

Discussion # 10

Name:

Cross Validation and Regularization

1. You build a model with two regularization hyperparameters λ and γ . You have 4 good candidate values for λ and 3 possible values for γ , and you are wondering which λ, γ pair will be the best choice. If you were to perform five-fold cross-validation, how many validation errors would you need to calculate?
2. In the typical setup of k-fold cross validation, we use a different parameter value on each fold, compute the mean squared error of each fold and choose the parameter whose fold has the lowest loss.
☐ A. True
☐ B. False

Questions 3, 4, 5, and 6 are all connected.

3. Elastic Net is a regression technique that combines L_1 and L_2 regularization. It is preferred in many situations as it possesses the benefits of both LASSO and Ridge Regression. Minimizing the L2 loss using Elastic Net is as follows, where $\lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = \lambda, \lambda > 0$.

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_i (y_i - \theta x)^2 + \lambda_1 \sum_{j=1}^p |\theta_j| + \lambda_2 \sum_{j=1}^p \theta_j^2$$

Suppose our goal was to get sparse parameters, i.e. we want as many parameters as possible to be zero. Which of the following choices for λ_1, λ_2 are most consistent with this goal, assuming $\lambda = 1$? **There is only one correct answer.**

- ☐ A. $\lambda_1 = 0, \lambda_2 = 1$
☐ B. $\lambda_1 = 0.5, \lambda_2 = 0.5$
☐ C. $\lambda_1 = 1, \lambda_2 = 0$

4. What happens to bias and variance as we increase the value of λ ? Assume $\lambda_2 = \lambda_1$. **There is only one correct answer in each part.** You will be asked to justify why in the next question.

(a) Bias:

- ☐ A. Bias goes up
- ☐ B. Bias stays the same
- ☐ C. Bias goes down

(b) Variance:

- ☐ A. Variance goes up
- ☐ B. Variance stays the same
- ☐ C. Variance goes down

5. Justify why by marking the true statements. **Select all that apply for each part.**

(a) Bias:

- ☐ A. Bias goes down because increasing λ reduces over fitting.
- ☐ B. Bias goes down because bias is minimized when $\lambda_2 = \lambda_1$.
- ☐ C. Bias goes up because increasing λ penalizes complex models, limiting the set of possible solutions.
- ☐ D. Bias goes up because the loss function becomes non-convex for sufficiently large λ .
- ☐ E. None of the above

(b) Variance:

- ☐ A. Variance goes down because increasing λ encourages the value of the loss to decrease.
- ☐ B. Variance goes down because increasing λ penalizes large model weights.
- ☐ C. Variance goes up because because increasing λ increases bias.
- ☐ D. Variance goes up because increasing λ increases the magnitude of terms in the loss function.
- ☐ E. None of the above

6. What happens to the model parameters $\hat{\theta}$ as $\lambda \rightarrow \infty$, i.e. what is $\lim_{\lambda \rightarrow \infty} \hat{\theta}$? **Select all that apply.**

- ☐ A. Converge to 0.
- ☐ B. Diverge to infinity.
- ☐ C. Converge to values that minimize the L2 loss.
- ☐ D. Converge to equal but non-zero values.
- ☐ E. Converge to a sparse vector.

Logistic Regression

1. You have a classification data set consisting of two (x, y) pairs $(1, 0)$ and $(-1, 1)$. You decide that you want your feature vector \mathbf{x} (the input to a model) for each pair to be a two-element column vector $\begin{bmatrix} 1 & x \end{bmatrix}^T$.

You run an algorithm to fit a model for the probability of $Y = 1$ given \mathbf{x} :

$$\mathbb{P}(Y = 1 \mid \mathbf{x}) = \sigma(\theta \cdot \mathbf{x})$$

where $\sigma(t) = \frac{1}{1 + \exp(-t)}$. Your algorithm returns $\hat{\theta} = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}$.

- (a) Calculate $\hat{\mathbb{P}}(Y = 1 \mid \mathbf{x} = \begin{bmatrix} 1 & 0 \end{bmatrix}^T)$

- (b) The empirical cross-entropy loss (a.k.a. log loss) is given by:

$$L(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(z_i) + (1 - y_i) \log(1 - z_i)]$$

where $z_i = \sigma(\theta \cdot \mathbf{X}_i)$. Let $\theta = [\theta_0 \quad \theta_1]$. Explicitly write out the empirical loss for the data set $(1, 0)$ and $(-1, 1)$ as a function of θ_0 and θ_1 . Note that in this problem, \mathbf{X}_i is the feature vector \mathbf{x} defined in the original problem statement.

- (c) Calculate the empirical loss for $\hat{\theta} = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}^T$ and the two observations $(1, 0)$ and $(-1, 1)$.