

Discussion #7 Exam Prep Solutions

Name:

1. Suppose in some universe, the true relationship between the measured luminosity of a single star Y can be written in terms of a single feature ϕ of that same star as

$$Y = \theta^* \phi + \epsilon$$

where $\phi \in \mathbb{R}$ is some non-random scalar feature, $\theta^* \in \mathbb{R}$ is a non-random scalar parameter, and ϵ is a random variable with $\mathbb{E}[\epsilon] = 0$ and $\text{var}(\epsilon) = \sigma^2$. For each star, you have a set of features $\Phi = [\phi_1 \ \phi_2 \ \dots \ \phi_n]^T$ and luminosity measurements $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_n]^T$ generated by this relationship. Your Φ may or may not include the feature ϕ described above. The ϵ_i for the various y_i have the same probability distribution and are independent of each other.

- (a) Suppose you have information about the exact ϕ value for each star, but try to fit a linear model for Y that includes an intercept term θ_0 .

$$Y = \theta_0 + \theta_1 \phi$$

Note the true relationship has no intercept term, so our model is not quite correct. Let $\hat{\theta}_0$ and $\hat{\theta}_1$ be the values that minimize the average L_2 loss. Let \mathbf{y} be the actual observed data and $\hat{\mathbf{y}} = \hat{\theta}_0 + \hat{\theta}_1 \phi$ be the fitted values.

- i. Which of the following could possibly be the value of $\hat{\theta}_0$ after fitting our model?
Select all that apply; at least one is correct.

☐ A. -1 ☐ B. 0 ☐ C. 1 ☐ D. 10

Solution: There are no restrictions on $\hat{\theta}_0$ given our assumptions.

- ii. Which of the following could possibly be the residual vector for our model?
Select all that apply; at least one is correct.

☐ A. $[-2 \ -4 \ 6]^T$ ☐ B. $[0.0001 \ 0.0003 \ -0.0005]^T$
☐ C. $[3 \ 12 \ -9]^T$ ☐ D. $[1 \ 1 \ 1]^T$

Solution: Since we are including an intercept/bias term, Φ has a column of 1s, which we denote with a boldface $\mathbf{1}$. Optimality requires orthogonality of the residual vector with the column space of Φ , which requires $\mathbf{1}^T \mathbf{e} = \sum_{i=1}^n 1 \times e_i = \sum_{i=1}^n e_i = 0$. (A) is the only choice satisfying this condition.

2. Throughout this section we refer to "least squares regression", which is the process of minimizing the average L2 loss using a linear regression model. Ordinary least squares is the version of least squares regression where **we do not use regularization**. Assume throughout that **our model includes a bias term**.

(a) What is always true about the residuals in least squares regression? Select all that apply.

- ☐ A. They are orthogonal to the column space of the features.
- ☐ B. They represent the errors of the predictions.
- ☐ C. Their sum is equal to the mean squared error.
- ☐ D. Their sum is equal to zero.
- ☐ E. None of the above.

Solution: (a), (b)

(c) is supposed to be a trick since the mean squared error is the *mean* of the sum of the *squares* of the residuals. So I guess this tests whether they understand what the acronym MSE represents or if they just regurgitate it mindlessly.

(e) is wrong since (c) is wrong obviously

(b) Which are true about the predictions made by OLS? Select all that apply.

- ☐ A. They are projections of the observations onto the column space of the features.
- ☐ B. They are linear in the chosen features.
- ☐ C. They are orthogonal to the residuals.
- ☐ D. They are orthogonal to the column space of the features.
- ☐ E. None of the above.

Solution: (a), (b), (c)

(a) is correct because they are linear projections onto the column space. This fact also makes (c) correct and (e) incorrect and is what makes (d) incorrect.

(b) is also correct because even in e.g. polynomial regression the resulting predictions are linear in the new/transformed features. But admittedly this is somewhat awkwardly worded.

(c) Which of the following would be true if you chose mean absolute error (L1) instead of mean squared error (L2) as your loss function? Select all that apply.

- ☐ A. The results of the regression would be more sensitive to outliers.

- ☐ B. You would not be able to use gradient descent to find the regression line.
- ☒ C. You would not be able to use the normal equation to calculate your parameters.
- ☐ D. The sum of the residuals would now be zero.
- ☐ E. None of the above.

Solution: (e)

(a) is false because using $L1$ loss increases robustness to outliers.

(b) is false because you can still use (sub)gradient descent given the convexity of $L1$ loss.

(c) is true, the normal equation only works if we're minimizing the $L2$ loss.

(d) is false because the sum of the residuals was zero in OLS, so if this happened it wouldn't be a change from OLS.

3. Let $\hat{\mathbf{y}} \in \mathbb{R}^n$ be the vector of fitted values in the ordinary least squares regression of $\mathbf{y} \in \mathbb{R}^n$ on the full column-rank feature matrix $\Phi \in \mathbb{R}^{n \times d}$ with n much larger than d . Denote the fitted coefficients as $\hat{\beta} \in \mathbb{R}^d$ and the vector of residuals as $\mathbf{e} \in \mathbb{R}^n$.

(a) What is $\Phi(\Phi^T\Phi)^{-1}\Phi^T\mathbf{y}$?

