

DATA7202 Assignment 2

Longpeng Xu

April 20, 2023

Contents

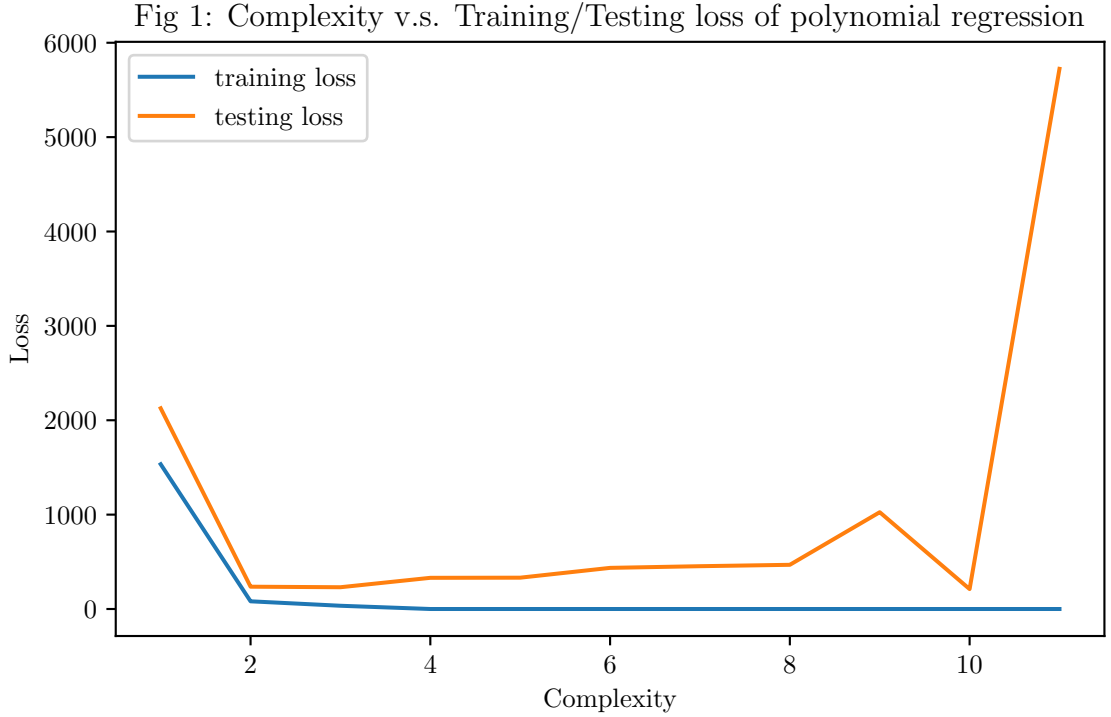
1	Question 1	1
2	Question 2	2
3	Question 3	3
4	Question 4	3
5	Question 5	3
6	Question 6	5
7	Question 7	6
8	Question 8	6
9	Appendix	9

1 Question 1

- Zero training loss can be achieved by a model which over-fits (the training data).
- An example is to generate $\tau = \{(x_i, y_i), i = 1, \dots, 6\}$ (x_i 's are unique) and varying the complexity (the highest degree k) of polynomial regression. When $k \geq 5$, the training loss is constant at zero (See Fig 1 and Code in Appendix).
- However, the model over-fits which is evidenced by the testing loss: When $k = 12$, the testing loss surges to a MSE of 5723.6 (See Fig 1). This is because when the training loss is zero i.e., when $k \geq 5$ and keeps increasing, the model becomes increasingly complex and fits to the training data very closely - It 'memorizes' the training data and disables itself to generalize to any test set.
- Hence, in general, for any training set $\tau = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ with unique \mathbf{x}_i values, one can propose a model $f : \mathbf{x}_i \rightarrow y_i$ such that $f(\mathbf{x}_i) = y_i$ holds for every unique \mathbf{x}_i . The training loss, of any types,

becomes zero. An example of training loss in the type of MSE:

$$\text{training loss} = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - y_i)^2 = 0.$$



2 Question 2

The squared-error loss is in the form of $\sum_{i=1}^k (y_i - u)^2$, set $u = h(\mathbf{x})$ be some unknown prediction function. To find the u that minimizes the squared error loss, set the derivative w.r.t. u to be zero.

$$\begin{aligned}
 \frac{d}{du} \sum_{i=1}^k (y_i - u)^2 &= \sum_{i=1}^k \frac{d}{du} (y_i - u)^2 \\
 &= \sum_{i=1}^k 2(y_i - u) \frac{d}{du} (y_i - u) \\
 &= \sum_{i=1}^k 2(y_i - u)(-1) \\
 &= -2 \sum_{i=1}^k (y_i - u) \\
 &= -2 \left(\sum_{i=1}^k y_i - ku \right). \\
 \frac{d}{du} \sum_{i=1}^k (y_i - u)^2 &= -2 \left(\sum_{i=1}^k y_i - ku \right) = 0 \Rightarrow \sum_{i=1}^k y_i - ku = 0 \Rightarrow u^* = \frac{1}{k} \sum_{i=1}^k y_i.
 \end{aligned}$$

Since $\frac{d}{du} \sum_{i=1}^k (y_i - u)^2 = -2 \left(\sum_{i=1}^k y_i - ku \right) \begin{cases} < 0, & u \in (-\infty, \frac{1}{k} \sum_{i=1}^k y_i) \\ > 0, & u \in (\frac{1}{k} \sum_{i=1}^k y_i, +\infty) \end{cases}$. Hence, u^* is the local and global minimum.

Hence, $h^w(\mathbf{x}) := \frac{1}{k} \sum_{i=1}^k y_i = u^*$ minimizes the squared-error loss.

