

EECS 553 Homework 2 Solution (FA24)

1. Full Rank Matrices (3 points)

Grading Rubrics

1. 3 points for fully correct answer with proper justification. Partial credit – Add one point for each:
 - 1 point for performing SVD on \mathbf{B} .
 - 1 point for showing \implies .
 - 1 point for showing \impliedby .
2. 0 if no effort.

Let $\mathbf{B} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be a thin SVD of \mathbf{B} where $\mathbf{U} \in \mathbb{R}^{p \times q}$, $\mathbf{V} \in \mathbb{R}^{q \times q}$, and $\mathbf{\Sigma} \in \mathbb{R}^{q \times q}$ and is diagonal. So $\mathbf{B}^T \mathbf{B} = \mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T$ is also an SVD.

$\text{rank}(\mathbf{B}) = q \iff$ all elements on the diagonal of $\mathbf{\Sigma}$ are nonzero \iff all elements on the diagonal of $\mathbf{\Sigma}^2$ are nonzero $\iff \mathbf{B}^T \mathbf{B} = \mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T$ is invertible.

Note: Alternative solution may receive full or partial credit.

Alternative Solution: \mathbf{B} has full rank $\iff \forall \mathbf{x} \neq \mathbf{0}, \mathbf{B}\mathbf{x} \neq \mathbf{0} \iff \forall \mathbf{x} \neq \mathbf{0}, \|\mathbf{B}\mathbf{x}\| > 0 \iff \forall \mathbf{x} \neq \mathbf{0}, \|\mathbf{B}\mathbf{x}\|^2 > 0 \iff \forall \mathbf{x} \neq \mathbf{0}, \mathbf{x}^T \mathbf{B}^T \mathbf{B} \mathbf{x} > 0 \iff \mathbf{B}^T \mathbf{B}$ is PD $\iff \mathbf{B}^T \mathbf{B}$ is invertible (since $\mathbf{B}^T \mathbf{B}$ is symmetric),

2. Weighted Least Squares (3 points)

Grading Rubrics

1. 3 points for fully correct answer with proper justification. Partial credit – Add one point for each:
 - 1 point for reformulating the problem as a quadratic programming problem.
 - 1 point for identifying the convexity of the problem.
2. 0 if no effort

Let

$$\mathbf{X} = \begin{bmatrix} 1 & \mathbf{x}_1^T \\ 1 & \mathbf{x}_2^T \\ \dots & \\ 1 & \mathbf{x}_n^T \end{bmatrix} \quad (1)$$

and

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix} \quad (2)$$

The objective can be written as $\mathbf{J}(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T \mathbf{C}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \boldsymbol{\theta}^T (\mathbf{X}^T \mathbf{C} \mathbf{X}) \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{C} \mathbf{y} + \mathbf{y}^T \mathbf{C} \mathbf{y}$.

$$\nabla_{\theta} J(\theta) = 2(\mathbf{X}^T \mathbf{C} \mathbf{X}) \theta - 2\mathbf{X}^T \mathbf{C} \mathbf{y}.$$

$$\nabla_{\theta}^2 J(\theta) = 2(\mathbf{X}^T \mathbf{C} \mathbf{X}).$$

$\nabla_{\theta}^2 J(\theta)$ is PSD, so the objective is convex. The solution can be determined letting $\nabla_{\theta} J(\theta) = 0$, leading to $\theta^* = (\mathbf{X}^T \mathbf{C} \mathbf{X})^{-1}(\mathbf{X}^T \mathbf{C} \mathbf{y})$, assuming $\mathbf{X}^T \mathbf{C} \mathbf{X}$ is invertible.

3. Grading Rubrics

1. 3 points for fully correct answer with proper justification. Partial credit – Add one point for each:

- 1 point for correct statements on the probabilistic assumptions.
- 1 point for applying the Bayes rule.
- 1 point for correct w and b .

2. 0 if no effort.

