

# Homework 6

Han Yong Wunrow  
Student No: 1877282  
AMATH 584

Due: Monday, December 14, 2020

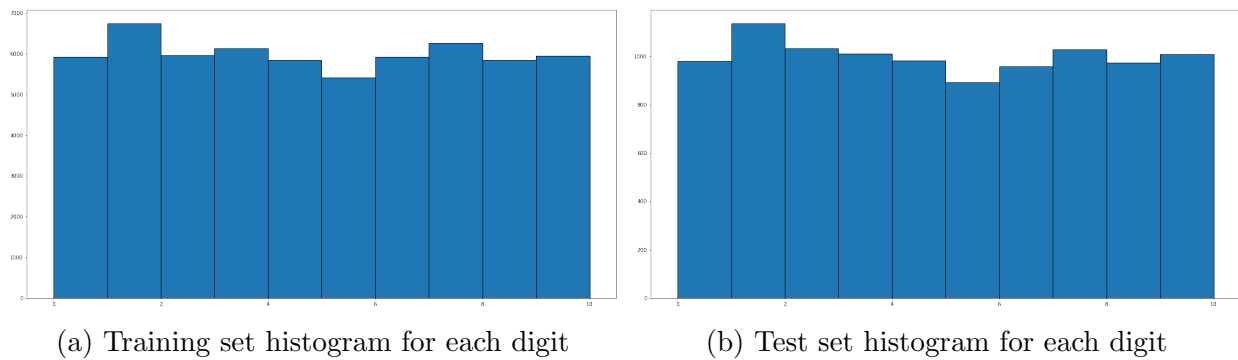
## Abstract

In this report, we build a linear classifier on the MNIST database of handwritten digits. We classify an image as either a 1,2,3,4,5,6,7,8,9, or 0 using logistic regression with both LASSO (using the  $\ell_1$ -norm), Ridge (using the  $\ell_2$ -norm), and Elastic Net (using both the  $\ell_1$ - and  $\ell_2$  norm). We use Grid Search to tune the hyperparameters  $\lambda_1$ ,  $\lambda_2$ , and `l1_ratio` using  $k$ -fold crossvalidation with  $k = 3$ . We rank each pixel using the absolute value of the coefficients for each of the models and subset the training and test sets to the top 100 pixels. When evaluating on the full test set, LASSO had the highest test accuracy with 92.62%, followed by Ridge with 92.54% and Elastic Net with 92.52%. However, after subsetting to the top 100 pixels, Elastic Net had the highest accuracy of 87.5% followed by Ridge with 88.3% and LASSO with 79.7%. We used Python 3 and Google Collaboratory for this analysis, and code is provided on [GitHub](#).

The MNIST database contains 60,000 images of handwritten digits (0 through 9) that are all labeled in a training set and 10,000 labeled images in a test set. Each image is 28 by 28 pixels.

We reshape each image in our training set as a  $784 \times 1$  column vector  $x_i$  and load into a matrix  $A$ . So each column in  $A$  represents an image, and  $A$  has 784 rows and 60,000 columns. We perform a one-hot encoding of the labels so that  $y_i = e_k$  (the standard basis vector) corresponds to image  $i$  having been labeled as the integer  $(k - 1)$  for  $0 \leq k \leq 9$ . We load each column vector  $y_i$  into a matrix  $B$ . So  $B$  has 10 rows and 60,000 columns. So  $A^T X = B^T$  is an overdetermined system where  $X$  has 784 rows and 10 columns. Since we are including an intercept in our model, we add an additional column  $(1, 1, \dots, 1)^T$  to  $A^T$  and  $X$ . So our feature space are the 784 pixels.

Both the training and test sets are uniformly distributed with similar counts for each digit so we do not have to worry about balancing. Also these data sets are generalizable to real-world scenarios where each digit is equally likely to appear.



Each image has been centered and the pixel value ranges from 0 to 255. We normalize the training and test sets by dividing by the max value 255.



Figure 2: First 10 digits in training set

## 2 Logistic Regression with LASSO, Ridge, and Elastic Net

We classify each image in our training set as either a 1,2,3,4,5,6,7,8,9, or 0 using logistic regression with LASSO (using the  $\ell_1$ -norm), Ridge (using the  $\ell_2$ -norm), and Elastic Net (using both the  $\ell_1$ - and  $\ell_2$  norm). The respective optimization problems for LASSO, Ridge, and Elastic Net are

$$\max_{\beta} \sum_{i=1}^n \sum_{k=1}^9 \mathbb{1}_{y_i=k} \log P(y_i = k \mid x_i, \beta) - \lambda_1 \|\beta\|_1 \quad (1)$$

$$\max_{\beta} \sum_{i=1}^n \sum_{k=1}^9 \mathbb{1}_{y_i=k} \log P(y_i = k \mid x_i, \beta) - \lambda_1 \|\beta\|_2^2 \quad (2)$$

$$\max_{\beta} \sum_{i=1}^n \sum_{k=1}^9 \mathbb{1}_{y_i=k} \log P(y_i = k \mid x_i, \beta) - \lambda_1 \|\beta\|_1 - \lambda_2 \|\beta\|_2^2 \quad (3)$$

where the conditional probabilities are computed using softmax

$$P(y = k \mid x, \beta) = \frac{\exp(\beta_k^T x)}{\sum_{j=1}^9 \exp(\beta_j^T x)}.$$

Our coefficients  $\beta$  are used to predict the label  $\hat{y}$  for a given image  $x$  via

$$\hat{y} = \operatorname{argmax}_k \frac{\exp(\hat{\beta}_k^T x)}{\sum_{j=1}^9 \exp(\hat{\beta}_j^T x)}$$

Equations (1),(2),and(3) are rewritten using scikit-learn's objectives

$$\max_{\beta} \sum_{i=1}^n \sum_{k=1}^9 \mathbb{1}_{y_i=k} \log P(y_i = k \mid x_i, \beta) - \frac{1}{C} \|\beta\|_1 \quad (4)$$

where  $\lambda_1 = 1/C$ ,

$$\max_{\beta} \sum_{i=1}^n \sum_{k=1}^9 \mathbb{1}_{y_i=k} \log P(y_i = k \mid x_i, \beta) - \frac{1}{C} \|\beta\|_2^2 \quad (5)$$

where  $\lambda_2 = 1/C$ , and

$$\max_{\beta} \sum_{i=1}^n \sum_{k=1}^9 \mathbb{1}_{y_i=k} \log P(y_i = k \mid x_i, \beta) - \frac{\text{11\_ratio}}{C} \|\beta\|_1 - \frac{1 - \text{11\_ratio}}{C} \|\beta\|_2^2 \quad (6)$$

where  $\lambda_1 = \text{11\_ratio}/C$  and  $\lambda_2 = (1 - \text{11\_ratio})/C$ . We use Grid Search to tune the hyperparameters  $C$  and `11_ratio` using  $k$ -fold crossvalidation with  $k = 3$ . The parameter space that we searched was  $C \in \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2, 10^3, 10^4\}$  and `11_ratio`  $\in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$ . We plot the 3-fold mean validation accuracy in Figure 3. Elastic Net had the highest validation accuracy of 91.97%, followed by LASSO and Ridge both with 91.95% accuracy. The optimal parameters for each model are provided in the caption of Figure 3.

