

Q1.1) given $\{1, \dots, n\} (x_i, y_i), x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$

$$\|\omega_{opt}\| = 1, \gamma > 0, y_i \cdot (\omega_{opt}^T x_i) \geq \gamma, \forall i \in \{1, \dots, n\}$$

(3)

$$\|x_i\|_2 \leq R, R > 0$$

$$\text{if } y_i (\omega_k^T x_i) < 0 \Rightarrow \omega_{k+1}^T \omega_{opt} \geq \omega_k^T \omega_{opt} + \gamma \|\omega_{opt}\|$$

Proof:

a. considering perceptron algorithm misclassifies (x_i, y_i) in k -th iteration \Rightarrow we have

$$y_i (\omega_k^T x_i) < 0 \quad \& \text{ updated vector of } \omega_{k+1} = \omega_k + y_i x_i$$

(2)

(1)

b. we are adding a dot product of ω_{opt} on both sides of equation (2):

$$\omega_{k+1}^T \omega_{opt} = (\omega_k + y_i x_i)^T \omega_{opt}$$

c. expanding step b:

$$\omega_{k+1}^T \omega_{opt} = \omega_k^T \omega_{opt} + y_i x_i^T \omega_{opt}$$

(4)

rules used (side Note)
$(a+b) \cdot c = a \cdot c + b \cdot c$
$a \cdot b = a^T b$

$$y_i^T = y_i$$

d. given (1) we can say y_i and $(\omega_k^T x_i)$ have different signs.

using the given (3) assumption of $y_i (\omega_k^T \omega_{opt})$ has a opposite sign to $\gamma \Rightarrow$

$$y_i (\omega_k^T \omega_{opt}) \geq \gamma \|\omega_{opt}\|_2$$

(5)

e. using equations from steps c & d:

$$y_i(x_i^T \omega_{opt}) \geq \gamma \|\omega_{opt}\|_2$$

$$\omega_{k+1}^T \omega_{opt} = \omega_k^T \omega_{opt} + y_i x_i^T \omega_{opt} \rightarrow$$

$$(\omega_{k+1}^T \omega_{opt}) - (\omega_k^T \omega_{opt}) = y_i x_i^T \omega_{opt}$$

$$(\omega_{k+1}^T \omega_{opt}) - (\omega_k^T \omega_{opt}) \geq \gamma \|\omega_{opt}\|_2 \xrightarrow{\text{add}} (\omega_k^T \omega_{opt})$$

$$\therefore \omega_{k+1}^T \omega_{opt} \geq \omega_k^T \omega_{opt} + \gamma \|\omega_{opt}\|_2$$

to both sides

1.2)

$$(a+b)^T = a^T + b^T$$

1.2 According to the perceptron algorithm:

$$w_{k+1} = w_k + y_i x_i$$

general rule

$$\|w_{k+1}\|^2 = (w_{k+1})^T (w_{k+1})$$

$$\text{Calculate } \|w_{k+1}\|_2^2$$

$$\|w_{k+1}\|_2^2 = w_{k+1}^T w_{k+1} = (w_k + y_i x_i)^T (w_k + y_i x_i)$$

$$g_i^T = g_i$$

$$= (w_k^T + y_i x_i^T) (w_k + y_i x_i)$$

$$= w_k^T w_k + y_i w_k^T x_i + y_i x_i^T w_k + y_i^2 x_i^T x_i$$

$$\left\{ \begin{array}{l} w_k^T w_k = \|w_k\|_2^2 \\ x_i^T x_i = \|x_i\|_2^2 \end{array} \right.$$

Because $x_i^T w_k = w_k^T x_i$, $y_i^2 = 1$, and $\|x_i\|_2 \leq R$:

$$\|w_{k+1}\|_2^2 = \|w_k\|_2^2 + 2y_i w_k^T x_i + \|x_i\|_2^2 \leq \|w_k\|_2^2 + R^2$$

1.3) From problem 1.1, we know that $w_{k+1}^T w_{opt} \geq w_k^T w_{opt} + \gamma \|w_{opt}\|_2$

By accumulating M mistakes, we get:

$$w_{k+1}^T w_{opt} \geq w_0^T w_{opt} + \gamma M \|w_{opt}\|_2$$

Since w_0 is an all-zero vector:

$$w_{k+1}^T w_{opt} \geq \gamma M \|w_{opt}\|_2$$

According to the Cauchy-Schwartz inequality:

$$\gamma M \|w_{opt}\|_2 \leq w_{k+1}^T w_{opt} \leq \|w_{k+1}\|_2 \|w_{opt}\|_2$$

$$\Rightarrow \gamma M \|w_{opt}\|_2 \leq \|w_{k+1}\|_2 \|w_{opt}\|_2 \Rightarrow \underline{\gamma M \leq \|w_{k+1}\|_2}$$

$$a = w_{k+1}$$

$$b = w_{opt}$$

By applying the result from problem 1.2, we get

$$\|w_{k+1}\|_2^2 \leq \|w_k\|_2^2 + R^2 \Rightarrow \|w_{k+1}\|_2^2 \leq \|w_0\|_2^2 + R^2 M$$