The LDA model assumes a normal distribution for the class-conditional densities $f_{\mathbf{X}|Y}(\mathbf{x}|Y = k)$, for $k = 0, 1$. On the other hand, logistic regression assumes

$$\eta(\mathbf{x}) = \Pr(Y = 1|\mathbf{X} = \mathbf{x}) = \frac{1}{1 + \exp\{-(\mathbf{w}^T \mathbf{x} + b)\}}$$

Now, we show that the distributional assumption of LDA implies the distributional assumption of logistic regression in the case of binary classification. Using Bayes rule, we have

$$\begin{aligned} \Pr(Y = 1|\mathbf{X} = \mathbf{x}) &= \frac{f_{\mathbf{X}|Y}(\mathbf{x}|Y = 1) \Pr(Y = 1)}{f_{\mathbf{X}}(\mathbf{x})} \\ &= \frac{f_{\mathbf{X}|Y}(\mathbf{x}|Y = 1) \Pr(Y = 1)}{f_{\mathbf{X}|Y}(\mathbf{x}|Y = 1) \Pr(Y = 1) + f_{\mathbf{X}|Y}(\mathbf{x}|Y = 0) \Pr(Y = 0)} \\ &= \frac{1}{1 + \frac{f_{\mathbf{X}|Y}(\mathbf{x}|Y=0) \Pr(Y=0)}{f_{\mathbf{X}|Y}(\mathbf{x}|Y=1) \Pr(Y=1)}} \end{aligned}$$

Let us denote by π_0, π_1 the *a priori* probabilities of the class label Y . Then, using the LDA assumption that $f_{\mathbf{X}|Y}(\mathbf{x}|Y = k)$ for $k \in \{0, 1\}$ is normally distributed according to $\mathcal{N}(\mu_k, \Sigma)$, we get

$$\begin{aligned} \Pr(Y = 1|X = \mathbf{x}) &= \frac{1}{1 + \frac{\pi_0 \exp\{-\frac{1}{2}(\mathbf{x}-\mu_0)^T \Sigma^{-1}(\mathbf{x}-\mu_0)\}}{\pi_1 \exp\{-\frac{1}{2}(\mathbf{x}-\mu_1)^T \Sigma^{-1}(\mathbf{x}-\mu_1)\}}} \\ &= \frac{1}{1 + \frac{\pi_0}{\pi_1} \exp\{ -[\mathbf{x}^T \Sigma^{-1}(\mu_1 - \mu_0) + \frac{\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1}{2}] \}} \\ &= \frac{1}{1 + \exp\{ -[\mathbf{x}^T \Sigma^{-1}(\mu_1 - \mu_0) + \frac{\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1}{2} - \log \frac{\pi_0}{\pi_1}] \}} \end{aligned}$$

Thus showing that $\Pr(Y = 1|X = \mathbf{x})$ is of the form $\frac{1}{1+\exp\{-(\mathbf{w}^T \mathbf{x} + b)\}}$, where $\mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_0)$ and $b = \frac{\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1}{2} - \log \frac{\pi_0}{\pi_1}$.

4. Logistic regression objective function (3 pts each)

(a) Grading Rubrics

1. 3 points for fully correct answers/reasoning. Some partial credits you can (but not required to, as long as the final answer is essentially correct) consider:
 - 1 point for showing correct $P(y|\tilde{\mathbf{x}}; \boldsymbol{\theta})$
 - 0.5 point each for showing the correct probability for $y = -1$ and $y = 1$
 - 1 point for showing the correct loss function
2. 0 if no effort

Recall logistic regression is assuming the following likelihood function:

$$P(y = 1|\tilde{\mathbf{x}}; \boldsymbol{\theta}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \tilde{\mathbf{x}}}}$$

$$P(y = -1|\tilde{\mathbf{x}}; \boldsymbol{\theta}) = \frac{e^{-\boldsymbol{\theta}^T \tilde{\mathbf{x}}}}{1 + e^{-\boldsymbol{\theta}^T \tilde{\mathbf{x}}}}$$

$$= \frac{1}{1 + e^{\boldsymbol{\theta}^T \tilde{\mathbf{x}}}}$$

Alternatively we can write:

$$P(y|\tilde{\mathbf{x}}; \boldsymbol{\theta}) = \frac{1}{1 + e^{-y\boldsymbol{\theta}^T \tilde{\mathbf{x}}}}$$