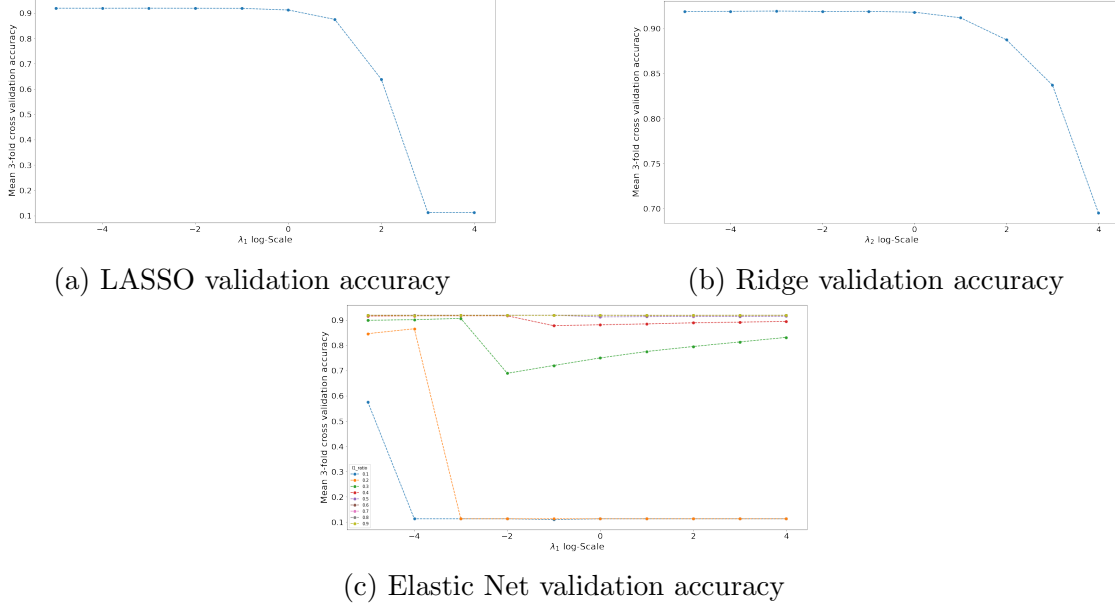


Figure 3: The parameters with the highest 3-fold mean validation accuracy was (a)  $\lambda_1 = 10^{-2}$  for LASSO, (b)  $\lambda_2 = 10^{-4}$  for Ridge, and (c)  $\lambda_1 = 10^{-2}$ , `11_ratio` = 0.8 for Elastic Net.

Since the  $\ell_1$ -norm induces sparsity, LASSO had the highest percent of features with coefficient 0 (18.62%) followed by Elastic Net with 15.22% and Ridge with 8.55% as seen in Figure 4. LASSO has advantages over Ridge regression since we can reduce the overdetermined system  $A^T X = B^T$  to have fewer columns with LASSO where we can remove the pixels that have a corresponding coefficient of zero.

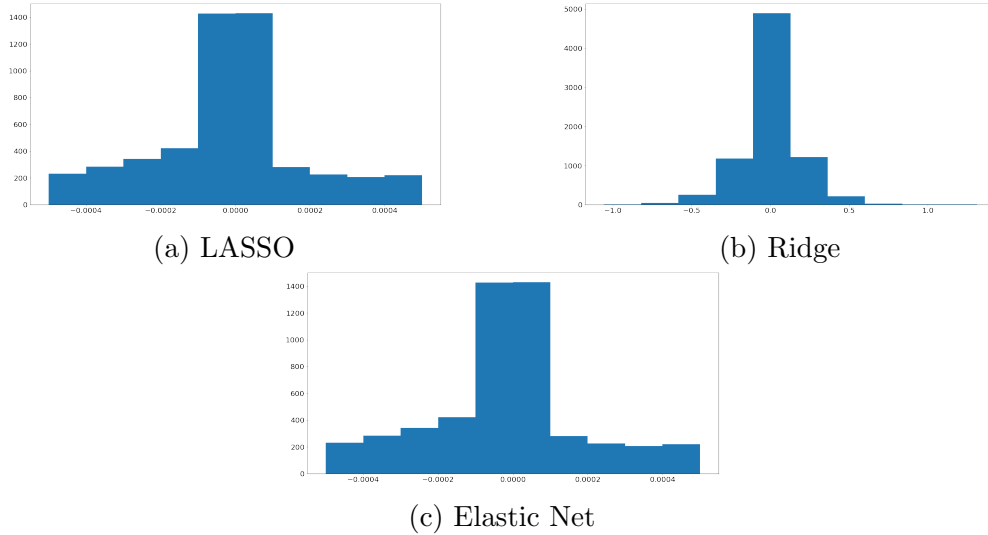


Figure 4: (a) 18.62% of features are 0 for LASSO, (b) 8.55% of features are 0 for Ridge, and (c) 15.22% of features are 0 for Elastic Net

We visualize the coefficient values in Figures 5, 6, and 7. We can see the where the digits are written in blue and the outline of the digit in red. As expected, the white pixels near the boarder to not contribute information towards categorizing the image.