Since w_0 is an all-zero vector:

$$\|w_{k+1}\|_2^2 \leq R^2 M \Rightarrow \underline{\|w_{k+1}\|_2 \leq R\sqrt{M}}$$

By combining the results above, we can conclude that:

$$\gamma M \leq \|w_{k+1}\|_2 \leq R\sqrt{M}$$

* side note (more steps if needed)

$$\begin{aligned} w_{k+1}^T w_{opt} &\geq w_k^T w_{opt} + \gamma \|w_{opt}\|_2 = \\ &\geq (w_{k-1}^T w_{opt} + \gamma \|w_{opt}\|_2) + \gamma \|w_{opt}\|_2 \\ &\geq [(w_{k-2}^T w_{opt} + \gamma \|w_{opt}\|_2) + \gamma \|w_{opt}\|_2] + \gamma \|w_{opt}\|_2 \\ &\vdots \\ &\geq w_0^T w_{opt} + \gamma M \|w_{opt}\|_2 \end{aligned}$$



side note (more steps if needed)

From 1.2 $\|w_{i+1}\|_2^2 \leq \|w_i\|_2^2 + R^2$

$$\begin{aligned} i=0 \quad & \|w_1\|_2^2 \leq \|w_0\|_2^2 + R^2 \\ i=1 \quad & \|w_2\|_2^2 \leq \|w_1\|_2^2 + R^2 \leq (\|w_0\|_2^2 + R^2) + R^2 = \|w_0\|_2^2 + 2R^2 \\ i=2 \quad & \|w_3\|_2^2 \leq \|w_2\|_2^2 + R^2 \leq (\|w_1\|_2^2 + R^2) + R^2 \leq (\|w_0\|_2^2 + R^2) + R^2 \\ & \vdots \\ i=M-1 \quad & \|w_M\|_2^2 \leq \|w_0\|_2^2 + MR^2 \quad (M \text{ mistakes}) \\ i=k \quad & \|w_{k+1}\|_2^2 \leq \|w_0\|_2^2 + MR^2 \quad \because w_{i+1} = w_i \text{ for } i \geq M \end{aligned}$$

Q 1.4)

a. given $\gamma \mu \leq \|w_{k+1}\|_2 \leq R\sqrt{\mu}$,
square sides

$$\gamma^2 \mu^2 \leq \|w_{k+1}\|_2^2 \leq R^2 \mu$$

the \leq sign does not change

b. Since $\gamma > 0 \Rightarrow \gamma > 0$, now
we divide all sides by γ^2

$$\frac{\mu^2}{\gamma^2} \leq \frac{\|w_{k+1}\|_2^2}{\gamma^2} \leq \frac{R^2 \mu}{\gamma^2}$$

c. divide by μ all sides

$$\frac{\mu}{\gamma^2 \mu} \leq \frac{\|w_{k+1}\|_2^2}{\gamma^2 \mu} \leq \frac{R^2}{\gamma^2} \quad \begin{matrix} \text{we have can} \\ \Rightarrow \therefore \end{matrix}$$
$$\mu \leq \frac{R^2}{\gamma^2}$$

2.1) Let $B(a_1, b_1, a_2, b_2)$ be the smallest rectangle returned by function f and B contains all positive examples in the training set S .

By using the realizability assumption, there exists a function f^* that correctly classifies all datapoints in the distribution D , meaning that $R(f^*) = 0$. Since the training set S is a subset of the dataset, $\hat{R}(f^*)$ is also zero. To prove f is an empirical risk minimizer, show that $\hat{R}(f) = 0$.