3 Question 3

- Note that τ has n points and τ^* has the same size as τ . For bootstrapping, we sample one point each time, by n times which are independent (i.e. with replacement)
- Hence, for every point, the probability that it is sampled in a single time is $\frac{1}{n}$. Hence, the probability that a point is not sampled in a single time is $1 - \frac{1}{n}$
- Hence, the probability that a point is not sampled in all the n times (i.e., a point is not in τ^*) is $(1 - \frac{1}{n})^n$
- For large n ,

$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n \stackrel[t \rightarrow +\infty \text{ as } n \rightarrow -\infty]{t = -n, t \rightarrow -\infty \text{ as } n \rightarrow +\infty}}{\lim_{t \rightarrow \infty} \left(1 + \frac{1}{t}\right)^{-t}} = \lim_{t \rightarrow \infty} \frac{1}{\left(1 + \frac{1}{t}\right)^t} = \frac{1}{e} \approx 0.37,$$

i.e., the probability that a point is not in τ^* converges to e^{-1} .

- Denote $X_i = \begin{cases} 1 & \text{if the } i\text{th point of } \tau \text{ is not in } \tau^* \\ 0 & \text{otherwise} \end{cases}, i = 1, \dots, n.$

Apparently, $X_i \sim^{\text{approx}} \text{Bern}(e^{-1})$ for large n .

Set $X = \sum_{i=1}^n X_i$ as # of the points not in τ^* . Since X_i 's are i.i.d., $X \sim^{\text{approx}} \text{Bin}(n, e^{-1})$ for large n .

Hence, for large n , the proportion of the points from τ but not in $\tau^* \approx \frac{\mathbb{E}X}{n} = \frac{ne^{-1}}{n} = e^{-1}$

\Leftrightarrow For large n , τ^* does not contain a fraction of about e^{-1} of the points from τ

4 Question 4

See random forest regressor code in Appendix.

The optimal parameter $m = 9$ with the highest $R^2 = 0.724804$.

5 Question 5

See gradient boosting code in Appendix. The plot, Fig 5-2, is in the next page.

Conclusion: As is shown by Fig 5-2,

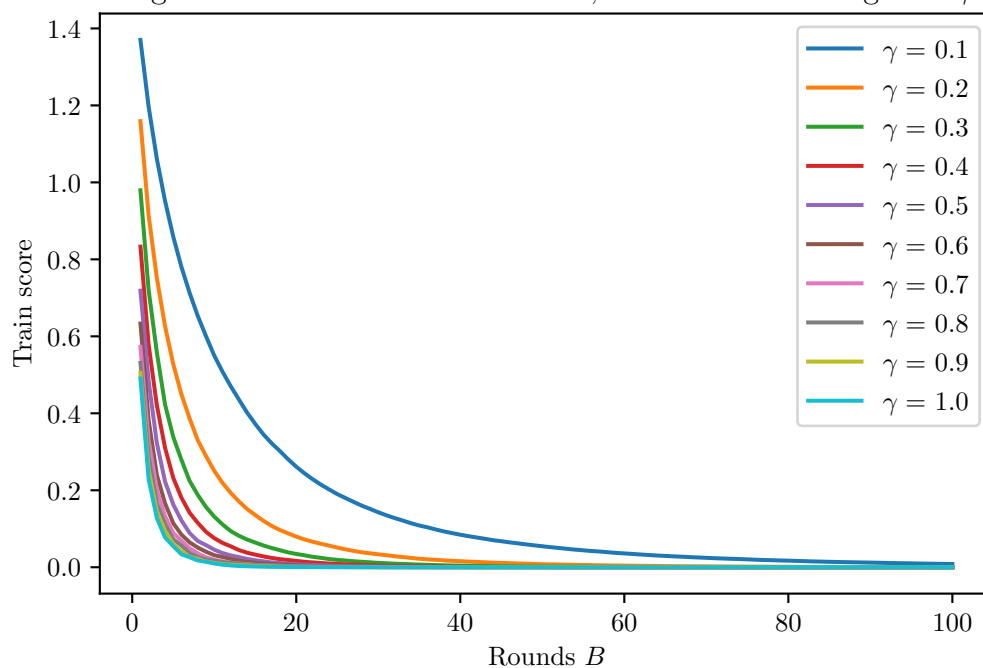
- The greater the B , the smaller the train scores (i.e. losses) at all levels of γ , and vice versa.
- The greater the γ , the smaller the train scores at all levels of B , and vice versa.

Note In `GradientBoostingClassifier`, `max_depth` is set so that there is no over-fit to the training set.

If not setting `max_depth`, it is default to be 3, which can lead to a sudden increase in train score where γ