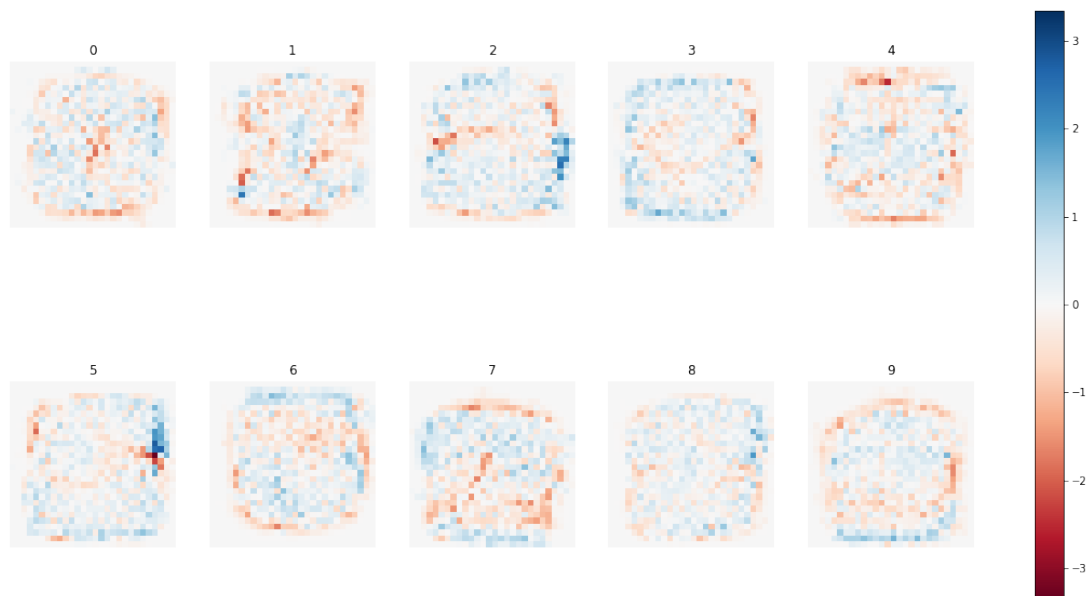


Figure 5: Coefficients for LASSO model

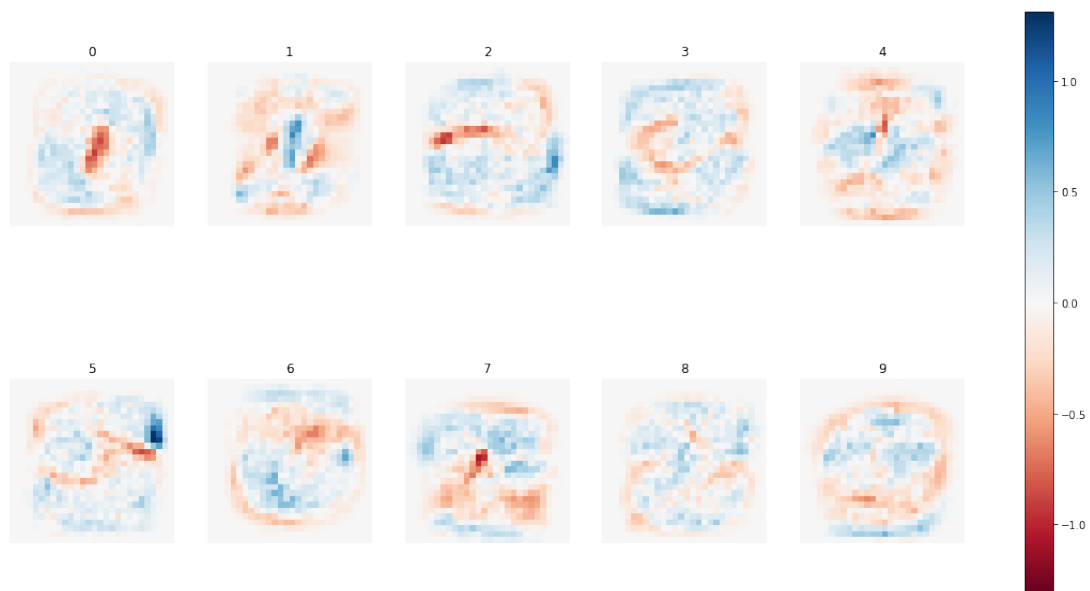


Figure 6: Coefficients for Ridge model

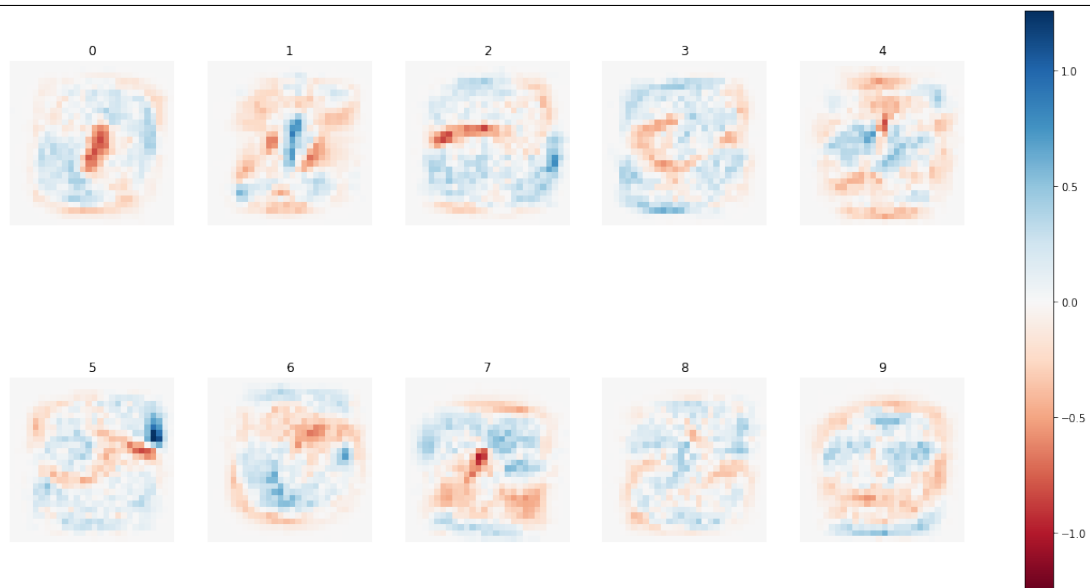


Figure 7: Coefficients for Elastic Net model

### 3 Ranking Pixels

For each model, we use the absolute value of the coefficient  $|\beta_{ij}|$  to rank each pixel  $j$  for digit  $i$ . In figures 8, 9, and 10 we plot the top 100 coefficients for each model.

To rank each pixel across all digits, we take the sum

$$\sum_{i=0}^9 |\beta_{ij}| \quad (7)$$

to rank pixel  $j$  for  $1 \leq j \leq 784$ . Again we plot the top 100 pixels for each method in Figures 11. Ridge and Elastic Net seem to rank pixels on the inside border higher, while LASSO ranks pixels on the outside boarder higher.

