

Question 4

a

Let D be a $M \times n \in \{0,1\}$ binary matrix representing which observations were generated from which Regression Model. For example the matrix:

$$\begin{pmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}$$

Is a 2×3 that represents a dataset of 3 observations belonging to one of 2 Regression Models. Observations 1 and 2 belong to Regression Model 2 and observation 3 belongs to Regression Model 1.

Let \hat{D} be the binary matrix of the above format labelling the predictions of Regression Models that each observation belongs to.

$$loss(\hat{D}, \hat{\beta}) = \sum_{j=1}^M \left[\sum_{i=1}^n [\hat{D}_{ji} * [(X_i^T \hat{\beta}_j - Y_i)^2 + (1 - D_{ji})]] \right]$$

The first term $(X_i^T \hat{\beta}_j - Y_i)^2$ is a penalty for having errors in $\hat{\beta}_j$. The second term $1 - D_{j,i}$ is a penalty for prediction observation i in the incorrect D_j . Note that the observation is skipped and no loss is applied if the observation does not belong to the predicted regression model, which is why the loss for each observation is multiplied by \hat{D}_{ji} .

In words, the loss function reads: For each observation i that is predicted to be in a regression model \hat{D}_j , find the mean squared error of that observation and add a penalty if the observation's predicted \hat{D} differed from its actual D .

When $M = 1$, this just simplifies down to:

$$loss(\hat{\beta}) = \sum_{i=1}^n ((X_i^T \hat{\beta} - Y_i)^2)$$

b

```
# First partition all observations
For each observation i:
  For each model m:
    -> calculate the absolute distance between the model (m) prediction
        on this observation (i) and the actual value of this observation

    -> assign this observation to the model that it
        had the \emph{smallest} error with

# Next find the coefficients of each model
For each partitioned dataset d:
  -> fit a linear regression model that uses MSE as its loss function
      and has the dataset (d) as its training data

# We now have values for D and \(\hat{B}\) that are minimised and the loss
# can be calculated with my function in a.
```

c

```
20 def total_loss(X, y, Z, models): | You, 24 minutes ago • code for 4
21     loss = 0
22     M = len(models)
23     n = X.shape[0]
24
25     for j in range(M):
26         model = models[j]
27
28         for i in range(n):
29             # If the observation is not partitioned to model "j", we skip
30             if Z[i] != j:
31                 continue
32
33             # Get penalty of the coefficients for this model on this observation
34             observation = X[i]
35             actual = y[i]
36             prediction = model.predict(np.array([observation]))
37             coef_penalty = (prediction[0] - actual)**2
38
39             # Cannot know partition_penalty since partition predictions
40             # are not provided to this function. Not really sure of the
41             # point of this function. Just set penalty to 0.
42             partition_penalty = 0
43
44             loss += coef_penalty + partition_penalty
45
46     return loss
47
48 def q4c():
49     data = import_data(Q4_TRAIN_DATA_DIR)
50     Xtrain = np.array(data.drop(columns=["Y"]))
51     ytrain = np.array(data["Y"])
52
53     mod = LinearRegression().fit(Xtrain, ytrain)
54     Z = np.zeros(shape=Xtrain.shape[0]) # all points would belong to a single partition.
55     print(total_loss(Xtrain, ytrain, Z, [mod])) # outputs 298.328178158043
56
```