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(f(x_i) \neq y_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(f(x_i) \neq f^*(x_i))$$

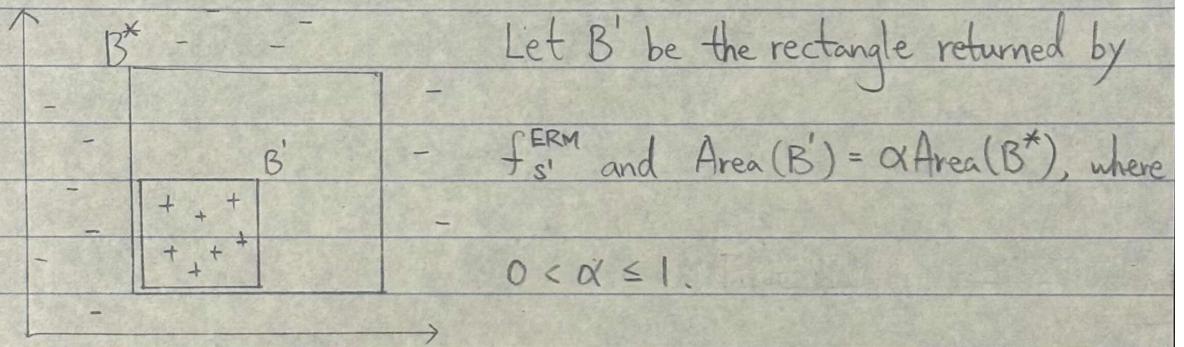
where n is the size of the training set S . Let n_{pos} be the # of mismatches for positive datapoints and n_{neg} for negative datapoints. \because all positive examples are enclosed within $B \therefore n_{\text{pos}} = 0$

$$\hat{R}(f) = \frac{1}{n} (n_{\text{pos}} \mathbb{I}(f(x_i) \neq 1) + n_{\text{neg}} \mathbb{I}(f(x_i) \neq 0))$$

$\because f^*$ exists \therefore no negative example exists within $B \Rightarrow n_{\text{neg}} = 0$

In conclusion, $\hat{R}(f) = \frac{1}{n}(0+0) = 0$, meaning that f is an ERM.

2.2) Assume the area of B^* is $\text{Area}(B^*)$ and the area outside of B^* is $\text{Area}(B^{\text{out}})$. Next, assume each datapoint is evenly distributed in the distribution D , the probability of picking a positive example is $P_{\text{pos}} = \frac{\text{Area}(B^*)}{\text{Area}(B^*) + \text{Area}(B^{\text{out}})} > 0$



Prove
 $R(f_{s'}^{\text{ERM}}) \geq 0.5$
 is possible

$$R(f_{s'}^{\text{ERM}}) = \frac{\text{Area}(B^*) - \text{Area}(B')}{\text{Area}(B^*) + \text{Area}(B^{\text{out}})} = \frac{(1-\alpha) \times \text{Area}(B^*)}{\text{Area}(B^*) + \text{Area}(B^{\text{out}})}$$

$$= (1-\alpha) P_{\text{pos}}$$

$\therefore 0 < \alpha \leq 1$ and $P_{\text{pos}} > 0 \quad \therefore R(f_{s'}^{\text{ERM}}) \geq 0.5$ is possible

2.3)

Training set f_s^{ERM} → gets a small error

if $n \geq \frac{4 \log(4/\delta)}{\epsilon}$ → prob at least: $R(f_s^{\text{ERM}}) \leq \epsilon$

Assumption:



$$a_i \geq a_i^* \rightarrow P(B_1) = \frac{\epsilon}{4}, B_1 = B(a_1^*, a_1, a_2^*, b_2^*)$$