Thus the negative log-likelihood function:

$$-\ell(\boldsymbol{\theta}) = -\sum_{i=1}^n \log P(y_i|\tilde{\mathbf{x}}_i; \boldsymbol{\theta}) = \sum_{i=1}^n \log(1 + \exp(-y_i \boldsymbol{\theta}^T \tilde{\mathbf{x}}_i))$$

Hence with the new notation of $\phi(t) = \log(1 + \exp(-t))$, the logistic regression regularized negative log-likelihood may be written

$$J(\boldsymbol{\theta}) = \sum_{i=1}^n \phi(y_i \boldsymbol{\theta}^T \tilde{\mathbf{x}}_i) + \lambda \|\mathbf{w}\|^2.$$

(b) Grading Rubrics

1. 3 points for correct gradient (in either form given below)
2. 2 points for answers that only make minor mistakes including: sign errors, missing regularization terms, etc.
3. 1 point if mostly wrong
4. 0 if no effort

First by chain rule, we have:

$$\nabla_{\boldsymbol{\theta}} \phi(y_i \boldsymbol{\theta}^T \tilde{\mathbf{x}}_i) = \phi'(y_i \boldsymbol{\theta}^T \tilde{\mathbf{x}}_i) y_i \tilde{\mathbf{x}}_i$$

where $\phi'(t) = \frac{-\exp(-t)}{1 + \exp(-t)} = -\frac{1}{1 + \exp(t)}$

Then by linearity of gradient, we have:

$$\nabla J(\boldsymbol{\theta}) = 2\lambda[0, \mathbf{w}^T]^T + \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \phi(y_i \boldsymbol{\theta}^T \tilde{\mathbf{x}}_i)$$

$$= 2\lambda[0, \mathbf{w}^T]^T - \sum_{i=1}^n y_i \left(\frac{1}{1 + \exp(y_i \boldsymbol{\theta}^T \tilde{\mathbf{x}}_i)} \right) \tilde{\mathbf{x}}_i$$

Alternatively answer:

$$\nabla J(\theta) = 2\lambda[0, \mathbf{w}^T]^T - \sum_{i=1}^n y_i \left(\frac{\exp(-y_i \theta^T \tilde{\mathbf{x}}_i)}{1 + \exp(-y_i \theta^T \tilde{\mathbf{x}}_i)} \right) \tilde{\mathbf{x}}_i$$

(c) Grading Rubrics

1. 3 points for correct Hessian
2. 2 points for answers that only make minor mistakes including: sign errors, missing regularization terms, etc.
3. 1 point if mostly wrong
4. 0 point if no effort

The Hessian

$$\begin{aligned} \mathbf{H} &= \frac{\partial}{\partial \theta^T} \left(\frac{\partial J(\theta)}{\partial \theta} \right) \\ &= \frac{\partial}{\partial \theta^T} \left\{ 2\lambda[0, \mathbf{w}^T]^T - \sum_{i=1}^n y_i \left(\frac{1}{1 + \exp(y_i \theta^T \tilde{\mathbf{x}}_i)} \right) \tilde{\mathbf{x}}_i \right\} \\ &= 2\lambda \begin{bmatrix} 0 & 0 \\ 0 & \mathcal{I} \end{bmatrix} + \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T y_i^2 \left(\frac{\exp(y_i \theta^T \tilde{\mathbf{x}}_i)}{[1 + \exp(y_i \theta^T \tilde{\mathbf{x}}_i)]^2} \right) \\ &= 2\lambda \begin{bmatrix} 0 & 0 \\ 0 & \mathcal{I} \end{bmatrix} + \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \left(\frac{\exp(y_i \theta^T \tilde{\mathbf{x}}_i)}{[1 + \exp(y_i \theta^T \tilde{\mathbf{x}}_i)]^2} \right) \end{aligned}$$

where \mathcal{I} is a $d \times d$ identity matrix.