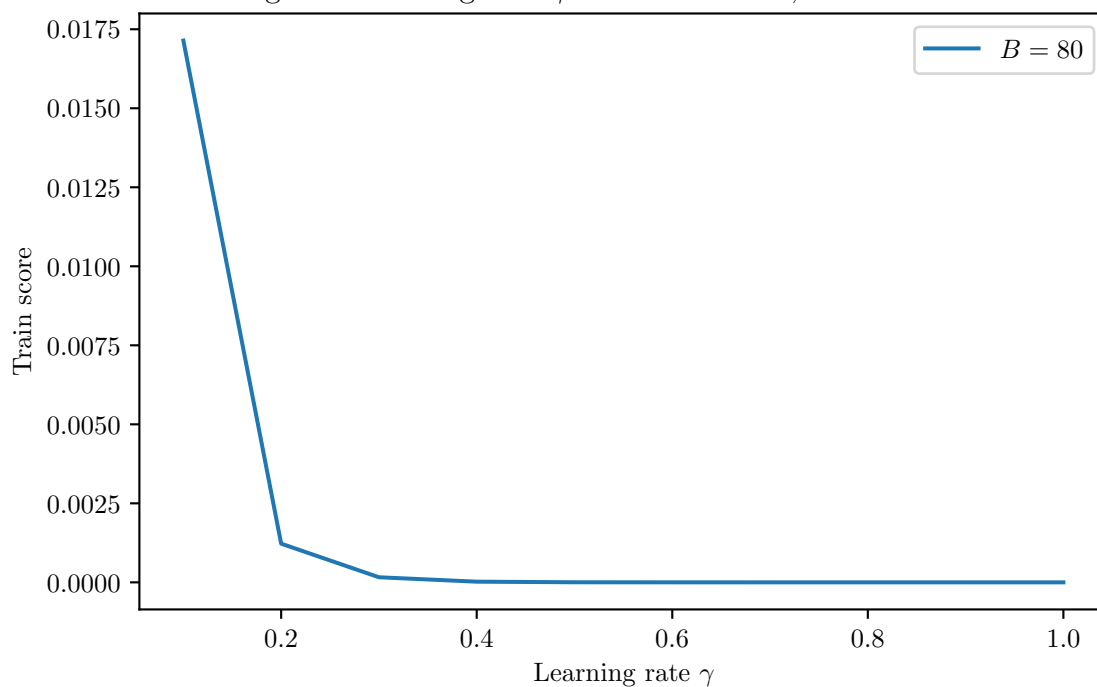
is large in this case.

Fig 5-2: Round B v.s. Train score, at different learning rate γ



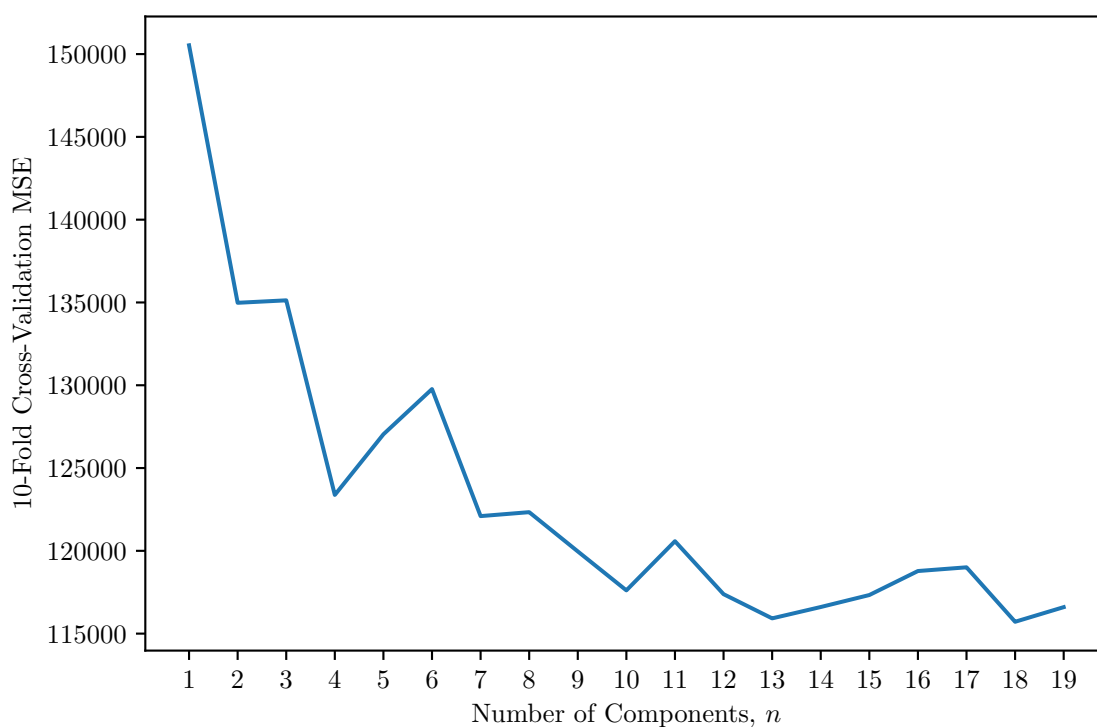
When B is fixed to 80 rounds, the plot of γ v.s. train score is shown in Fig 5-1. (Not sure whether this is required.)