- ☐ A. $\mathbf{0}$ ☒ B. $\hat{\mathbf{y}}$ ☐ C. \mathbf{e} ☐ D. $\hat{\beta}$ ☐ E. 1 ☐ F. None of the above

Solution: We discussed this in discussion 6, where we called $\Phi(\Phi^T\Phi)^{-1}\Phi^T$ the hat matrix. It projects \mathbf{y} into the feature space.

(b) What is $\Phi(\Phi^T\Phi)^{-1}\Phi^T\hat{\mathbf{y}}$? Notice: This problem has a hat in $\hat{\mathbf{y}}$.

- ☐ A. $\mathbf{0}$ ☒ B. $\hat{\mathbf{y}}$ ☐ C. \mathbf{e} ☐ D. $\hat{\beta}$ ☐ E. 1 ☐ F. None of the above

Solution: Since $\hat{\mathbf{y}}$ is already in the feature space, projecting it into the feature space has no effect.

Suppose $\mathbf{e} \neq \mathbf{0}$. Define a new feature matrix Ψ by appending the residual vector \mathbf{e} to the feature matrix Φ . In other words,

$$\Psi = \begin{bmatrix} | & | & \vdots & | & | \\ \Phi_{:,1} & \Phi_{:,2} & \cdots & \Phi_{:,d} & \mathbf{e} \\ | & | & \vdots & | & | \end{bmatrix}$$

- (c) We now want to fit the model $\mathbf{y} = \Psi\boldsymbol{\gamma} = \gamma_1\Phi_{:,1} + \gamma_2\Phi_{:,2} + \cdots + \gamma_d\Phi_{:,d} + \gamma_{d+1}\mathbf{e}$ by choosing $\hat{\boldsymbol{\gamma}} = [\hat{\gamma}_1 \dots \hat{\gamma}_{d+1}]^T$ to minimize the L_2 loss. What is $\hat{\gamma}_{d+1}$?

- ☐ A. 0 ☒ B. 1 ☐ C. $\mathbf{e}^T \mathbf{y}$ ☐ D. $1 - \hat{\boldsymbol{\beta}}^T \hat{\boldsymbol{\beta}}$
☐ E. $(\Phi^T \Phi)^{-1} \Phi^T$ ☐ F. None of the above

Solution: We're effectively memorizing all of our regression values here. This is the equivalent (in a roundabout way) of using someone's weight, age, and height to predict their height. It'll work perfectly, but the model is useless.

4. We collect some data $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ and decide to model the relationship between \mathbf{X} and \mathbf{y} as

$$\mathbf{y} = \beta_1 \Phi_{:,1} + \beta_2 \Phi_{:,2}$$

where $\Phi_{i,:} = [1 \ x_i]$. We found the estimates $\hat{\beta}_1 = 2$ and $\hat{\beta}_2 = 5$ for the coefficients by minimizing the L_2 loss. Given that $\Phi^T \Phi = \begin{bmatrix} 4 & 2 \\ 2 & 5 \end{bmatrix}$, answer the following problems. If not enough information is given, write "Cannot be determined."

- (a) What was the sample size n ? Hint: Consider the form of the feature matrix.

Solution:

$$[\Phi^T \Phi]_{1,1} = \sum_{i=1}^n 1 \times 1 = n = 4$$

- (b) What must $\Phi^T \mathbf{y}$ be for this data set?

Solution: $\hat{\boldsymbol{\beta}}$ comes from the normal equations

$$\Phi^T \Phi \hat{\boldsymbol{\beta}} = \Phi^T \mathbf{y}$$

Therefore, we have

$$\Phi^T \mathbf{y} = \begin{bmatrix} 4 & 2 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 18 \\ 29 \end{bmatrix}$$

5. Consider the following loss function based on data x_1, \dots, x_n with mean \bar{x} :

$$\ell(\beta) = \log \beta + \frac{\bar{x}}{\beta} + \frac{1}{n} \sum_{i=1}^n e^{-x_i/\beta}$$

Given an estimate $\beta^{(t)}$, write out the update $\beta^{(t+1)}$ after one iteration of gradient descent with step size α .

Solution: The update is

$$\beta^{(t+1)} \leftarrow \beta^{(t)} - \alpha \ell'(\beta^{(t)}),$$

where

$$\begin{aligned} \ell'(\beta) &= \frac{1}{\beta} \left(1 - \frac{\bar{x}}{\beta} + \frac{1}{n\beta} \sum_{i=1}^n x_i e^{-x_i/\beta} \right) \\ &= \frac{1}{\beta} - \frac{\bar{x}}{\beta^2} + \frac{1}{n\beta^2} \sum_{i=1}^n x_i e^{-x_i/\beta} \end{aligned}$$

Alternate notation:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} - \alpha \left. \frac{\partial \ell}{\partial \beta} \right|_{\beta=\beta^{(t)}}$$

With everything substituted in:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} - \alpha \left(\frac{1}{\beta^{(t)}} - \frac{\bar{x}}{\beta^{(t)2}} + \frac{1}{n\beta^{(t)2}} \sum_{i=1}^n x_i e^{-x_i/\beta^{(t)}} \right)$$