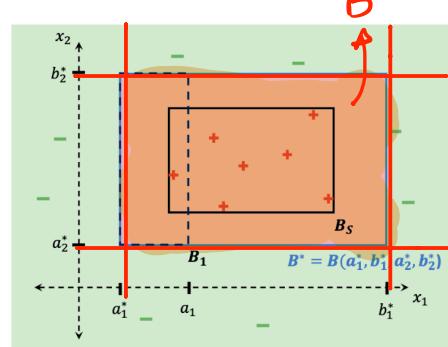
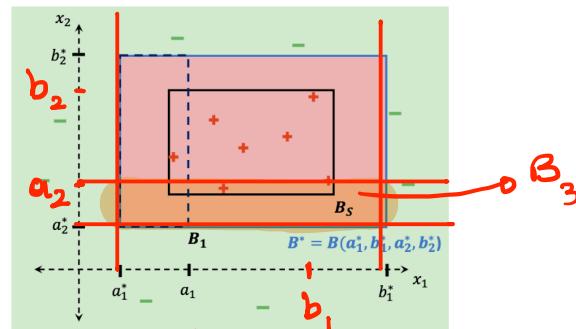
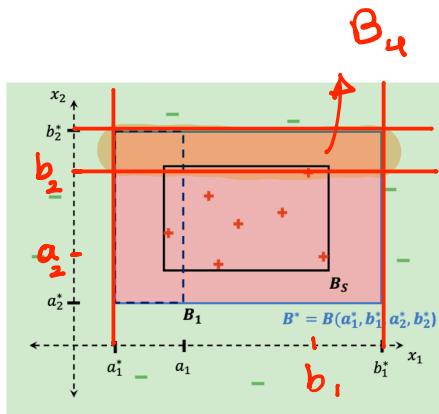
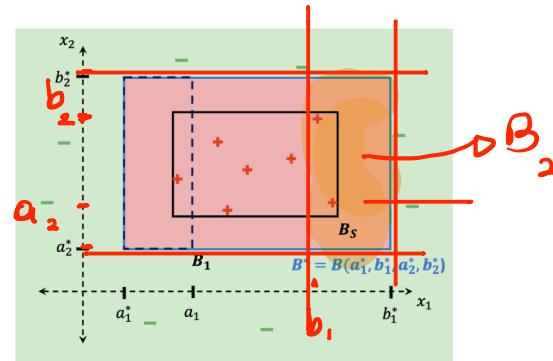
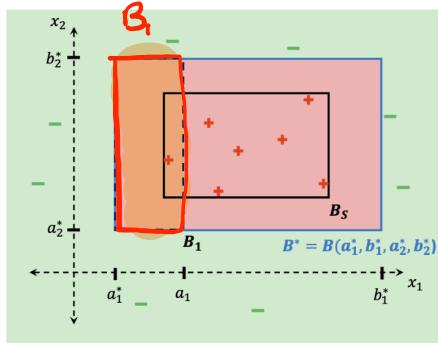
$$\begin{cases} b_1 \leq b_1^* \\ a_2 \geq a_2^* \\ b_2 \leq b_2^* \end{cases} \rightarrow P(B_2) = \frac{\epsilon}{4}, B_2 = B(b_1, b_1^*, a_2^*, b_2^*)$$

B_S → rectangle
on training

f_s^{ERM}

$$\rightarrow P(B_3) = \frac{\epsilon}{4}, B_3 = B(a_1^*, b_1^*, a_2^*, a_2^*)$$

$$\rightarrow P(B_4) = \frac{\epsilon}{4}, B_4 = B(a_1^*, b_1^*, b_2^*, b_2^*)$$



$$B^* = B(a_1^*, b_1^*, a_2^*, b_2^*)$$

for $f(a_1^*, b_1^*, a_2^*, b_2^*)$

- ①
- Assume we have a distribution D for (x, y) and $B^* = B(a_1^*, b_1^*, a_2^*, b_2^*)$ be the rectangle for $f(a_1^*, b_1^*, a_2^*, b_2^*)$. By the realizability assumption, B^* perfectly separates \oplus and \ominus cases.
 - we showed (earlier section) B_s returns all positive examples and it is the smallest rectangle.
 - According to the steps a and b which is given or proved in other sections $\Rightarrow B_s \subseteq B^*$
- ②
- According to given conditions (*) above all B_1, B_2, B_3, B_4 are inside B^* and we showed $B_s \subseteq B^*$
 - If s contains positive examples in all rectangles B_1, B_2, B_3, B_4 , then B_s will contain all positive examples in B_1, B_2, B_3, B_4 .
 - we know mass prob of each B_i is $\frac{\epsilon}{4}$.
 - from steps a,b,c $\Rightarrow R(f_i) \leq \frac{\epsilon}{4} \forall i \in \{1, 2, 3, 4\}$
 $\Rightarrow R(f_s^{ERM}) \leq 4 \left(\frac{\epsilon}{4} \right) \Rightarrow \therefore R(f_s^{ERM}) \leq \epsilon$

Probability sample
belong $B^* - B_s$
Set diff
 $\sum_{i=1}^4 \frac{\epsilon}{4} = 4 \left(\frac{\epsilon}{4} \right)$
*

$$\textcircled{3} \quad \text{a) } P(\text{sample has no point from any } B_i) = 1 - P(\text{sample } \in B_i) \\ = 1 - \left(\frac{E}{4} \right)$$

b) expanding for all rectangles $B_i, i \in \{1, 2, 3, 4\}$

$$P(\text{sample } \notin \text{ for all rectangles}) = \left(1 - \frac{E}{4}\right)^n \times e^{-\frac{E}{4}n}$$

we know $1 - E \approx e^{-E} \dots$

$$E = \frac{E}{4}$$

④ a) we know $R(f_s^{\text{ERM}}) \leq \epsilon$ and know a bad sample may be found in S .