Fig 5-1: Learning rate γ v.s. Train score, $B = 80$ rounds



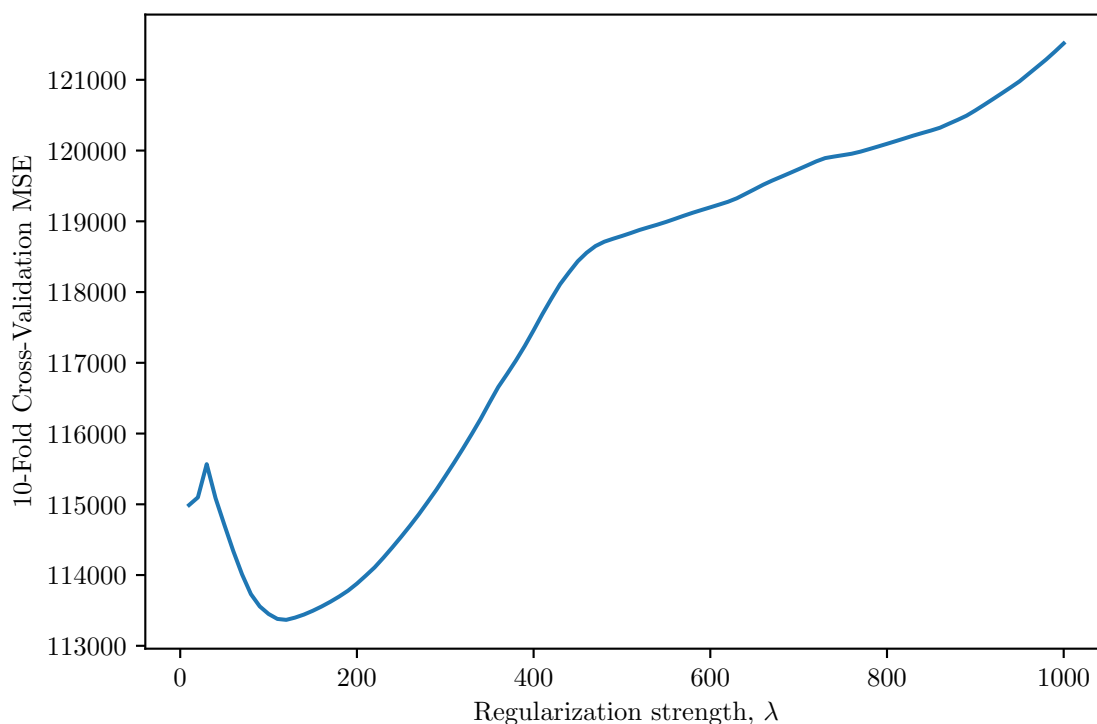
6 Question 6

(a) See PCR and 10-Fold Cross-Validation code in Appendix.



The optimal number of component is $n = 18$.

(b) See Lasso and 10-Fold Cross-Validation code in Appendix.



The best λ is 120, with the lowest MSE 113366.62.

Note The problem does not specify whether a constant should be included. However, if not including

the constant, the 10-Fold CV MSEs will be extremely high. Hence, it remains for Question 6.

7 Question 7

- (a) See Poisson regression code in Appendix.

	coef	lower bound (0.025)	upper bound (0.975)
type	-0.2237	-0.317	-0.130
construction	0.3714	0.255	0.488
operation	0.7680	0.567	0.969
months	8.095×10^{-5}	7.54×10^{-5}	8.65×10^{-5}

Note The problem does not specify whether a constant should be included. However, the added constant term can be regarded as insignificant (p -value > 0.1). Hence, it is removed for Question 7.

- (b) See bootstrapping code (with coefficients, standard errors and normal 95% CIs) in Appendix.

Difference: CIs with bootstrapping are narrower than CIs without bootstrapping.

Reason: Bootstrapping, the sampling with replacement, makes the bootstrapped sampled datasets similar to the original dataset. Hence, the coefficients of each bootstrapped dataset are similar to the coefficients of the original dataset. (Also, possibly, bootstrapping 34 observations for building 1000 bootstrapped datasets brings repeated, identical observations, i.e. even more similarity, to these datasets). By averaging the coefficients across the 1000 bootstrapped datasets, the estimates of the four coefficients are even more stable, leading to the lower standard deviations of the coefficients and hence, narrower CIs.

The CIs with bootstrapping:

	lower bound (0.025)	upper bound (0.975)
type	-0.215245	-0.199097
construction	0.393864	0.414539
operation	0.574860	0.621786
months	0.000096	0.000101

8 Question 8

- (a) See code for multiple regression in Appendix.

The fitted model is $Time = 2.3412 + 1.6159 \cdot Cases + 0.0144 \cdot Distance$

The estimated residual standard deviation is 3.259473

The p -value for overall model is 4.69×10^{-16} ;

The p -value for Cases is 3.254932×10^{-9} ;

The p -value for Distance is 6.312469×10^{-4}

Note The problem does not specify whether a constant should be included. However, the added constant term can be regarded as significant ($p\text{-value} < 0.05$). Hence, it remains for Question 8.

(b) See code for residual plot/histogram in Appendix.

The residual plots (Fig 8-1, ..., 8-5) are shown below.

- The linear relationship between the explanatory variables and the explained variable fails to be guaranteed, because the residuals does not show a random scatter of points around zero (the dashed horizontal line). In fact, some residuals significantly deviates from zero (Fig 8-1, 8-2, 8-3).
- The constant variance for the residuals fails to be guaranteed, since a as the predicted Y increases, the residuals are increasingly divergent, i.e. a funnel shape in Fig 8-3 which indicates heteroscedasticity.
- Independence among the residuals fail to be guaranteed, since the residuals plotted, by the order in which y_i were measured, are not randomly scattered but show some degree of correlation in Fig 8-4.
- The residuals fail to follow a normal distribution, since the residual histogram is not normally distributed (Fig 8-5).

