

Solutions

1 K-Nearest Neighbors

1. **Select all that apply:** Please select all that apply about k-NN in the following options:

- ☐ k-NN works great with a small amount of data, but it is too slow when the amount of data becomes large.
- ☐ k-NN is sensitive to outliers; therefore, in general, we decrease k to avoid overfitting.
- ☐ k-NN can only be applied to classification problems, and it cannot be used to solve regression problems.
- ☐ We can always achieve zero training error (perfect classification) on a consistent data set with k-NN, but it may not generalize well in testing.

Option 1: True: Curse of dimensionality

Option 2: False: we increase k to avoid overfitting

Option 3: False: K-NN regression

Option 4: True: by setting $k = 1$

2. (1 point) **Select one:** A k-Nearest Neighbor model with a large value of k is analogous to...

- ☐ A *short* Decision Tree with a *low* branching factor
- ☐ A *short* Decision Tree with a *high* branching factor
- ☐ A *long* Decision Tree with a *low* branching factor
- ☐ A *long* Decision Tree with a *high* branching factor

A short Decision Tree with a low branching factor

3. (1 point) **Select one.** Imagine you are using a k -Nearest Neighbor classifier on a data set with lots of noise. You want your classifier to be *less* sensitive to the noise. Which is more likely to help and with what side-effect?

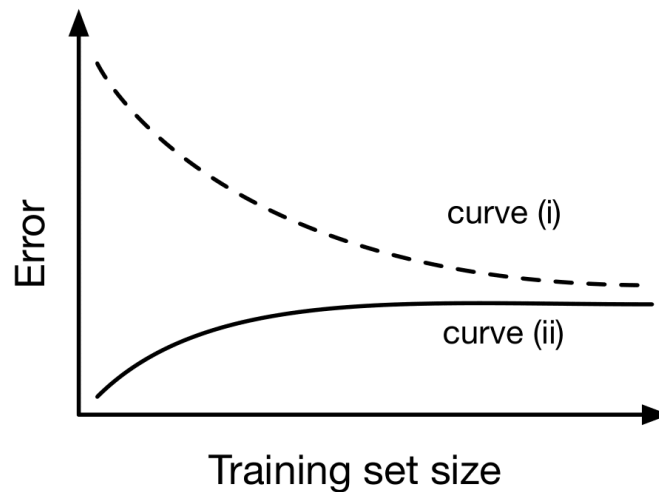
- ☐ Increase the value of $k \rightarrow$ Increase in prediction time
- ☐ Decrease the value of $k \rightarrow$ Increase in prediction time
- ☐ Increase the value of $k \rightarrow$ Decrease in prediction time
- ☐ Decrease the value of $k \rightarrow$ Decrease in prediction time

Increase the value of $k \rightarrow$ Increase in prediction time

2 Model Selection and Errors

1. **Training Sample Size:** In this problem, we will consider the effect of training sample size N on a linear regression problem with M features.

The following plot shows the general trend for how the training and testing error change as we increase the training sample size N . Your task in this question is to analyze this plot and identify which curve corresponds to the training and test error. Specifically:



1. Which curve represents the training error? **Please provide 1–2 sentences of justification.** Curve (ii) is the training set. Training error increases as the training set increases in size (more points to account for). However, the increase tapers out when the model generalizes well. Evidently, curve (i) is testing, since larger training sets better form generalized models, which reduces testing error.
2. In one word, what does the gap between the two curves represent? **Overfitting**

3 Linear Regression

1. (1 point) **Select one:** The closed form solution for linear regression is $\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. Suppose you have $N = 35$ training examples and $M = 5$ features (excluding the bias term). Once the bias term is now included, what are the dimensions of \mathbf{X} , \mathbf{y} , $\hat{\theta}$ in the closed form equation?

- ☐ \mathbf{X} is 35×6 , \mathbf{y} is 35×1 , $\hat{\theta}$ is 6×1
☐ \mathbf{X} is 35×6 , \mathbf{y} is 35×6 , $\hat{\theta}$ is 6×6
☐ \mathbf{X} is 35×5 , \mathbf{y} is 35×1 , $\hat{\theta}$ is 5×1
☐ \mathbf{X} is 35×5 , \mathbf{y} is 35×5 , $\hat{\theta}$ is 5×5

A.