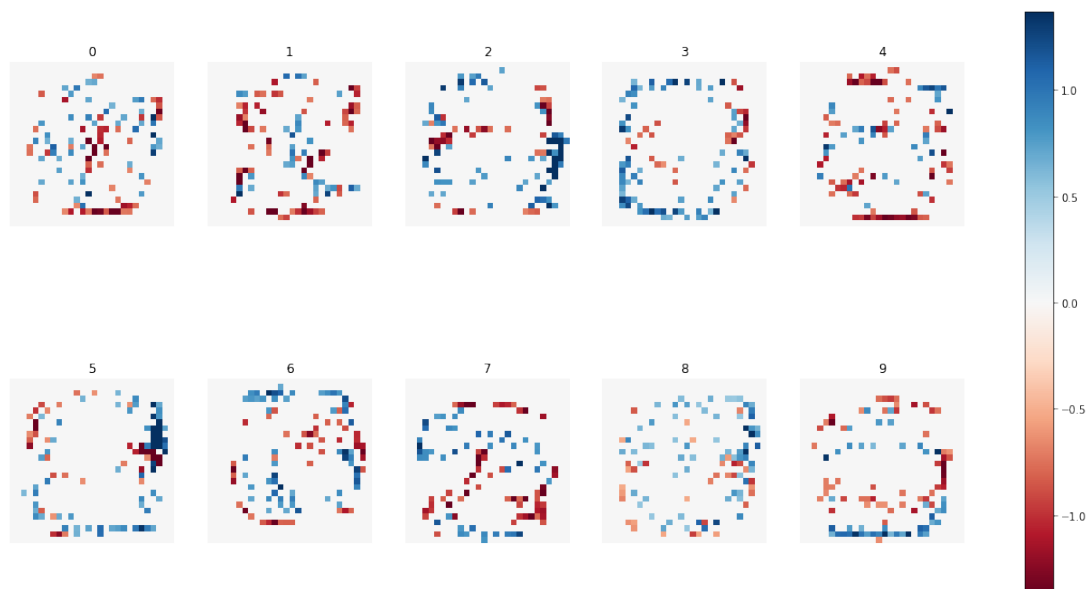


Figure 8: Each digit top 100 coefficients for LASSO model

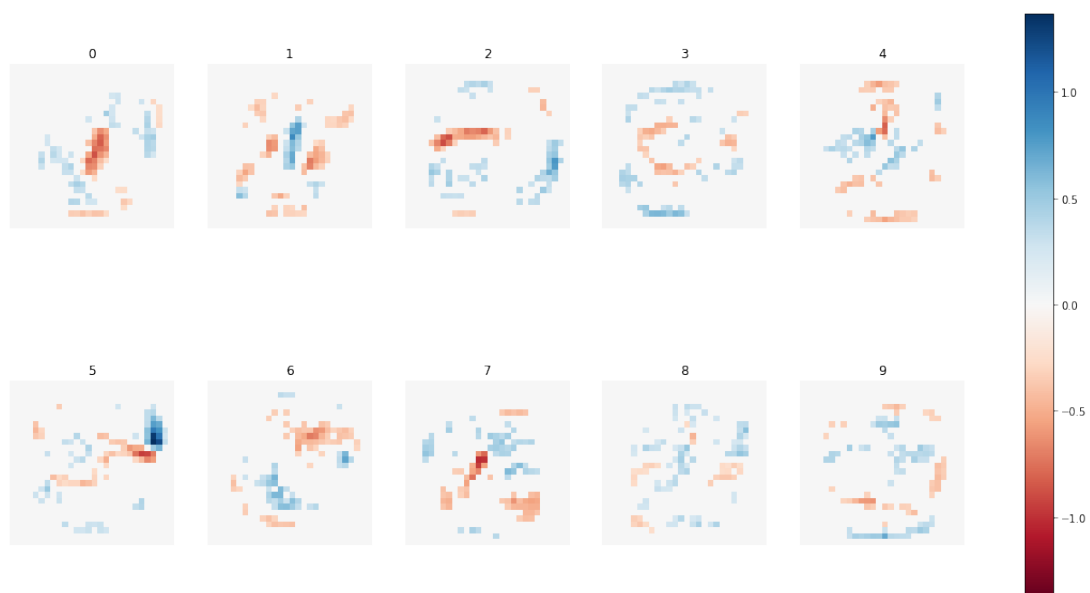


Figure 9: Each digit top 100 coefficients for Ridge model

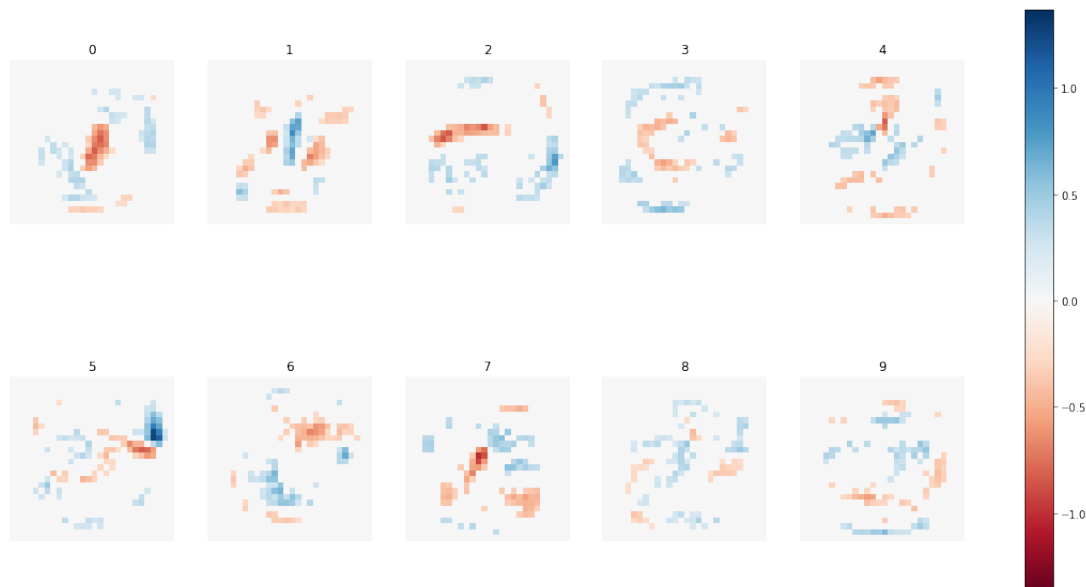
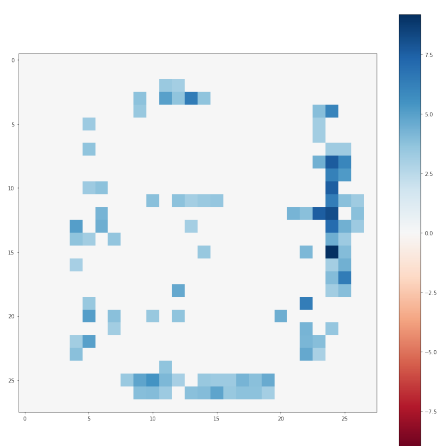
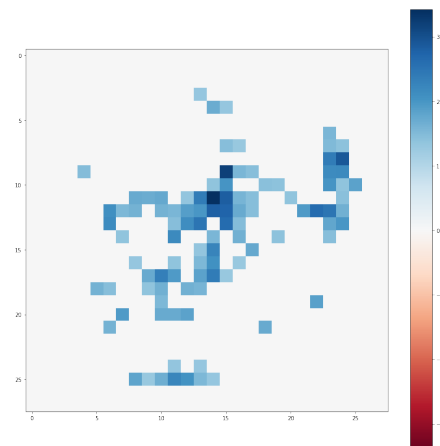


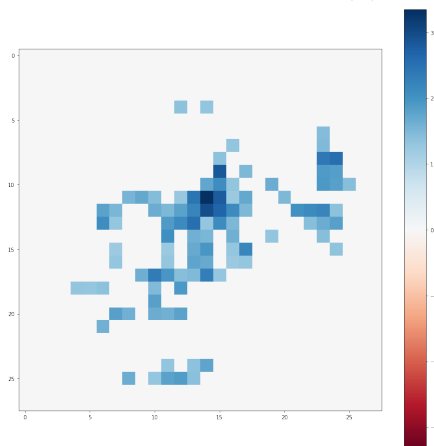
Figure 10: Each digit top 100 coefficients for Elastic Net model



(a) LASSO Top 100 pixels



(b) Ridge Top 100 pixels



(c) Elastic Net Top 100 pixels

Figure 11: Top 100 pixels for each method using equation (7) to rank



## 4 Test Set Evaluation

First, we evaluate the full models on the test set of 10,000 images. LASSO had the highest test accuracy with 92.62%, followed by Ridge with 92.54% and Elastic Net with 92.52%. We also compute the precision, recall, and  $F_1$  Scores for each digit where