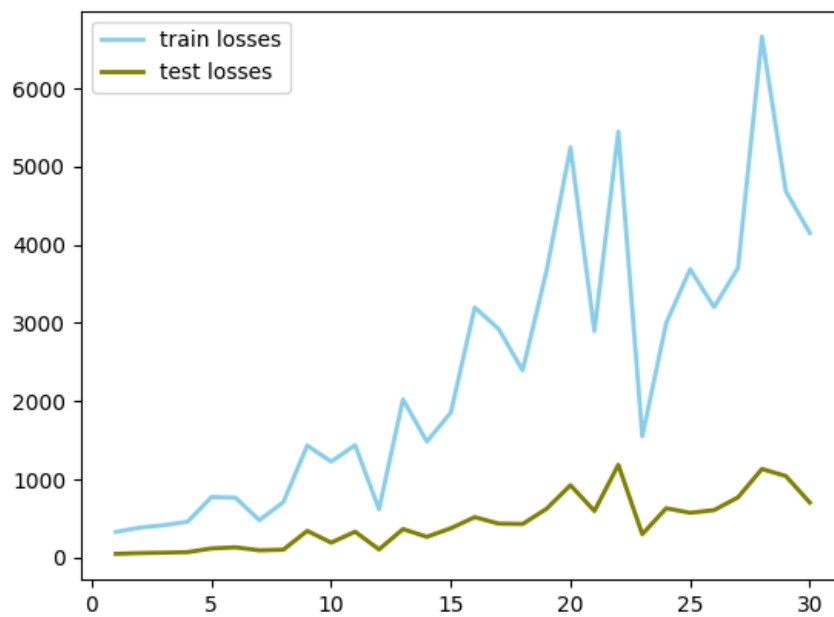
d

```
57 def find_partitions(X, y, models):      You, 25 minutes ago • code for
58     M = len(models)
59     n = X.shape[0]
60
61     Z = []
62     for i in range(n):
63         observation = X[i]
64
65         # predict this observation on each model. Whichever model
66         # predicted the closest to the actual we assume this
67         # datapoint belongs to this model.
68         pred_diffs = []
69         for model in models:
70             pred = model.predict(np.array([observation]))
71             actual = y[i]
72
73             # Could also square this to get a positive number
74             pred_diffs.append(abs(pred - actual))
75
76         # Append index of smallest value to list
77         Z.append(pred_diffs.index(min(pred_diffs)))
78
79     return np.array(Z)
80
81 def q4d():
82     data = import_data(Q4_TEST_DATA_DIR)
83     Xtest = np.array(data.drop(columns=["Y"]))
84     ytest = np.array(data["Y"])
85
86     # Run with one model
87     mod = LinearRegression().fit(Xtest, ytest)
88
89     Z = find_partitions(Xtest, ytest, [mod])
90     print(total_loss(Xtest, ytest, Z, [mod]))
```

e

See code for this question here 0.1

M	train	test
5	773.2139301322158	116.60294714118255
10	1226.975556530428	190.77503856677967
15	1854.7758161852978	374.3318184723467
20	5245.153194657786	924.8003322754047
25	3687.1371589765313	572.1619888958278
30	4150.443300036157	703.5325244593158



Based on the results, I would guess that the most reasonable value for M is around 6 or 8. Even though the train losses were lower when $M < 5$, the test losses were approximately the same for M from 0 to 8 and a lower value for M would be favoured since each of the models gets trained on more data.

Appendix

0.1 q4e

```
97 # So that we don't have to train thousands of models
98 def generate_Z(M, n):    You, seconds ago • code for 4
99
100     # np.tile is such a misdirect, you have to find how many
101     # times you want to repeat and then append two lists together.
102     #
103     # Sorry, I've been going at this exam for years and am getting
104     # to the end of a tether.
105     Z = []
106     for i in range(n):
107         Z.append(i % M)
108
109     return np.array(Z)
110
111 def q4e():    You, a minute ago • code for 4
112     # import train and test data
113     data_train = import_data(Q4_TRAIN_DATA_DIR)
114     Xtrain = np.array(data_train.drop(columns=["Y"]))
115     ytrain = np.array(data_train["Y"])
116
117     data_test = import_data(Q4_TEST_DATA_DIR)
118     Xtest = np.array(data_test.drop(columns=["Y"]))
119     ytest = np.array(data_test["Y"])
120
121     # model
122     model = LinearRegression()
123
124     # Iterate over values of M
125     train_losses = []
126     test_losses = []
127
128     print("M |          train          |          test          |")
129     MAX_M = 30
130     for M in range(1, MAX_M + 1):
131         # Partition data (done in a mechanical way for exam testing)
132         Ztrain = generate_Z(M, Xtrain.shape[0])
133         Ztest = generate_Z(M, Xtest.shape[0])
134
135         # Train models on their own data
136         train_models = []
137         for i in range(M):
138             indices = np.where(Ztrain == i)
139             fitted = model.fit(Xtrain[indices], ytrain[indices])
140             train_models.append(fitted)
```

```

141
142     test_models = []
143     for i in range(M):
144         indices = np.where(Ztest == i)
145         fitted = model.fit(Xtest[indices], ytest[indices])
146         test_models.append(fitted)
147
148     # Calculate loss
149     train_loss = total_loss(Xtrain, ytrain, Ztrain, train_models)
150     test_loss = total_loss(Xtest, ytest, Ztest, test_models)
151
152     train_losses.append(train_loss)
153     test_losses.append(test_loss)
154
155     if M % 5 == 0:
156         padding=""
157         if M == 5:
158             padding = " "
159         print(F"{M}{padding} | {train_loss} | {test_loss} |")
160
161     # Plot everything
162     X = list(range(1, MAX_M + 1))
163     plt.plot(X, train_losses, color="skyblue", label="train losses", linewidth=2)
164     plt.plot(X, test_losses, color="olive", label="test losses", linewidth=2)
165     plt.legend()
166     plt.savefig("outputs/q4e.png")

```