Fig 8-1: Cases (X_1) v.s. Residuals

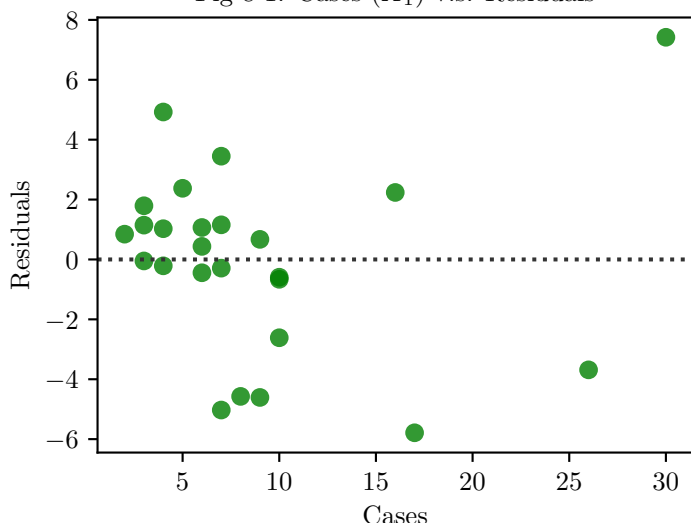


Fig 8-2: Distance (X_2) v.s. Residuals

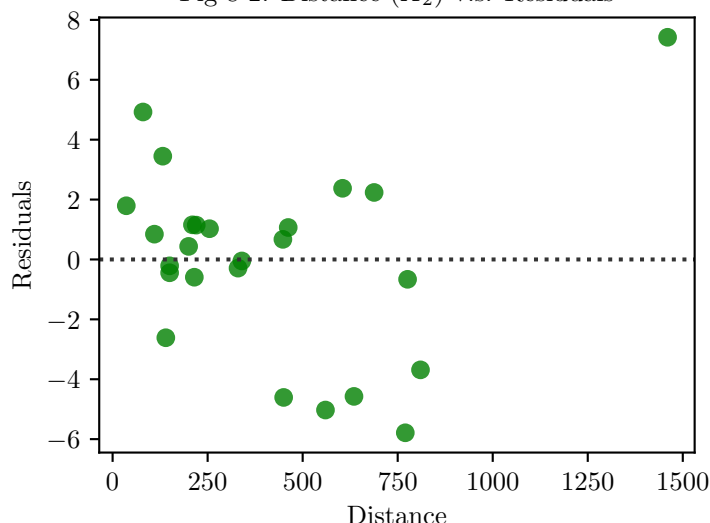


Fig 8-3: Predicted Time (Y) v.s. Residuals

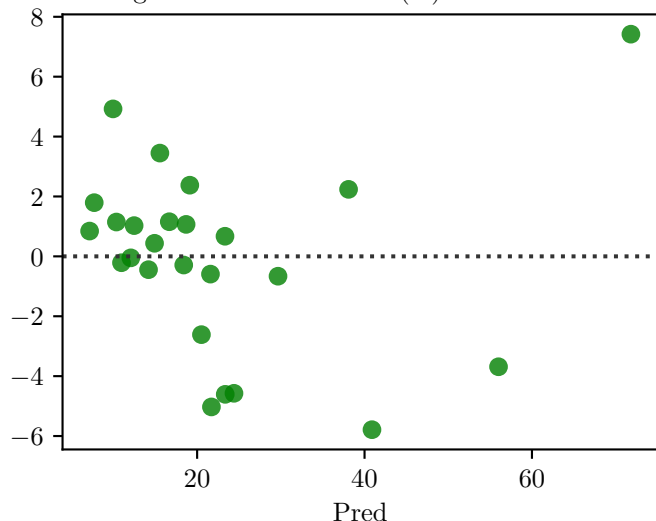
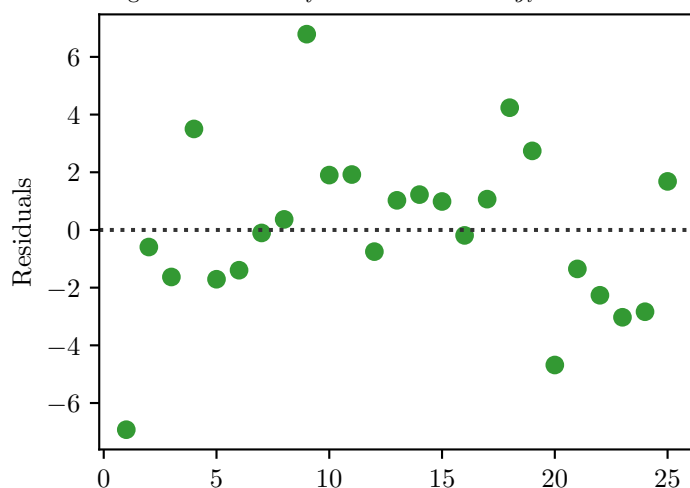
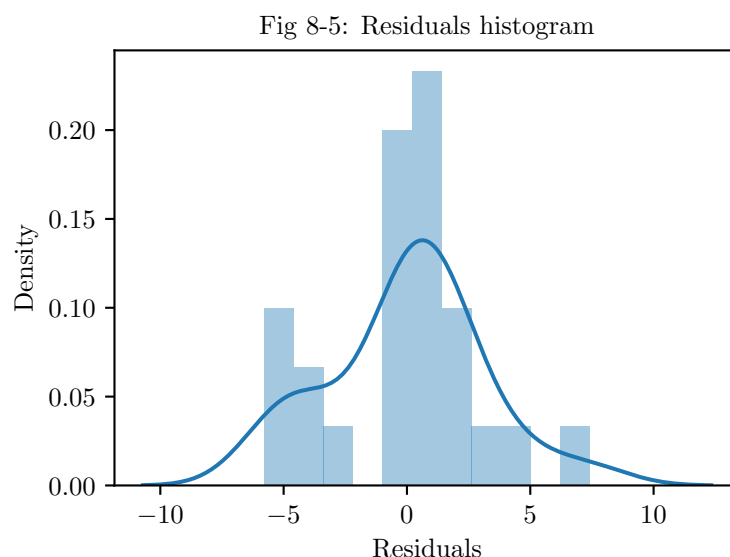


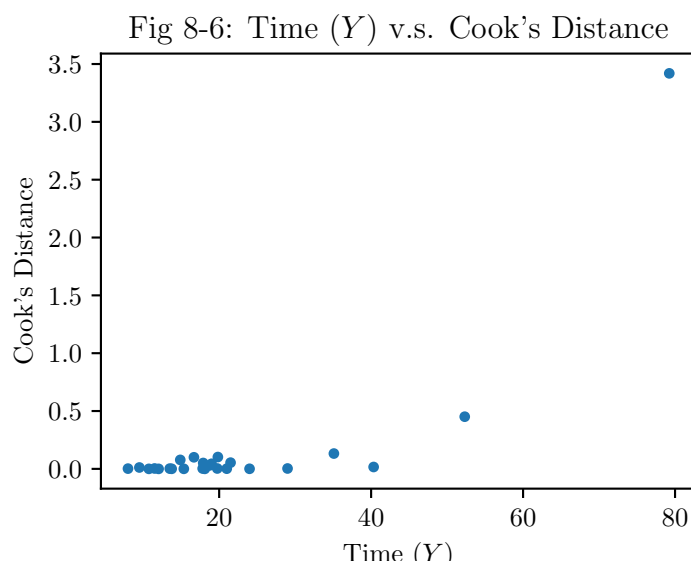
Fig 8-4: Residuals by the order in which y_i were measured





(c) See code for Cook's distance plot in Appendix.

