  $\delta_{cc'}$  denotes the Dirac delta function and is equal to one if  $c = c'$  and zero otherwise.

$$\begin{aligned}
H_{c,c'}(W) &= \frac{\partial}{\partial \mathbf{w}_c} g_{c'}(W) \\
&= \frac{\partial}{\partial \mathbf{w}_c} \sum_{i=1}^n \left( y_{ic'} \mathbf{x}_i - \frac{\exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)}{\sum_{c''=1}^C \exp(\mathbf{w}_{c''}^\top \mathbf{x}_i)} \right) \\
&= - \sum_{i=1}^n \frac{\partial}{\partial \mathbf{w}_c} \frac{\exp(\mathbf{w}_{c'}^\top \mathbf{x}_i)}{\sum_{c''=1}^C \exp(\mathbf{w}_{c''}^\top \mathbf{x}_i)} \\
&= - \sum_{i=1}^n \frac{\delta_{cc'} \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \sum_{c''=1}^C \exp(\mathbf{w}_{c''}^\top \mathbf{x}_i) - \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \exp(\mathbf{w}_c^\top \mathbf{x}_i)}{\left( \sum_{c''=1}^C \exp(\mathbf{w}_{c''}^\top \mathbf{x}_i) \right)^2} \\
&= - \sum_{i=1}^n (\delta_{cc'} \mu_{ic} - \mu_{ic'} \mu_{ic}) \mathbf{x}_i \mathbf{x}_i^\top \\
&= \sum_{i=1}^n \mu_{ic} (\mu_{ic'} - \delta_{cc'}) \mathbf{x}_i \mathbf{x}_i^\top
\end{aligned}$$

## 4 Perceptron Mistake Bounds (20 Points) (Xun)

1. Expand  $\mathbf{w}^{(t)}$ :

$$\langle \mathbf{w}^{(t)}, \mathbf{w} \rangle = \langle \mathbf{w}^{(t-1)} + y^{(t)} \mathbf{x}^{(t)}, \mathbf{w} \rangle \quad (1)$$

$$= \langle \mathbf{w}^{(t-1)}, \mathbf{w} \rangle + \langle y^{(t)} \mathbf{x}^{(t)}, \mathbf{w} \rangle \quad (2)$$

$$\geq \langle \mathbf{w}^{(t-1)}, \mathbf{w} \rangle + \gamma, \quad (3)$$

where the inequality holds due to the separability assumption. Expand recursively until  $\mathbf{w}^{(0)}$ , we get the desired inequality.

2. Similarly,

$$\|\mathbf{w}^{(t)}\|_2^2 = \|\mathbf{w}^{(t-1)}\|_2^2 + 2 \cdot \langle y^{(t)} \mathbf{x}^{(t)}, \mathbf{w}^{(t-1)} \rangle + \|\mathbf{x}^{(t)}\|_2^2 \quad (4)$$

$$\leq \|\mathbf{w}^{(t-1)}\|_2^2 + M^2, \quad (5)$$

where we use the fact that  $(\mathbf{x}^{(t)}, y^{(t)})$  is a misclassified example, therefore  $\langle y^{(t)} \mathbf{x}^{(t)}, \mathbf{w}^{(t-1)} \rangle \leq 0$ . Result follows by recursively expand the inequality until  $\mathbf{w}^{(0)}$ .

3. Square both sides of the first inequality and apply Cauchy-Schwartz:

$$t^2 \gamma^2 \leq \langle \mathbf{w}^{(t)}, \mathbf{w} \rangle^2 \leq \|\mathbf{w}^{(t)}\|_2^2 \cdot \|\mathbf{w}\|_2^2 \leq t M^2. \quad (6)$$

4. False. If data is linearly separable by a margin, then there are infinitely many classifiers that achieve zero error. Depending on the ordering of the data, Perceptron may stop at any of the hyperplanes.

If a stronger argument is favored, we can easily construct a 2d example that achieves zero error even at  $\mathbf{w}^{(0)}$ , which may not have margin  $\gamma$ .

## 5 Logistic Regression for Image Classification (20 Points) (Xun)

### 5.1 Exploring the data

1. Each image is sized  $28 \times 28$ , so expressed as a  $784 \times 1$  vector.

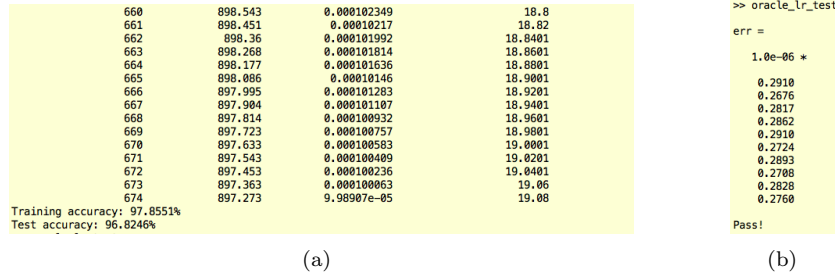


Figure 2: Results of binary logistic regression without regularization.

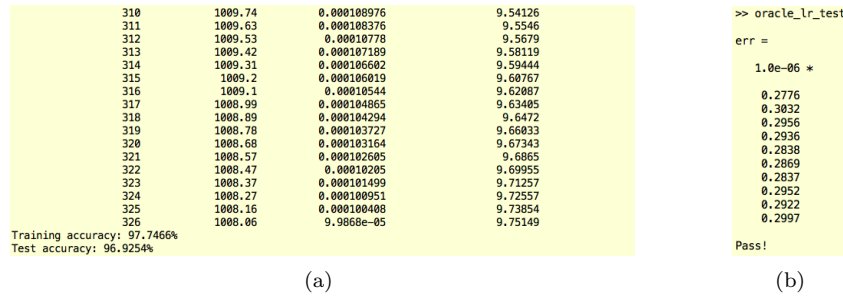


Figure 3: Results of binary logistic regression with regularization.

2. 0 to 9.
3. 0 to 1.
4. 3.5698 and 17.1790.
5. Nonzero fraction is 0.1912. Depending on the sparsity level, you could either argue it is sparse or dense.
6. Close to uniform.

## 5.2 Binary logistic regression

1. See the reference code.
2. See the reference code.
3. See the reference code. Figure 2(b) shows the result.
4. See the reference code.
5. See the reference code.
6. Figure 2(a) and (b) show the result produced by the reference code.
7. Figure 3(a) and (b) show the result produced by the reference code.

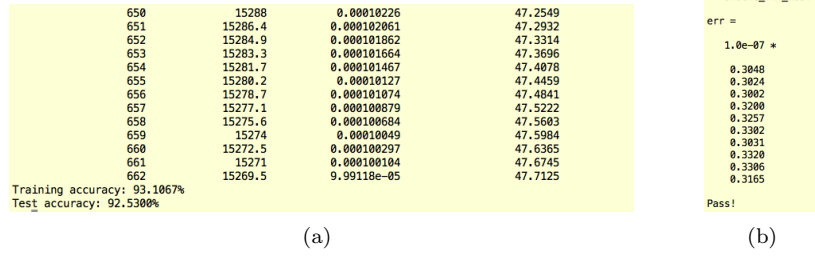


Figure 4: Results of multiclass logistic regression.

### 5.3 Multiclass logistic regression

1. See the reference code.
2. See the reference code.
3. See the reference code.
4. Figure 4(a) and (b) show the result produced by the reference code.