2. Consider linear regression on N 1-dimensional points $x^{(i)} \in \mathbb{R}$ with labels $y^{(i)} \in \mathbb{R}$. We apply linear regression in both directions on this data, i.e., we first fit y with x and get $y = \beta_1 x$ as the fitted line, then we fit x with y and get $x = \beta_2 y$ as the fitted line. Discuss the relations between β_1 and β_2 :

True or False: The two fitted lines are always the same, i.e. we always have $\beta_2 = \frac{1}{\beta_1}$.

- ☐ True
☐ False

False. $\beta_1 = \frac{x^T y}{x^T x}$ and $\beta_2 = \frac{y^T x}{y^T y}$

3. Please circle **True** or **False** for the following questions, providing brief explanations to support your answer.

- (i) [3 pts] Consider a linear regression model with only one parameter, the bias, i.e., $y = b$. Then given N data points $(x^{(i)}, y^{(i)})$ (where $x^{(i)}$ is the feature and $y^{(i)}$ is the output), minimizing the sum of squared errors results in b being the median of the $y^{(i)}$ values.

Circle one: **True** **False**

Brief explanation:

False. $\sum_{i=1}^N (y^{(i)} - b)^2$ is the training cost, which when differentiated and set to zero gives $b = \frac{\sum_{i=1}^N y^{(i)}}{N}$, the mean of the $y^{(i)}$ values.

- (ii) [3 pts] Given data $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$, we obtain $\hat{\theta}$, the parameters that minimize the training error cost for the linear regression model $y = \theta^T \mathbf{x}$ we learn from \mathcal{D} .

Consider a new dataset \mathcal{D}_{new} generated by duplicating the points in \mathcal{D} and adding 10 points that lie along $y = \hat{\theta}^T \mathbf{x}$. Then the $\hat{\theta}_{\text{new}}$ that we learn for $y = \theta^T \mathbf{x}$ from

\mathcal{D}_{new} is equal to $\hat{\theta}$.

Circle one: **True** **False**

Brief explanation:

True. The new squared error can be written as $2\epsilon_1 + \epsilon_2$, where ϵ_1 is the old squared error. $\epsilon_2 = 0$ for the 10 points that lie along the line, the lowest possible value for ϵ_2 . And $2\epsilon_1$ is least when ϵ_1 is least, which is when the parameters don't change.

4. We have an input x and we want to estimate an output y using linear regression.

Consider the dataset S plotted in Fig. 1 along with its associated regression line. For each of the altered data sets S^{new} plotted in Fig. 2, indicate which regression line (relative to the original one) in Fig. 3 corresponds to the regression line for the new data set. Write your answers in the table below.

Dataset	(a)	(b)	(c)	(d)	(e)
Regression line					

Dataset	(a)	(b)	(c)	(d)	(e)
Regression line	(b)	(c)	(b)	(a)	(a)

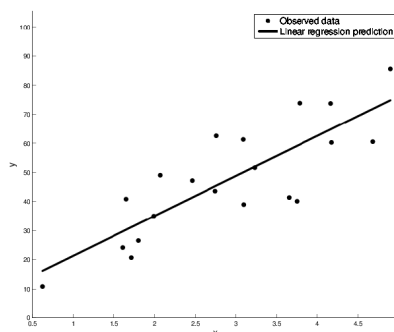
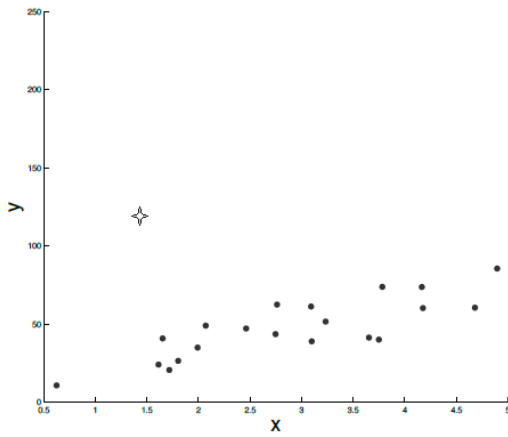
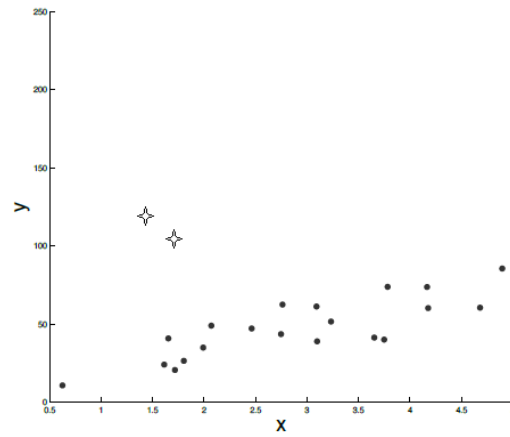