The Cook's distance plot is below.



- The extremely influential observation, with Cook's distance 3.419318, is the row index 8 of `softdrink.csv`, with $Time = 79.24$, $Cases = 30$, $Distance = 1460$
- The next most influential, with Cook's distance 0.451045, is the row index 21 of `softdrink.csv`, with $Time = 52.32$, $Cases = 26$, $Distance = 810$

9 Appendix

Listing 1: Code for Questions

```

1  # -----
2  # Q1
3  # -----
4
5  import numpy as np
6  from sklearn.preprocessing import PolynomialFeatures
7  from sklearn.linear_model import LinearRegression
8  from sklearn.metrics import mean_squared_error
9  from sklearn.model_selection import train_test_split
10 import matplotlib.pyplot as plt
11
12 np.random.seed(seed=1)
13 x_tr = np.random.normal(0, 1, 6)
14 y_tr = 1*x_tr**6 + 2*x_tr**5 + 3*x_tr**4 + 4*x_tr**3 + 5*x_tr**2 + 6*x_tr
15 np.random.seed(seed=2)
16 x_te = np.random.normal(0, 1, 6)
17 y_te = 1.5*x_te**6 + 2*x_te**5 + 3*x_te**4 + 4*x_te**3 + 5*x_te**2 + 6*x_te + np.random.uniform(-2,2,6)
18
19 train_loss = []
20 test_loss = []
21
22 for i in range(1,12):
23     poly = PolynomialFeatures(degree=i, include_bias=False)
24     poly_features_tr = poly.fit_transform(x_tr.reshape(-1, 1))
25     poly_features_te = poly.fit_transform(x_te.reshape(-1, 1))
26     poly_reg_model = LinearRegression()
27     pr = poly_reg_model.fit(poly_features_tr, y_tr)
28     train_loss += [mean_squared_error(y_tr, pr.predict(poly_features_tr))]
29     test_loss += [mean_squared_error(y_te, pr.predict(poly_features_te))]
30
31 print('train_loss: ', train_loss, '\n test_loss: ', test_loss)
32
33 plt.rcParams.update({'text.usetex': True, 'font.family': 'serif'})
34 fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True, dpi=600, facecolor='white')
35
36 ax.plot( range(1,12), train_loss)
37 ax.plot( range(1,12), test_loss)
38
39 ax.legend(labels = ['training loss','testing loss'], loc='best')
40 ax.set_xlabel('Complexity')
41 ax.set_ylabel('Loss')
42 ax.set_title('Fig 1: Complexity v.s. Training/Testing loss of polynomial regression')
43 plt.savefig('7202a2fig1.pdf')
44 plt.show()
45
46
47
48 # -----

```

```

49 # Q4
50 # -----
51
52 import numpy as np
53 from sklearn.datasets import make_friedman1
54 from sklearn.ensemble import RandomForestRegressor
55 from sklearn.model_selection import train_test_split
56 from sklearn.metrics import r2_score
57 # create regression problem
58 n_points = 1000
59 x, y = make_friedman1(n_samples = n_points, n_features = 20 ,
60                       noise = 2.0 , random_state = 100)
61 # split to train /test set
62 x_train , x_test , y_train , y_test = \
63 train_test_split(x, y, test_size = 0.33 , random_state = 100)
64
65 def best_pred_size(m):
66     for i in m:
67         rf = RandomForestRegressor(n_estimators=500, #oob_score = True,
68                                   max_features=i, random_state=100)
69         rf.fit(x_train,y_train)
70         yhatrf = rf.predict(x_test)
71         print('max_features = ', i, ' '*4,
72               'RF R^2 score = ', round(r2_score(y_test, yhatrf), 6))
73     return None
74
75 best_pred_size(m= list( range(1,21)))
76
77
78
79 # -----
80 # Q5
81 # -----
82
83 from sklearn.datasets import make_blobs
84 from sklearn.metrics import zero_one_loss
85 from sklearn.model_selection import train_test_split
86 import numpy as np
87 import matplotlib.pyplot as plt
88 from sklearn.ensemble import GradientBoostingClassifier
89 import matplotlib.pyplot as plt
90 import numpy as np
91
92 if __name__ == "__main__":
93     X_train, y_train = make_blobs(n_samples=1000, n_features=20, centers=5,
94                                   random_state=100, cluster_std=6)
95
96 def grad_b_c_1(gamma, round_list):
97     train_score = []
98     for i in round_list:
99         breg = GradientBoostingClassifier(learning_rate = gamma,
100                                           n_estimators = i,