b) $P(R(f_s^{\text{ERM}}) \leq \epsilon) = P(S \text{ is a good sample})$
 $= 1 - P(S \text{ is bad sample})$

side note:

on side we know $\Pr[B_1 \cup B_2 \cup B_3 \cup B_4] \leq \Pr \sum_{i=1}^4 \Pr[B_i]$

$$\geq 1 - \sum_{i=1}^4 e^{-\frac{\epsilon}{n} n}$$

$$\geq 1 - 4e^{-\frac{\epsilon}{n} n}$$

c) we want to show $P(R(f_s^{\text{ERM}}) \leq \epsilon)) \geq 1 - \delta$

$$P(R(f_s^{\text{ERM}}) \leq \epsilon)) \geq 1 - \delta$$

$$1 - 4(e^{-\frac{n\epsilon}{4}}) \geq 1 - \delta \quad \text{multiply both sides by } (-1)$$

$$4(e^{-\frac{n\epsilon}{4}}) \leq \delta$$

$$e^{\frac{n\epsilon}{4}} \geq \frac{4}{\delta} \quad \text{take ln both sides}$$

$$n(\frac{\epsilon}{4}) \geq \ln(\frac{4}{\delta})$$

$$\therefore n \geq \frac{4 \ln(\frac{4}{\delta})}{\epsilon}$$

Q3

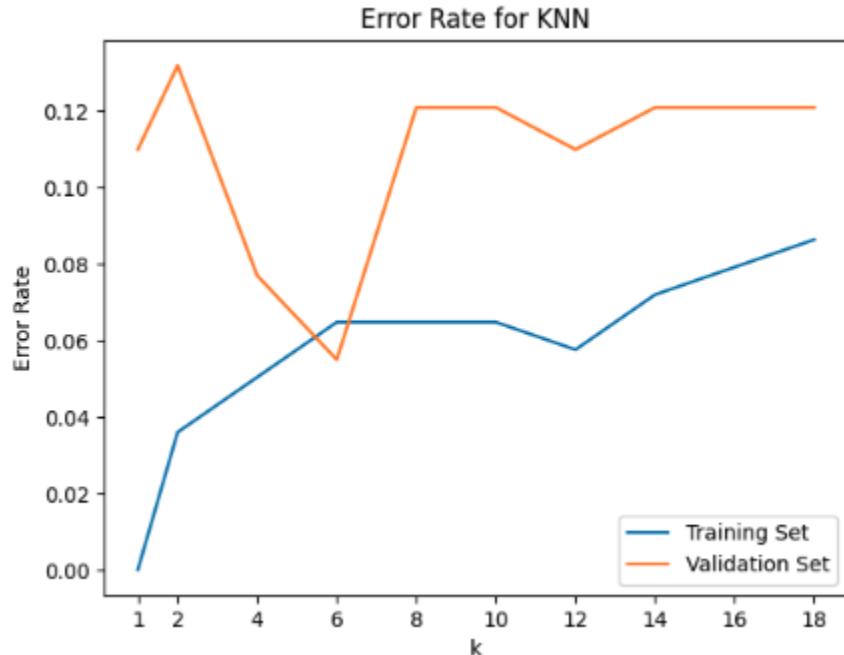
Print statements from code:

The validation error rate is 0.07692307692307693 in Problem **Set 1.1**

The validation error rate is 0.04395604395604396 in Problem Set **1.2 when using normalization**

The validation error rate is 0.04395604395604396 in Problem Set **1.2 when using minmax_scaling**

The validation error rate is 0.04395604395604396 in Problem Set **1.3, which use cosine distance**



In Problem Set 1.4, we use the best $k = 6$ with the best validation error rate 0.054945054945054944

Using the best k , the final test error rate is 0.07100591715976332

3.4 report:

(1) Report and draw a curve based on the error rate of your model on the training set for each k. What do you observe? (2pts)

- a. When $k = 1$, we observe the lowest error.
- b. As k increases and finding the k -nearest neighbors (knn) becomes more complex, the error rate is higher compared to the first few iterations.
- c. There is a drop at $k = 12$. For $6 \leq k \leq 20$, $k = 12$ exhibits the lowest error.
- d. The error line remains flat between $k = 6$ and $k = 10$.
- e. As the value of k increases in KNN, the model may become less sensitive to local patterns and more influenced by global trends. This could impact the interpretability and generalization of the model.
- g. When $k=0$, the error is the lowest equal to zero. However, this is not what we are looking for.