Note:

$$\begin{aligned} \frac{\exp(y_i \theta^T \tilde{\mathbf{x}}_i)}{[1 + \exp(y_i \theta^T \tilde{\mathbf{x}}_i)]^2} &= \frac{1}{[1 + \exp(y_i \theta^T \tilde{\mathbf{x}}_i)][1 + \exp(-y_i \theta^T \tilde{\mathbf{x}}_i)]} \\ &= \frac{1}{2 + \exp(y_i \theta^T \tilde{\mathbf{x}}_i) + \exp(-y_i \theta^T \tilde{\mathbf{x}}_i)} \end{aligned}$$

So any form of above are correct answers.

(d) Grading Rubrics

1. 3 points for fully correct answer that associated the (semi) positive-definiteness of the Hessian to λ . Some partial credits you can (but not required to, as long as the final answer is essentially correct) consider:
 - 1 point for computing bilinear form $\mathbf{z}^T \mathbf{H} \mathbf{z}$ correctly
 - 1 point for each for correctly concluding convex cases and strictly convex cases respectively
2. 0 point if no effort

Letting $a_i = \frac{\exp(y_i \theta^T \tilde{\mathbf{x}}_i)}{[1 + \exp(y_i \theta^T \tilde{\mathbf{x}}_i)]^2} > 0$ regardless of $\tilde{\mathbf{x}}_i$ and y_i , we have for any $\mathbf{z} \in \mathbb{R}^{d+1}$ such

that $\mathbf{z} \neq 0$:

$$\begin{aligned}\mathbf{z}^T \mathbf{H} \mathbf{z} &= \mathbf{z}^T \left(\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T a_i + 2\lambda \begin{bmatrix} 0 & 0 \\ 0 & \mathcal{I} \end{bmatrix} \right) \mathbf{z} \\ &= \sum_{i=1}^n a_i (\mathbf{z}^T \tilde{\mathbf{x}}_i) (\tilde{\mathbf{x}}_i^T \mathbf{z}) + 2\lambda (\mathbf{z}^T \mathbf{z} - z_1^2) \\ &= \sum_{i=1}^n a_i (\mathbf{z}^T \tilde{\mathbf{x}}_i)^2 + 2\lambda (\|\mathbf{z}\|^2 - z_1^2)\end{aligned}$$

Observe:

- 1) when $\lambda \geq 0$, we have $\mathbf{z}^T \mathbf{H} \mathbf{z} \geq 0, \forall \mathbf{z}$ (i.e Hessian is PSD everywhere), hence the problem is convex.
- 2) when $\lambda > 0$, $\forall \mathbf{z} \neq 0$, if for all i , $z_i = 0$ except for $z_1 \neq 0$, $\mathbf{z}^T \mathbf{H} \mathbf{z} > \sum_{i=1}^n a_i (\mathbf{z}^T \tilde{\mathbf{x}}_i)^2 = z_1^2 \sum_{i=1}^n a_i > 0$. Otherwise, $\mathbf{z}^T \mathbf{H} \mathbf{z} \geq 2\lambda (\|\mathbf{z}\|^2 - z_1^2) = 2\lambda \|\mathbf{z}\|^2 > 0$.

5. Logistic Regression for Fashion Classification (3 points each)

(a) Grading Rubrics

1. 3 points if fully correct:
 - 1 point for correct test error (within $\pm 1\%$)
 - 1 point for correct number of iterations (within ± 1)
 - 1 point for correct reported objective function after convergence (within ± 10)
2. 0 point if no effort

Test error = 3.6%

Number of iterations = 8

Value of objective function after convergence = 451.2632670172336

hw2 Note: the values might not exactly match. Answers within 1% of test error, ± 1 number of iterations, or ± 10 for the objective are acceptable.

(b) Grading Rubrics

1. 3 points if fully correct:
 - 1 point for explaining what is the notion of confidence used (anything related to distance from boundary or value of $x \cdot \theta + b$)
 - 2 points for showing the 20 misclassified examples (figures may not be exactly the same, give full credits for submissions that resemble the solution)
2. 0 point if no effort

See Figure 1 for the figure of the misclassified images. We define confidence as the distance to the learned hyperplane. The further a point is away from the hyperplane, the more confident the classifier is.

(c) Grading Rubrics

1. 3 points for submitting code (as long as students are not submitting only print functions, full credits will be given)
2. 0 point if not code submitted



Figure 1: P5 Figure