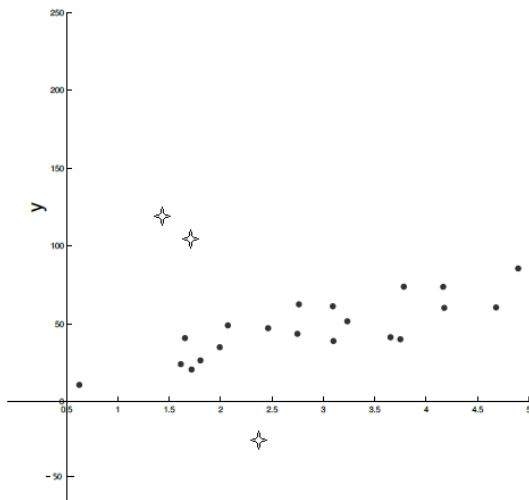
Figure 1: An observed data set and its associated regression line.



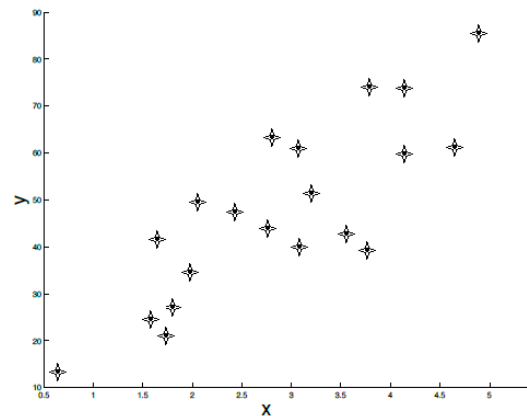
(a) Adding one outlier to the original data set.



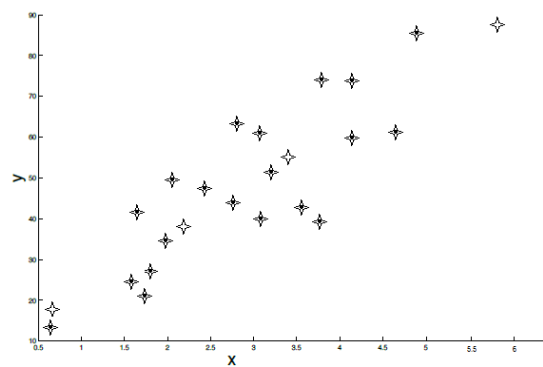
(b) Adding two outliers to the original data set.



(c) Adding three outliers to the original data set. Two on one side and one on the other side.

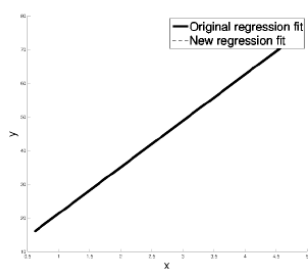


(d) Duplicating the original data set.

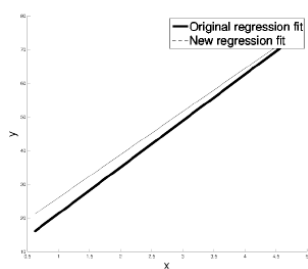


(e) Duplicating the original data set and adding four points that lie on the trajectory of the original regression line.

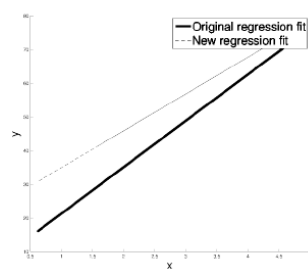
Figure 2: New data set S^{new} .



(a) Old and new regression lines.



(b) Old and new regression lines.



(c) Old and new regression lines.

Figure 3: New regression lines for altered data sets S^{new} .