(2) Report and draw a curve based on the error rate of your model on the validation set for each k. What is your best k? (2pts) \n

- a. The best k is 6 having the lowest Error Rate.
- b. For $8 \leq k \leq 10$ and $14 \leq k \leq 18$ the error rate stays the same. Potentially, we can have an early stop after 10 if the error rate does not improve.
- c. From going $k=2$ to $k=6$, the error rate improves

(3) What do you observe by comparing the difference between the two curves? (2pts) \n

- a. Validation graph may provide better insight into the decision-making process leading to the conclusion than training graph.
- b. The training curve tends to increase for $2 \leq k \leq 6$ whereas the validation curve has a drop.
- d. A similar curve pattern has been observed for both training and validation curves when $8 \leq k \leq 14$. However, training and validation for each k have a different error rate.
- e. For $k=6$ (the best k - question 2) the validation error rate is lower than training error rate.
- f. broadly speaking, as k increases in the training set, the error rate increases or remains the same.

g. In the Validation graph, we have a convex function for k between 2 and 8. k=6 is the global minimum.

h. Both validation and training have a convex function for k between 10 and 14, with a local min of k=12

(4) What is the final test set error rate you get using your best-k? (1pt) \n

a. In Problem Set 1.4, we use the best k = 6 with the best validation error rate 0.054945054945054944

Using the best k, the final test error rate is 0.07100591715976332

(5) Comment on these results from the perspective of overfitting, generalization and hyper-parameter tuning. (3pts).\n

Overfitting:

When k=1, the error rate on the training set is low and higher on the validation set, this shows the model does not perform well on unseen data due to overfitting. \n

For small values of k, such as k=1 in KNN, the prediction depends solely on the nearest neighbor, which can lead to overfitting. The model becomes highly sensitive to noise in the training data, capturing local fluctuations rather than the true underlying structure of the data. \n

Generalization:

For k=6, the training and validation error rates are very close. This shows the model is generalized enough to perform the same for unseen data.

k=6 is less sensitive to a single point of data and still preserves the local pattern in the data.

k=6 provides a good balance between overfitting and underfitting in terms of generalization.

hyper-parameter tuning:

Using different Ks to find the best k is a hyper parameter tuning which we implemented in this problem. As a result, we found the best k is 6, which has the lowest error rate. \n

We learned a hyperparameter tuning is a necessary step for finding the best parameters in a classification problem.

It is crucial to have validation set in hyperparameter tuning to avoid overfitting and detect cases such as k=1, as mentioned above.