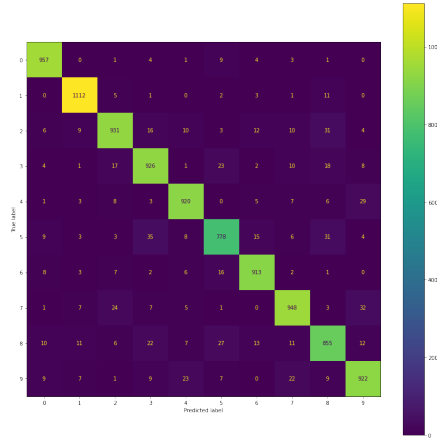
$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

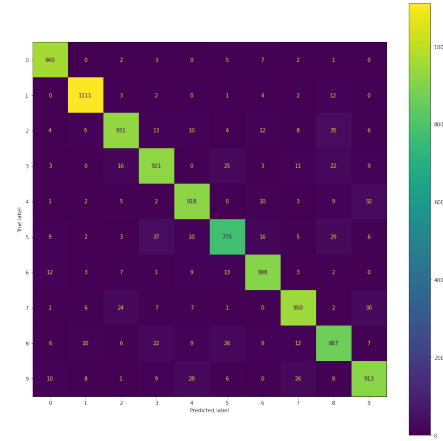
and

$$F_1 = 2 \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}.$$

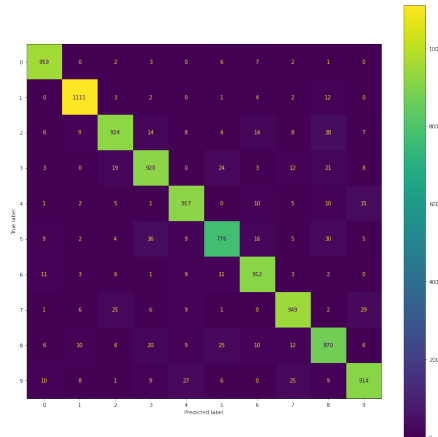
TP stands for true positive, FP stands for false positive, and FN stands for false negative. LASSO had a higher  $F_1$  score for 4s, 6s, and 9s than Ridge and Elastic Net.



(a) LASSO Confusion Matrix



(b) Ridge Confusion Matrix



(c) Elastic Net Confusion Matrix

Figure 12: Evaluated on full test set

	precision	recall	f1-score		precision	recall	f1-score		precision	recall	f1-score
0	0.952	0.977	0.964	0	0.954	0.980	0.967	0	0.953	0.979	0.966
1	0.962	0.980	0.971	1	0.965	0.979	0.972	1	0.965	0.979	0.972
2	0.928	0.902	0.915	2	0.933	0.902	0.917	2	0.929	0.895	0.912
3	0.903	0.917	0.910	3	0.906	0.912	0.909	3	0.909	0.911	0.910
4	0.938	0.937	0.937	4	0.926	0.935	0.931	4	0.928	0.934	0.931
5	0.898	0.872	0.885	5	0.905	0.869	0.887	5	0.909	0.870	0.889
6	0.944	0.953	0.949	6	0.937	0.948	0.942	6	0.934	0.952	0.943
7	0.929	0.922	0.926	7	0.930	0.924	0.927	7	0.928	0.923	0.925
8	0.885	0.878	0.881	8	0.878	0.890	0.884	8	0.874	0.893	0.884
9	0.912	0.914	0.913	9	0.910	0.905	0.908	9	0.914	0.906	0.910

Table 1: LASSO

Table 2: Ridge

Table 3: Elastic Net

Table 4: Classification metrics for full models

#### 4.1 Top 100 pixel Test Set Evaluation

We analyze how our models perform on the top 100 pixels as show in Figure 11. We subset both the training and test sets to the top 100 models and train three models using LASSO, Ridge, and Elastic Net regularization using the parameters as stated in the caption of Figure 3. We evaluated our models on the sparse test sets and compute the confusion matrix, accuracy, precision, recall, and  $F_1$  scores. In this setting, Elastic Net has the highest accuracy of 87.5% followed by Ridge with 88.3% and LASSO with 79.7%. Ridge had a higher  $F_1$  score than LASSO and Elastic Net for all digits.

Subsetting our images to these top 100 pixels decreases the total number of pixels 87.2%, which when scaled will decrease storage requirements and runtime. However, with this reduction we saw a reduction of approximately 5% in our test accuracy for Ridge Regression.

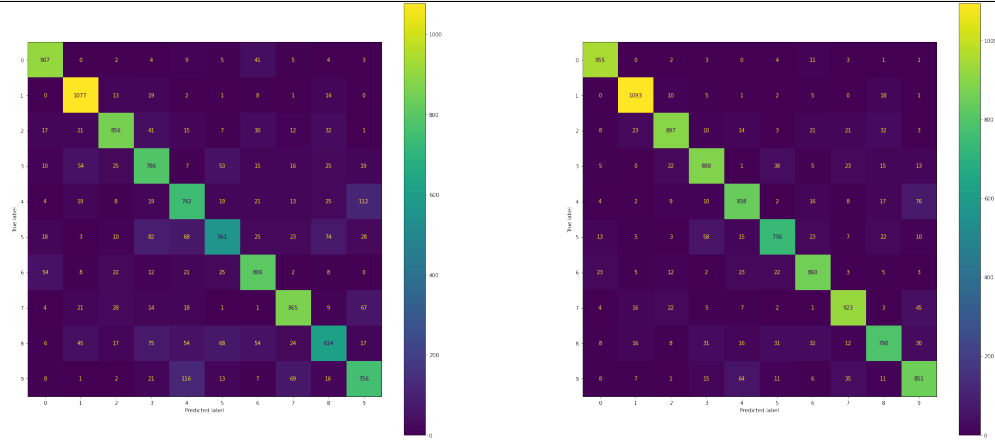
	precision	recall	f1-score		precision	recall	f1-score		precision	recall	f1-score
0	0.882	0.926	0.903	0	0.929	0.974	0.951	0	0.932	0.971	0.951
1	0.862	0.949	0.904	1	0.937	0.963	0.950	1	0.927	0.969	0.947
2	0.871	0.829	0.850	2	0.910	0.869	0.889	2	0.896	0.840	0.867
3	0.733	0.778	0.755	3	0.865	0.879	0.872	3	0.875	0.857	0.866
4	0.705	0.756	0.730	4	0.856	0.853	0.855	4	0.836	0.838	0.837
5	0.745	0.629	0.682	5	0.865	0.825	0.845	5	0.853	0.835	0.844
6	0.800	0.841	0.820	6	0.878	0.898	0.888	6	0.885	0.887	0.886
7	0.840	0.841	0.841	7	0.892	0.898	0.895	7	0.880	0.894	0.887
8	0.748	0.630	0.684	8	0.864	0.811	0.837	8	0.849	0.813	0.831
9	0.754	0.749	0.751	9	0.824	0.843	0.833	9	0.802	0.825	0.813