```

```

101                                     max_depth=2, random_state=100)
102     breg.fit(X_train,y_train)
103     train_score += [breg.train_score_[-1]]
104     return round_list, train_score
105
106 gamma_list = np.arange(0.1, 1.1, 0.1)
107 round_list = np.arange(1,101)
108
109 plt.figure(figsize=(4,3),
110             dpi=600, facecolor='white')
111 for j in gamma_list:
112     xpoints, ypoints = grad_b_c_1(gamma=j, round_list= round_list)
113     plt.plot(xpoints, ypoints)
114 plt.legend(labels=['$\gamma = \ $' + str( round(i,1)) for i in gamma_list], loc='best')
115 plt.xlabel('Rounds $B$')
116 plt.ylabel('Train score')
117 plt.title('Fig 5-2: Round $B$ v.s. Train score, at different learning rate $\gamma$')
118 plt.savefig('7202a2fig5-2.pdf')
119 plt.show()
120
121 def grad_b_c(gamma, rounds):
122     train_score = []
123     for i in gamma:
124         breg = GradientBoostingClassifier(learning_rate = i,
125                                           n_estimators = rounds,
126                                           max_depth=2, random_state=100)
127         breg.fit(X_train,y_train)
128         train_score += [breg.train_score_[-1]]
129     return gamma, train_score
130
131 grad_b_c(gamma_list, 80)
132
133 plt.rcParams.update({'text.usetex': True, 'font.family': 'serif'})
134 fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True, dpi=600, facecolor='white')
135 xpoints, ypoints = grad_b_c(gamma_list, 80)
136 ax.plot(xpoints, ypoints)
137 ax.legend(labels = ['$B = 80$'], loc='best')
138 ax.set_xlabel('Learning rate $\gamma$')
139 ax.set_ylabel('Train score')
140 ax.set_title('Fig 5-1: Learning rate $\gamma$ v.s. Train score, $B=80$ rounds')
141 plt.savefig('7202a2fig5-1.pdf')
142 plt.show()
143
144
145
146 # -----
147 # Q6
148 # -----
149
150 import pandas as pd
151 import numpy as np
152 from sklearn import preprocessing

```

```

153 from sklearn.linear_model import LinearRegression
154 from sklearn.model_selection import KFold
155 from sklearn.metrics import mean_squared_error
156 from sklearn.decomposition import PCA
157 import matplotlib.pyplot as plt
158 import warnings
159 warnings.filterwarnings("ignore")
160 from sklearn import linear_model
161
162 hitters = pd.read_csv('Hitters.csv', header=0)
163 le = preprocessing.LabelEncoder()
164 change = ['League', 'Division', 'NewLeague']
165 for colname in change:
166     hitters[colname] = le.fit_transform(hitters[colname])
167 X_ = hitters.drop(['Salary'], axis=1)
168 y_ = hitters['Salary']
169
170 def Validate(X,Y):
171     kf = KFold(n_splits = 10)
172     kf.get_n_splits(X)
173     mse_10fold = []
174     for i in range(1, X_.shape[1] + 1):
175         mse_fold = []
176         for train_index, test_index in kf.split(X):
177             X_train, X_test = X_.iloc[train_index, :], X_.iloc[test_index, :]
178             y_train, y_test = Y.iloc[train_index], Y.iloc[test_index]
179             X_train_pca, X_test_pca = pca_user(X_train=X_train, X_test=X_test, ncomp=i)
180             lm = LinearRegression().fit(X_train_pca, y_train)
181             y_pred = lm.predict(X_test_pca)
182             mse_fold.append(mean_squared_error(y_test, y_pred))
183         mse_10fold.append( round(np.mean(mse_fold),2))
184     return pd.DataFrame({'ncomp': range(1, X_.shape[1] + 1), 'mse_10fold': mse_10fold})
185
186 def pca_user(X_train, X_test, ncomp):
187     pca = PCA(n_components = ncomp)
188     X_train_pca = pca.fit_transform(X_train)
189     X_test_pca = pca.transform(X_test)
190     return X_train_pca, X_test_pca
191
192 Q6a = Validate(X_, y_)
193 Q6a
194
195 plt.rcParams.update({'text.usetex': True, 'font.family': 'serif'})
196 fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True, dpi=600, facecolor='white')
197
198 xpoints, ypoints = Q6a.iloc[:,0], Q6a.iloc[:,1]
199 ax.plot(xpoints, ypoints)
200
201 ax.set_xlabel('Number of Components, $n$')
202 ax.set_ylabel('10-Fold Cross-Validation MSE')
203 plt.xticks(xpoints, xpoints)
204 plt.savefig('7202a2fig3.pdf')

```

```

205 plt.show()
206
207 Q6a[Q6a.mse_10fold == Q6a.mse_10fold. min()]
208
209 def Validate_1(X,Y):
210     kf = KFold(n_splits = 10)
211     kf.get_n_splits(X)
212     mse_10fold = []
213     for i in np.arange(10,1010,10):
214         mse_fold = []
215         lasso = linear_model.Lasso(fit_intercept=True, alpha=i)
216         for train_index, test_index in kf.split(X):
217             X_train, X_test = X.iloc[train_index, :], X.iloc[test_index, :]
218             y_train, y_test = Y.iloc[train_index], Y.iloc[test_index]
219             lasso.fit(X_train, y_train)
220             y_pred = lasso.predict(X_test)
221             mse_fold.append(mean_squared_error(y_test, y_pred))
222         mse_10fold.append( round(np.mean(mse_fold),2))
223     return pd.DataFrame({'Lambda': np.arange(10,1010,10), 'Ten_Fold_CV_MSE': mse_10fold})
224
225 Q6b = Validate_1(X_, y_)
226 Q6b
227
228 plt.rcParams.update({'text.usetex': True, 'font.family': 'serif'})
229 fig, ax = plt.subplots(figsize=(6, 4), tight_layout=True, dpi=600, facecolor='white')
230
231 xpoints, ypoints = Q6b.iloc[:,0], Q6b.iloc[:,1]
232 ax.plot(xpoints, ypoints)
233
234 ax.set_xlabel(r'Regularization strength,  $\lambda$ ')
235 ax.set_ylabel('10-Fold Cross-Validation MSE')
236 plt.savefig('7202a2fig4.pdf')
237 plt.show()
238
239 Q6b[Q6b.Ten_Fold_CV_MSE == Q6b.Ten_Fold_CV_MSE. min()]
240
241
242
243 # -----
244 # Q7
245 # -----
246
247 import pandas as pd
248 import statsmodels.api as sm
249
250 ships = pd.read_csv('ships.csv', header=0)
251 ships
252
253 exog, endog = ships.drop('damage', axis=1), ships['damage']
254 pois = sm.GLM(endog, exog, family=sm.families.Poisson())
255 res = pois.fit()
256 res.summary()