Q4

4.1

$F(w_{LS}) = 217.48452613174004$ on training data

$F(w_0) = 78885.82819617869$ on training data

$F(w_{LS}) = 294.06836989399164$ on testing data

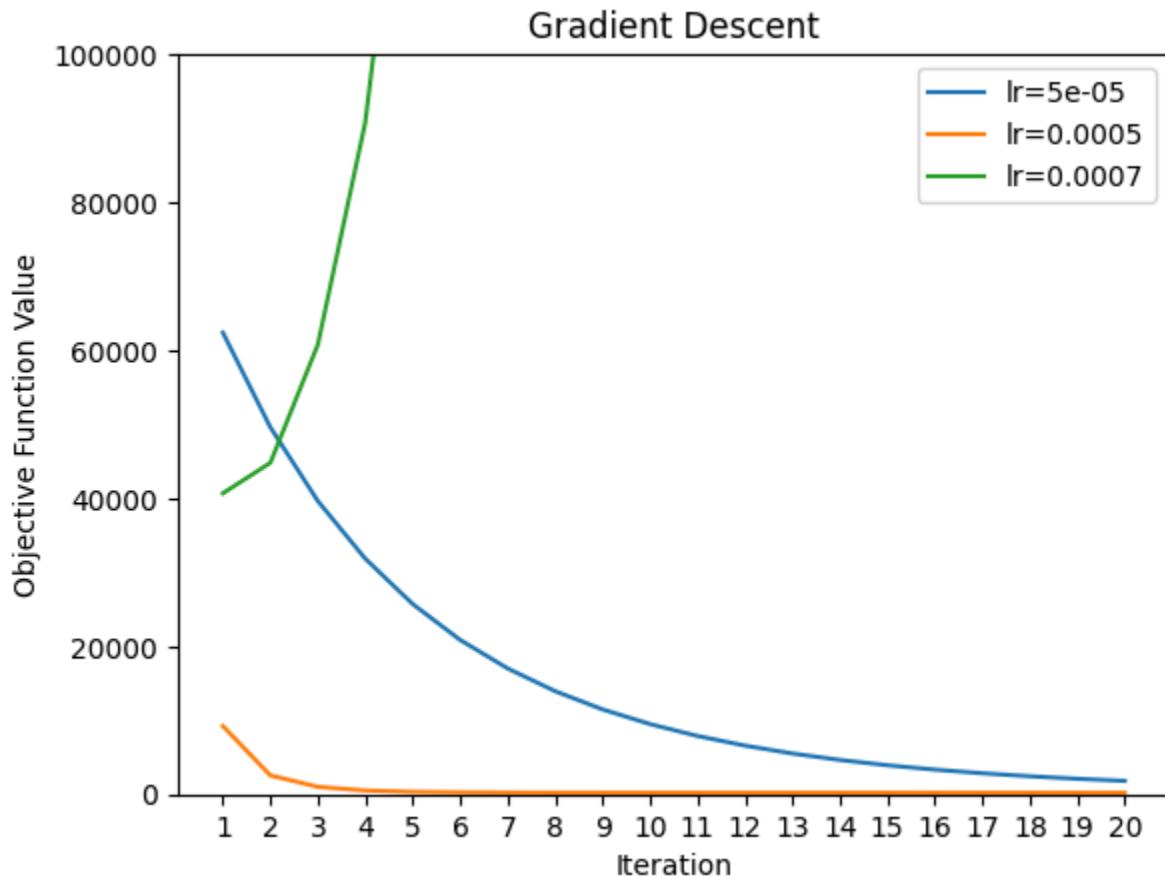
Gap: $294.068 - 217.485 = 76.583$

The test objective function value with weights w_{LS} performed much better than setting the weights w_0 to a zero vector, showing the effectiveness of w_{LS} in minimizing the residual sum of squares (RSS). However, there was still a gap between the training and test data, suggesting that the model with weights w_{LS} may be slightly overfitting the training data.