Table 5: LASSO

Table 6: Ridge

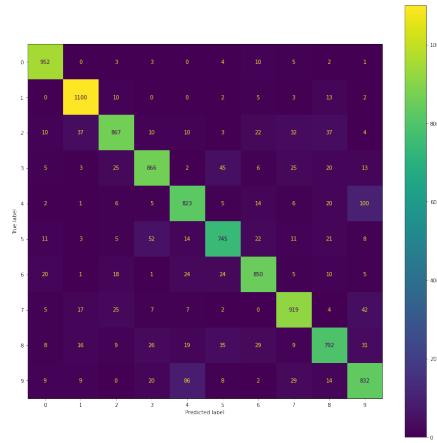
Table 7: Elastic Net

Table 8: Classification metrics for sparse top 100 pixel models



(a) LASSO Confusion Matrix

(b) Ridge Confusion Matrix



(c) Elastic Net Confusion Matrix

Figure 13: Evaluated on sparse top 100 pixel test set

## 5 Appendix Code

All analysis was completed using Python 3 in a Google Collaboratory Jupyter Notebook. Libraries used include `google.colab` for file management; `numpy` for linear algebra; `matplotlib` for plotting; `tensorflow` for loading the mnist dataset; and `scikitlearn` for machine learning. Code for this analysis is provided below and on [GitHub](#).

```

1 # -*- coding: utf-8 -*-
2 """ hwunrow_amath584_hw6.ipynb
3
4 Automatically generated by Colaboratory.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1okPbSGqe5fgeb_keR9biC92qE0xrbCd3r
8 """
9
10 # Commented out IPython magic to ensure Python compatibility.
11 import numpy as np
12 import pandas as pd
13 import tensorflow as tf
14
15 import matplotlib.pyplot as plt
16 # %matplotlib inline
17
18 from tensorflow.keras.datasets.mnist import load_data
19 from tensorflow.keras.utils import to_categorical
20
21 from sklearn.linear_model import LogisticRegression
22 from sklearn.model_selection import GridSearchCV
23 from sklearn.metrics import confusion_matrix
24 from sklearn.metrics import plot_confusion_matrix
25 from sklearn.metrics import accuracy_score
26 from sklearn.metrics import classification_report
27
28 from google.colab import drive
29 drive.mount('/content/drive')
30
31 """# Data Processing and EDA"""
32
33 data_dir = "/content/drive/MyDrive/School/AMATH584/repos/amath584/hw6/data/mnist.npz"
34 (x_train, y_train), (x_test, y_test) = load_data(path=data_dir)
35
36 print(f"training set dimensions: {x_train.shape}")
37 print(f"test set dimensions: {x_test.shape}")
38
39 x_train = x_train.reshape(60000, 784)
40 x_test = x_test.reshape(10000, 784)
41
42 # one-hot encode
43 y_train_encode = to_categorical(y_train, 10)
44 y_test_encode = to_categorical(y_test, 10)
45
46 """Plot the first 10 numbers in training set"""
47
48 fig=plt.figure(figsize=(20,10))
49
50 rows = 2
51 columns = 5
52
53 for i in range(10):
54     ax = fig.add_subplot(rows, columns, i+1)
55     plt.imshow(x_train[i].reshape(28,28))
56     plt.axis('off')
57

```

```

58 x_train.max()
59
60 # normalize
61 x_train = x_train/255
62 x_test = x_test/255
63
64 np.histogram(y_test, bins=[0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
65
66 fig=plt.figure(figsize=(20,10))
67
68 plt.hist(y_train, bins=[0,1,2,3,4,5,6,7,8,9, 10], ec='black')
69
70 fig=plt.figure(figsize=(20,10))
71
72 plt.hist(y_test, bins=[0,1,2,3,4,5,6,7,8,9,10], ec='black')
73
74 """# Logistic Regression with LASSO, Ridge, and Elastic Net
75
76 ## Hyperparameter Tuning
77 We use Grid Search to tune the type of regularization (L1 norm or Elastic Net) and
78 size of the penalty.
79 """
80 lr = LogisticRegression(tol=0.1)
81
82 param_grid = [
83     {'C': [10**x for x in range(-5,5)], 'penalty': ['l1'], 'solver': ['saga']},
84     {'C': [10**x for x in range(-5,5)], 'l1_ratio': [0.1*(x+1) for x in range(9)], '
85         penalty': ['elasticnet'], 'solver': ['saga']},
86 ]
87 clf = GridSearchCV(lr, param_grid, cv=3, scoring='accuracy', n_jobs=-1)
88
89 clf.fit(x_train, y_train)
90
91 # import pickle
92 # filename = "/content/drive/MyDrive/School/AMATH584/repos/amath584/hw6/lr_norm_model
93     .sav"
94 # pickle.dump(clf, open(filename, 'wb'))
95
96 # clf = pickle.load(open(filename, 'rb'))
97
98 """### LASSO"""
99 lasso_accuracy = clf.cv_results_['mean_test_score'][0:10]
100
101 lasso_accuracy
102
103 plt.figure(figsize=(20,10))
104 #flip because the hyperparameter in scikit learn C = 1/lambda
105 plt.plot(list(range(-5,5)), np.flip(lasso_accuracy), '—o')
106 plt.xlabel(r"$\lambda^{-1}$ log-Scale", fontsize=22)
107 plt.ylabel("Mean 3-fold cross validation accuracy", fontsize=22)
108
109 plt.xticks(fontsize=20)
110 plt.yticks(fontsize=20)
111

```

```

112 lasso_best_params = clf.cv_results_['params'][np.argmax(lasso_accuracy)]
113 clf_lasso = LogisticRegression(**lasso_best_params)
114
115 clf_lasso.fit(x_train, y_train)
116
117 lasso_best_params
118
119 lasso_coefs = clf_lasso.coef_
120
121 print(f"{np.mean(lasso_coefs == 0) * 100:.2f} % of features are zero!")
122
123 plt.figure(figsize=(20,10))
124 plt.hist(lasso_coefs.flatten())
125
126 plt.xticks(fontsize=20)
127 plt.yticks(fontsize=20)
128
129 fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(20,10))
130
131 scale = np.abs(lasso_coefs).max()
132
133 for i, ax in enumerate(axes.flat):
134     ax.set_axis_off()
135     ax.set_title(f"{i}")
136     img = lasso_coefs[i]
137     img = img.reshape(28,28)
138     im = ax.imshow(img, interpolation='nearest', vmin=-scale, vmax=scale, cmap=plt.cm.
        RdBu)
139
140 cbar = fig.colorbar(im, ax=axes.ravel().tolist())
141
142 """### Elastic Net"""
143
144 elastic_net_accuracy = clf.cv_results_['mean_test_score'][10:101]
145 elastic_net_params = clf.cv_results_['params'][10:101]
146
147 # rows represent same C, columns are for same l1_ratio
148 elastic_net_accuracy = elastic_net_accuracy.reshape(9,10)
149 #flip because the hyperparameter in scikit learn C = 1/lambda
150 elastic_net_accuracy = np.flip(elastic_net_accuracy, axis=1)
151
152 plt.figure(figsize=(20,10))
153 #flip because the hyperparameter in scikit learn C = 1/lambda
154 plt.plot(list(range(-5,5)), elastic_net_accuracy.T, '—o')
155 plt.xlabel(r"$\lambda_1$ log-Scale", fontsize=22)
156 plt.ylabel("Mean 3-fold cross validation accuracy", fontsize=22)
157
158 plt.legend([round(0.1*(x+1),1) for x in range(9)], title="l1_ratio")
159
160 plt.xticks(fontsize=20)
161 plt.yticks(fontsize=20)
162
163 clf.best_params_
164
165 coefs = clf.best_estimator_.coef_
166
167 print(f"{np.mean(coefs == 0) * 100:.2f} % of features are zero!")