```

```

257
258 n_estimators = 1000
259 params = pd.DataFrame(columns=['type', 'construction', 'operation', 'months'])
260
261 for i in range(n_estimators):
262     np.random.seed(seed=i)
263     ids = np.random.choice( range( len(ships)), size = len(ships), replace=True)
264     exog, endog = ships.loc[ids,:].drop('damage', axis=1), ships.loc[ids,'damage']
265     pois_res = sm.GLM(endog, exog, family=sm.families.Poisson()).fit()
266     params.loc[i,:] = list(pois_res.params)
267
268 params
269
270 lower_bound = list(params.mean(axis=0) - 1.96 * params.std(axis=0) / np.sqrt(1000))
271 upper_bound = list(params.mean(axis=0) + 1.96 * params.std(axis=0) / np.sqrt(1000))
272
273 params_1 = pd.DataFrame(index = ['type', 'construction', 'operation', 'months'],
274                           columns=['coef', 'std err', '[0.025', '0.975]'])
275 params_1.loc[:, 'coef'] = list(params.mean(axis=0))
276 params_1.loc[:, 'std err'] = list(params.std(axis=0))
277 params_1.loc[:, '[0.025'] = lower_bound
278 params_1.loc[:, '0.975'] = upper_bound
279 params_1
280
281
282
283 # -----
284 # Q8
285 # -----
286
287 import pandas as pd
288 import statsmodels.api as sm
289 import matplotlib.pyplot as plt
290 import seaborn as sns
291
292 softdrink = pd.read_csv('softdrink.csv', header=0)
293 softdrink
294
295 X = sm.add_constant(softdrink.drop('Time', axis=1))
296 Y = softdrink['Time']
297 lr = sm.OLS(Y, X).fit()
298 print(lr.summary())
299
300 lr.params
301 lr.pvalues
302 lr.resid
303
304 softdrink['Pred'] = lr.predict(X)
305 softdrink['Residuals'] = softdrink['Time'] - softdrink['Pred']
306 softdrink
307
308 # first method to calculate the estimated residual standard deviation

```

```

309 lr.mse_resid**0.5
310 # second method to calculate the estimated residual standard deviation
311 ( sum(softdrink['Residuals']**2) / (25-2-1))**0.5
312
313 plt.rcParams.update({'text.usetex': True, 'font.family': 'serif'})
314 plt.figure(figsize=(4,3),dpi=600, facecolor='white')
315 sns.residplot(x=softdrink['Cases'], y=lr.resid, lowess=False, color='g', label='a')
316 plt.title('Fig 8-1: Cases ($X_1$) v.s. Residuals', fontsize=10)
317 plt.savefig('7202a2fig8-1.pdf')
318 plt.show()
319
320 plt.figure(figsize=(4,3),dpi=600, facecolor='white')
321 sns.residplot(x=softdrink['Distance'], y=lr.resid, lowess=False, color='g', label='a')
322 plt.title('Fig 8-2: Distance ($X_2$) v.s. Residuals', fontsize=10)
323 plt.savefig('7202a2fig8-2.pdf')
324 plt.show()
325
326 plt.figure(figsize=(4,3),dpi=600, facecolor='white')
327 sns.residplot(x=softdrink['Pred'], y=lr.resid, lowess=False, color='g', label='a')
328 plt.title('Fig 8-3: Predicted Time ($Y$) v.s. Residuals', fontsize=10)
329 plt.savefig('7202a2fig8-3.pdf')
330 plt.show()
331
332 plt.figure(figsize=(4,3),dpi=600, facecolor='white')
333 sns.residplot(x= list( range(1, len(softdrink)+1)), y=lr.resid, lowess=False, color='g', label='a')
334 plt.title('Fig 8-4: Residuals by the order in which $y_i$ were measured', fontsize=8)
335 plt.savefig('7202a2fig8-4.pdf')
336 plt.show()
337
338 plt.figure(figsize=(4.2,3),dpi=600, facecolor='white')
339 sns.distplot(softdrink['Residuals'])
340 plt.title('Fig 8-5: Residuals histogram', fontsize=10)
341 plt.savefig('7202a2fig8-5.pdf')
342 plt.show()
343
344 cooks = lr.get_influence().cooks_distance
345 cooks
346 Q8c = pd.DataFrame({'Y': softdrink['Time'], 'Cooks': cooks[0]})
347 Q8c
348
349 plt.figure(figsize=(4,3),dpi=600, facecolor='white')
350 plt.scatter(Q8c.iloc[:,0], Q8c.iloc[:,1], marker='.')
351 plt.xlabel('Time ($Y$)')
352 plt.ylabel('Cook\'s Distance')
353 plt.title('Fig 8-6: Time ($Y$) v.s. Cook\'s Distance')
354 plt.savefig('7202a2fig8-6.pdf')
355 plt.show()
356
357 Q8c.sort_values('Cooks',ascending=False)
358 softdrink.iloc[Q8c.sort_values('Cooks',ascending=False).index[:2],:]

```