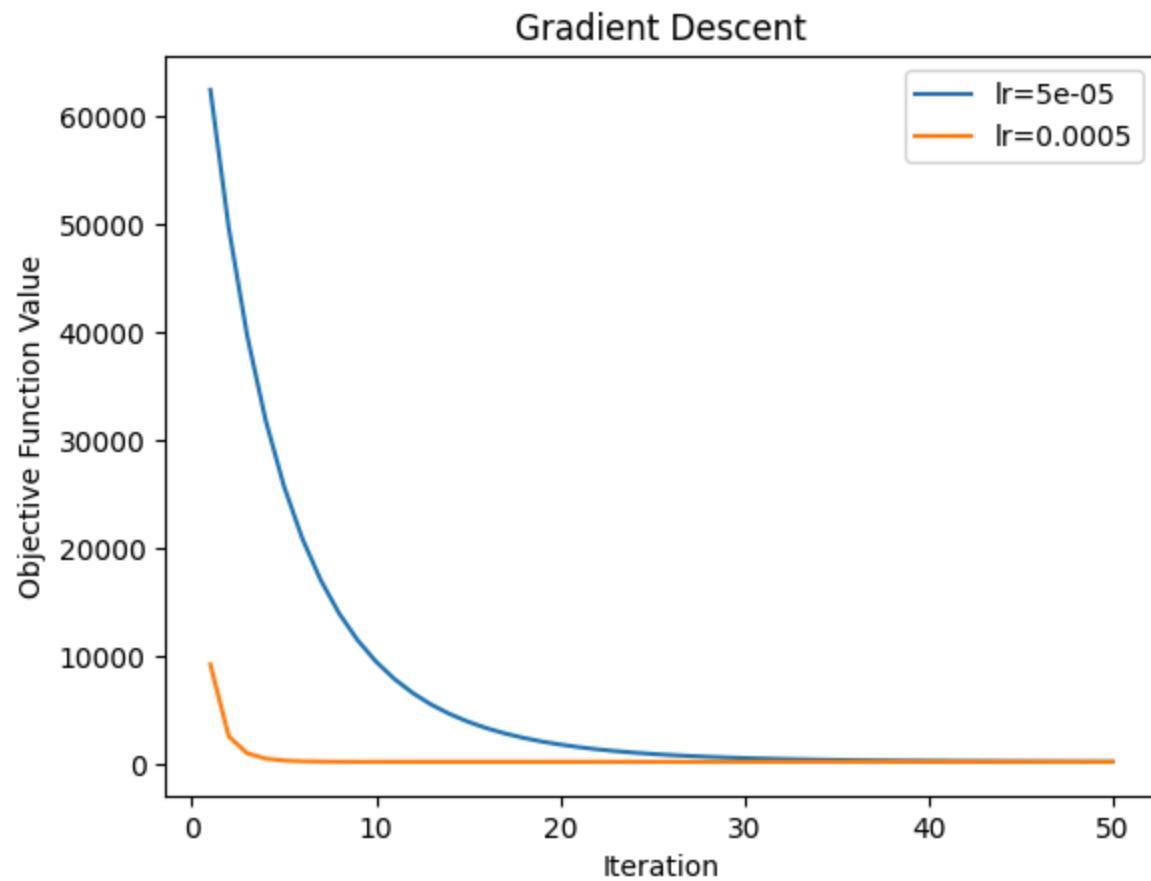
4.2

A larger step size led to faster convergence in gradient descent. For instance, with a step size of 0.0005, convergence occurred in approximately seven iterations. In contrast, when the step size was reduced to 0.00005, it took over 30 iterations to approach convergence, and the objective function value remained significantly higher than that of 0.0005. However, surpassing a step size of 0.000582 resulted in a divergence of the objective function value after 100 iterations. In conclusion, the step size should be neither too small nor too large, and increasing the number of iterations does not guarantee a lower objective function value.

Best final objective value: 217.486, where step size=0.0005

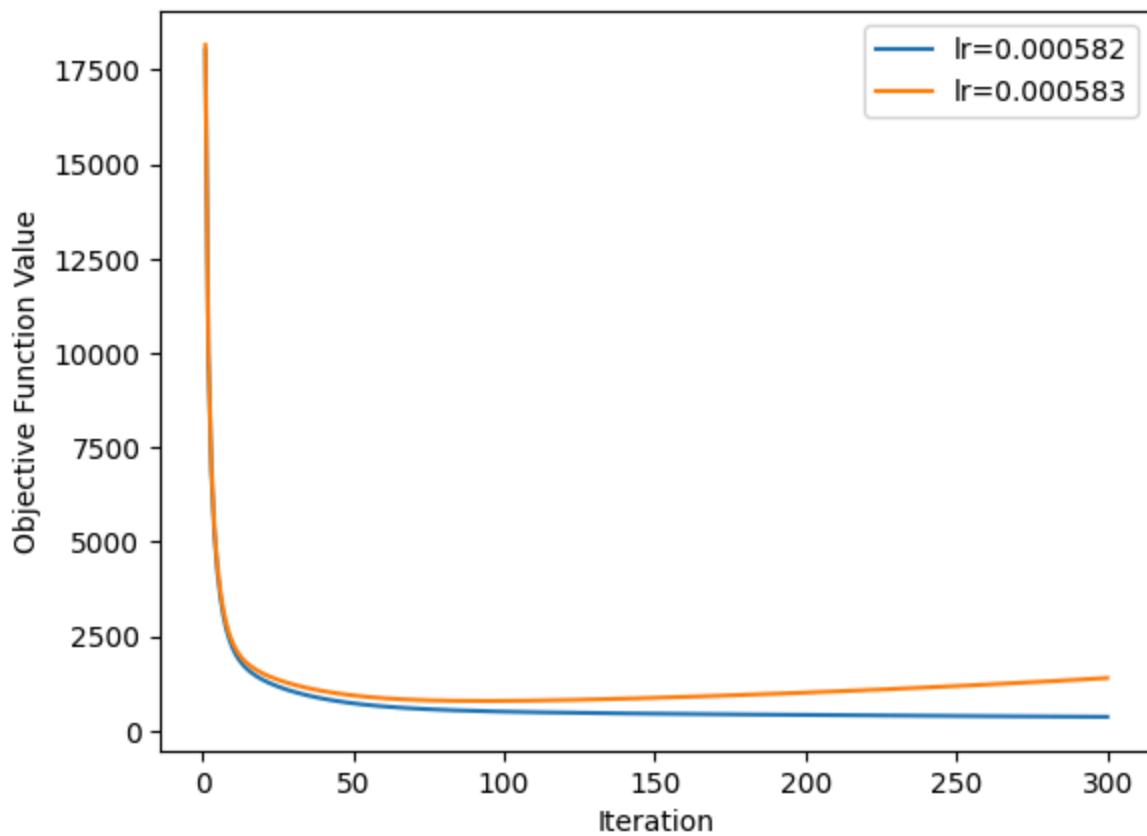


Comparison of gradient descent convergence for various step sizes



The convergence comparison graph for step sizes of 0.00005 and 0.0005

Gradient Descent



The divergence boundary with 300 iterations

4.3

Like gradient descent, stochastic gradient descent showed faster convergence with a larger step size, and the optimal step size should be neither too small nor too large. However, stochastic gradient descent had an erratic convergence behavior because only one random data point was used in each iteration instead of the entire data set. Also, we found that SGD required more iterations than GD to converge because randomly choosing a data point for each iteration resulted in noisy gradient estimates.

Best final objective value: 443.185, where step size = 0.005