```

```

168
169 plt.figure(figsize=(20,10))
170 plt.hist(coefs.flatten())
171
172 plt.xticks(fontsize=20)
173 plt.yticks(fontsize=20)
174
175 fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(20,10))
176
177 scale = np.abs(coefs).max()
178
179 for i, ax in enumerate(axes.flat):
180     ax.set_axis_off()
181     ax.set_title(f"{i}")
182     img = coefs[i]
183     img = img.reshape(28,28)
184     im = ax.imshow(img, interpolation='nearest', vmin=-scale, vmax=scale, cmap=plt.cm.
        RdBu)
185
186 cbar = fig.colorbar(im, ax=axes.ravel().tolist())
187
188 """### Ridge
189 For comparison, we also fit a logistic regression model with  $\ell_2$ -norm
        regularization.
190 """
191
192 lr = LogisticRegression(tol=0.1)
193
194 param_grid = [
195     {'C': [10**x for x in range(-5,5)], 'penalty': ['l2'], 'solver': ['saga']},
196 ]
197
198 clf_ridge = GridSearchCV(lr, param_grid, cv=3, scoring='accuracy', n_jobs=2)
199
200 clf_ridge.fit(x_train, y_train)
201
202 clf_ridge.best_params_
203
204 print(f"LASSO: {lasso_accuracy.max() * 100}")
205 print(f"Ridge: {clf_ridge.best_score_ * 100}")
206 print(f"Elastic Net: {clf.best_score_ * 100}")
207
208 ridge_accuracy = clf_ridge.cv_results_['mean_test_score']
209
210 plt.figure(figsize=(20,10))
211 #flip because the hyperparameter in scikit learn C = 1/lambda
212 plt.plot(list(range(-5,5)), np.flip(ridge_accuracy), '—o')
213 plt.xlabel(r" $\lambda$  log-Scale", fontsize=22)
214 plt.ylabel("Mean 3-fold cross validation accuracy", fontsize=22)
215
216 plt.xticks(fontsize=20)
217 plt.yticks(fontsize=20)
218
219 ridge_coefs = clf_ridge.best_estimator_.coef_
220
221 print(f"{np.mean(ridge_coefs == 0) * 100:.2f} % of features are zero!")
222

```

```

223 plt.figure(figsize=(20,10))
224 plt.hist(ridge_coefs.flatten())
225
226 plt.xticks(fontsize=20)
227 plt.yticks(fontsize=20)
228
229 fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(20,10))
230
231 scale = np.abs(ridge_coefs).max()
232
233 for i, ax in enumerate(axes.flat):
234     ax.set_axis_off()
235     ax.set_title(f"{i}")
236     img = ridge_coefs[i]
237     img = img.reshape(28,28)
238     im = ax.imshow(img, interpolation='nearest', vmin=-scale, vmax=scale, cmap=plt.cm.
        RdBu)
239
240 cbar = fig.colorbar(im, ax=axes.ravel().tolist())
241
242 """# Ranking Pixels
243 We rank the pixels by each its corresponding coefficient for each digit.
244
245 ### LASSO
246 """
247
248 # most important pixels for each digit
249 lasso_coefs.argmax(axis=1)
250
251 n = 100
252
253 lasso_coefs_copy = lasso_coefs.copy()
254
255 arr = np.empty((0,784), float)
256 for row in lasso_coefs_copy:
257     abs_row = abs(row)
258     abs_row.sort()
259     thresh = abs_row[-n]
260     filter = abs(row) < thresh
261     row[filter] = 0
262     arr = np.append(arr, np.array([row]), axis=0)
263
264
265 fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(20,10))
266
267 scale = np.abs(ridge_coefs).max()
268
269 for i, ax in enumerate(axes.flat):
270     ax.set_axis_off()
271     ax.set_title(f"{i}")
272     img = arr[i]
273     img = img.reshape(28,28)
274     im = ax.imshow(img, interpolation='nearest', vmin=-scale, vmax=scale, cmap=plt.cm.
        RdBu)
275
276 cbar = fig.colorbar(im, ax=axes.ravel().tolist())
277

```



```

278 lasso_all_digit_rank = abs(lasso_coefs_copy).sum(axis=0)
279
280 lasso_all_digit_copy = lasso_all_digit_rank.copy()
281 lasso_all_digit_copy.sort()
282 thresh = lasso_all_digit_copy[-n]
283 filter = lasso_all_digit_rank < thresh
284 lasso_all_digit_rank[filter] = 0
285
286 scale = np.abs(lasso_all_digit_rank).max()
287
288 fig, ax = plt.subplots(figsize=(15, 15), ncols=1)
289
290 im = ax.imshow(lasso_all_digit_rank.reshape(28,28), interpolation='nearest',
291              vmin=-scale, vmax=scale, cmap=plt.cm.RdBu)
292
293 fig.colorbar(im, ax=ax)
294
295 """### Ridge"""
296
297 n = 100
298
299 ridge_coefs_copy = ridge_coefs.copy()
300
301 arr = np.empty((0,784), float)
302 for row in ridge_coefs_copy:
303     abs_row = abs(row)
304     abs_row.sort()
305     thresh = abs_row[-n]
306     filter = abs(row) < thresh
307     row[filter] = 0
308     arr = np.append(arr, np.array([row]), axis=0)
309
310
311 fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(20,10))
312
313 scale = np.abs(ridge_coefs).max()
314
315 for i, ax in enumerate(axes.flat):
316     ax.set_axis_off()
317     ax.set_title(f"{i}")
318     img = arr[i]
319     img = img.reshape(28,28)
320     im = ax.imshow(img, interpolation='nearest', vmin=-scale, vmax=scale, cmap=plt.cm.
321                    RdBu)
322
323 cbar = fig.colorbar(im, ax=axes.ravel().tolist())
324
325 ridge_all_digit_rank = abs(ridge_coefs_copy).sum(axis=0)
326
327 ridge_all_digit_copy = ridge_all_digit_rank.copy()
328 ridge_all_digit_copy.sort()
329 thresh = ridge_all_digit_copy[-n]
330 filter = ridge_all_digit_rank < thresh
331 ridge_all_digit_rank[filter] = 0
332
333 scale = np.abs(ridge_all_digit_rank).max()

```

```

334 fig , ax = plt.subplots(figsize=(15, 15), ncols=1)
335
336 im = ax.imshow(ridge_all_digit_rank.reshape(28,28), interpolation='nearest',
337               vmin=-scale, vmax=scale, cmap=plt.cm.RdBu)
338
339 fig.colorbar(im, ax=ax)
340
341 """## Elastic Net"""
342
343 n = 100
344
345 en_coefs_copy = clf.best_estimator_.coef_.copy()
346
347 arr = np.empty((0,784), float)
348 for row in en_coefs_copy:
349     abs_row = abs(row)
350     abs_row.sort()
351     thresh = abs_row[-n]
352     filter = abs(row) < thresh
353     row[filter] = 0
354     arr = np.append(arr, np.array([row]), axis=0)
355
356
357 fig , axes = plt.subplots(nrows=2, ncols=5, figsize=(20,10))
358
359 scale = np.abs(ridge_coefs).max()
360
361 for i, ax in enumerate(axes.flat):
362     ax.set_axis_off()
363     ax.set_title(f"{i}")
364     img = arr[i]
365     img = img.reshape(28,28)
366     im = ax.imshow(img, interpolation='nearest', vmin=-scale, vmax=scale, cmap=plt.cm.
367                   RdBu)
368
369 cbar = fig.colorbar(im, ax=axes.ravel().tolist())
370
371 en_all_digit_rank = abs(en_coefs_copy).sum(axis=0)
372
373 en_all_digit_copy = en_all_digit_rank.copy()
374 en_all_digit_copy.sort()
375 thresh = en_all_digit_copy[-n]
376 filter = en_all_digit_rank < thresh
377 en_all_digit_rank[filter] = 0
378
379 scale = np.abs(en_all_digit_rank).max()
380
381 fig , ax = plt.subplots(figsize=(15, 15), ncols=1)
382
383 im = ax.imshow(en_all_digit_rank.reshape(28,28), interpolation='nearest',
384               vmin=-scale, vmax=scale, cmap=plt.cm.RdBu)
385
386 fig.colorbar(im, ax=ax)
387
388 """# Test set
389 ### LASSO

```

```

390 """
391
392 yhat = clf_lasso.predict(x_test)
393
394 print(f"Test accuracy {accuracy_score(y_test, yhat)*100}%")
395 confusion_matrix(y_test, yhat)
396
397 print(classification_report(y_test, yhat, target_names=['0', '1', '2', '3', '4', '5',
398               '6', '7', '8', '9'], digits=3))
399
400 fig, ax = plt.subplots(figsize=(15, 15))
401
402 plot_confusion_matrix(clf_lasso, x_test, y_test, ax=ax, values_format='.4g')
403
404 """### Ridge"""
405
406 confusion_matrix(y_test, clf_ridge.predict(x_test))
407
408 print(f"Test accuracy {accuracy_score(y_test, clf_ridge.predict(x_test))*100}%")
409
410 print(classification_report(y_test, clf_ridge.predict(x_test), target_names=['0', '1',
411               '2', '3', '4', '5', '6', '7', '8', '9'], digits=3))
412
413 fig, ax = plt.subplots(figsize=(15, 15))
414
415 plot_confusion_matrix(clf_ridge, x_test, y_test, ax=ax, values_format = '.4g')
416
417 """### Elastic Net"""
418
419 confusion_matrix(y_test, clf.best_estimator_.predict(x_test))
420
421 print(f"Test accuracy {accuracy_score(y_test, clf.best_estimator_.predict(x_test))
422       *100}%")
423
424 print(classification_report(y_test, clf.best_estimator_.predict(x_test), target_names
425       =['0', '1', '2', '3', '4', '5', '6', '7', '8', '9'], digits=3))
426
427 fig, ax = plt.subplots(figsize=(15, 15))
428
429 plot_confusion_matrix(clf.best_estimator_, x_test, y_test, ax=ax, values_format = '.4
430       g')
431
432 """## Sparse Test Evaluation
433 Using the top 100 pixels from each model we evaluate on the test set.
434
435 ### LASSO
436 """
437
438 lasso_sparse_subset = lasso_all_digit_rank.nonzero()
439
440 lasso_sparse_subset
441
442 # subset training and test sets to top 100 pixels
443 x_train_sparse_lasso = np.transpose(np.transpose(x_train)[lasso_sparse_subset])
444 x_test_sparse_lasso = np.transpose(np.transpose(x_test)[lasso_sparse_subset])
445

```

```

442 clf_lasso_sparse = LogisticRegression(C= 100, penalty='l1', solver='saga', tol=0.1)
443 clf_lasso_sparse.fit(x_train_sparse_lasso, y_train)
444
445 yhat = clf_lasso_sparse.predict(x_test_sparse_lasso)
446 print(f"Test accuracy {accuracy_score(y_test, yhat)*100}%")
447 print(classification_report(
448     y_test, yhat,
449     target_names=['0', '1', '2', '3', '4', '5', '6', '7', '8', '9'],
450     digits=3)
451 )
452
453 fig, ax = plt.subplots(figsize=(15, 15))
454
455 plot_confusion_matrix(clf_lasso_sparse, x_test_sparse_lasso, y_test, ax=ax,
456     values_format = '.4g')
457
458 """### Ridge"""
459
460 ridge_sparse_subset = ridge_all_digit_rank.nonzero()
461
462 # subset training and test sets to top 100 pixels
463 x_train_sparse_ridge = np.transpose(np.transpose(x_train)[ridge_sparse_subset])
464 x_test_sparse_ridge = np.transpose(np.transpose(x_test)[ridge_sparse_subset])
465
466 clf_ridge_sparse = LogisticRegression(C= 10000, penalty='l2', solver='saga', tol=0.1)
467 clf_ridge_sparse.fit(x_train_sparse_ridge, y_train)
468
469 yhat = clf_ridge_sparse.predict(x_test_sparse_ridge)
470 print(f"Test accuracy {accuracy_score(y_test, yhat)*100}%")
471 print(classification_report(
472     y_test, yhat,
473     target_names=['0', '1', '2', '3', '4', '5', '6', '7', '8', '9'],
474     digits=3)
475 )
476
477 fig, ax = plt.subplots(figsize=(15, 15))
478
479 plot_confusion_matrix(clf_ridge_sparse, x_test_sparse_ridge, y_test, ax=ax,
480     values_format = '.4g')
481
482 """### Elastic Net"""
483
484 en_sparse_subset = en_all_digit_rank.nonzero()
485
486 # subset training and test sets to top 100 pixels
487 x_train_sparse_en = np.transpose(np.transpose(x_train)[en_sparse_subset])
488 x_test_sparse_en = np.transpose(np.transpose(x_test)[en_sparse_subset])
489
490 clf.best_params_
491
492 clf_en_sparse = LogisticRegression(C= 100, penalty='elasticnet', l1_ratio=0.8, solver
493     ='saga', tol=0.1)
494 clf_en_sparse.fit(x_train_sparse_en, y_train)
495
496 yhat = clf_en_sparse.predict(x_test_sparse_en)
497 print(f"Test accuracy {accuracy_score(y_test, yhat)*100}%")
498 print(classification_report(

```

```
496     y_test , yhat ,  
497     target_names=['0' , '1' , '2' , '3' , '4' , '5' , '6' , '7' , '8' , '9'] ,  
498     digits=3)  
499 )  
500  
501 fig , ax = plt.subplots(figsize=(15, 15))  
502  
503 plot_confusion_matrix(clf_en_sparse , x_test_sparse_en , y_test , ax=ax, values_format =  
    '.4g')
```