# Final Exam, CPSC 8420, Spring 2022

### Sherrer, Will

### Due 05/06/2022, Friday, 11:59PM EST

## Problem 1 [15 pts]

Consider the elastic-net optimization problem:

$$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda[\alpha\|\beta\|_2^2 + (1-\alpha)\|\beta\|_1]. \tag{1}$$

1. Show the objective can be reformulated into a lasso problem, with a slightly different $\mathbf{X}, \mathbf{y}$.

   We can see that if alpha is 0, the ridge regression aspect of elastic net will go to 0 and we will have the Lasso formulation remaining as

   $$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|_1 \tag{2}$$

2. If we fix $\alpha = .5$, please derive the closed solution by making use of alternating minimization that each time we fix the rest by optimizing one single element in $\beta$. You need randomly generate $\mathbf{X}, \mathbf{y}$ and initialize $\beta_0$, and show the objective decreases monotonically with updates.
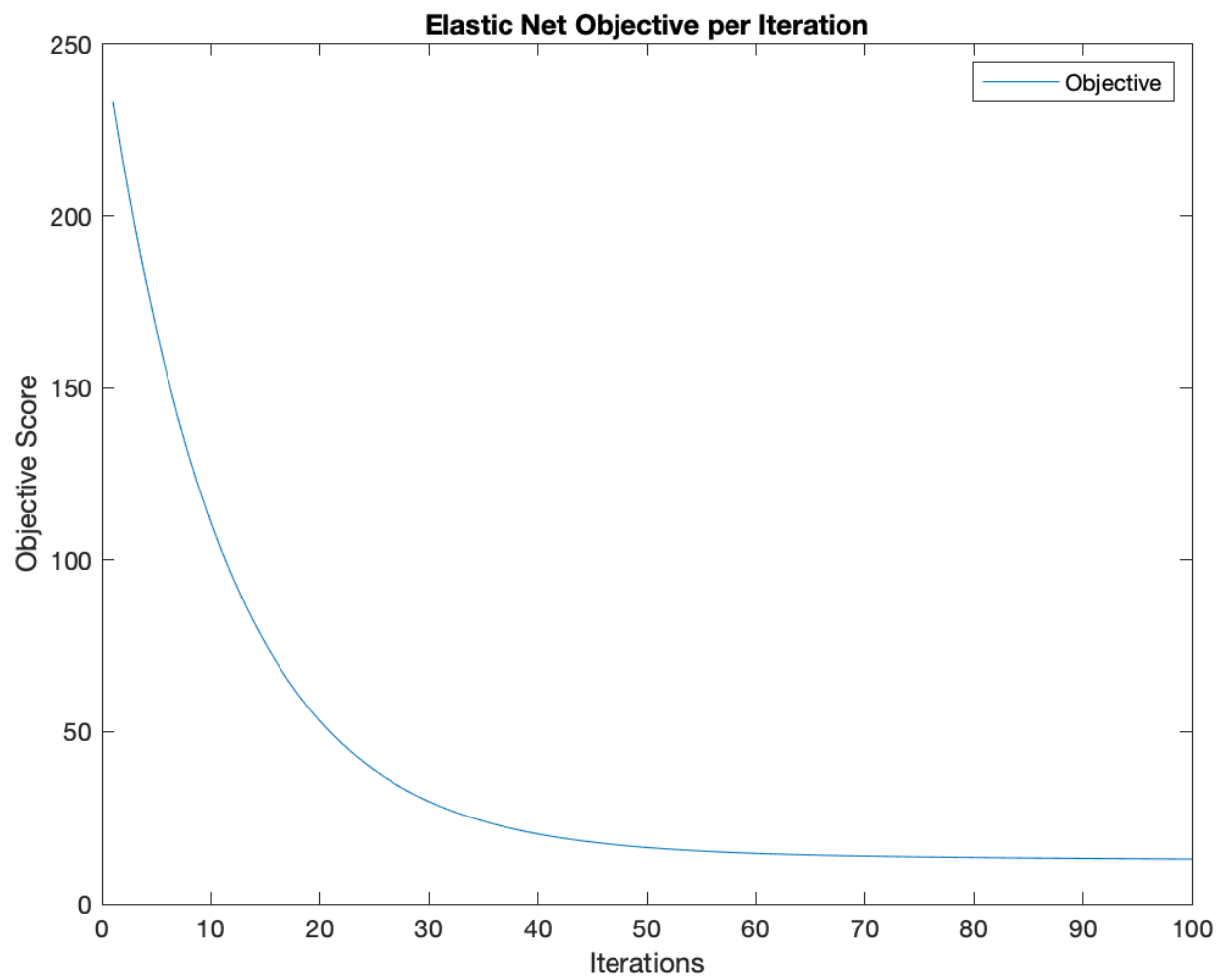
   First we can use reformulate the objective with $\alpha = 0.5$ to be:

   $$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \frac{\lambda}{2}\|\beta\|_2^2 + \frac{\lambda}{2}\|\beta\|_1. \tag{3}$$

   We can take the derivative w.r.t $\beta$, however since we are also using L1 norm, we must divide into positive and negative cases where: $\hat{\beta} = -2 * X^T(y - X\beta) + \frac{\lambda}{2}\beta$

   $$\hat{\beta}^{Lasso} = \begin{cases} \hat{\beta} - \frac{\lambda}{2} & if f\hat{\beta} > \lambda \\ \hat{\beta} + \frac{\lambda}{2} & if f\hat{\beta} \le \lambda \\ 0 & else \end{cases} \tag{4}$$

   For my implementation I am also randomly initializing $\beta$ and calculating for 100 iterations with a learning rate of 0.0001 and where $\lambda = 1$. Since I am not plotting the objective with different $\lambda$ since I only need to show that the objective is decreasing.

**Elastic Net Objective per Iteration**



From this graph I can see that the objective of Elastic Net Objective is decreasing per iteration.

# Problem 2 [15 pts]

- For PCA, the loading vectors can be directly computed from the $q$ columns of $\mathbf{U}$ where $[\mathbf{U}, \mathbf{S}, \mathbf{U}] = svd(\mathbf{X}^T\mathbf{X})$, please show that any $[\pm\mathbf{u}_1, \pm\mathbf{u}_2, \ldots, \pm\mathbf{u}_q]$ will be equivalent to $[\mathbf{u}_1, \mathbf{u}_2, \ldots, \mathbf{u}_q]$ in terms of the same variance while satisfying the orthonormality constraint.

  We know that the columns of U are orthogonal as that is in the definition of SVD. Orthogonal vectors are calculated by taking the dot product which takes the cosine of the angle between the vectors. Cosine has a symmetrical property and is an even function meaning that $cos(\theta) = cos(-\theta)$. Since the magnitude is not dependent of the signs of the vector, and the angle between the vectors is the same whether negative or positve, the dot product of both positive and negative columns will be the same which is 0. This means negative columns of U satisfies the orthonormality constraint.

  We can assume that $\phi = \sum_{i=1}^{p} \lambda_i U(:, i)$ and the constraint $||\phi||_2 = 1$ is equal to $\sum_{i=1}^{p} \lambda_i^2 = 1$. We can rewrite $X^T X = \sum_{i=1}^{p} S(i, i)U(:, i)U(:, i)^T$ and since we know the columns of U are orthogonal whether negative or positive, therefore, $||X\phi||_2^2 = \sum_{i=1}^{p} S(i, i)\lambda_i^2 \leq S(1, 1)$ Therefore the variance is the same in both cases.

- We consider the case when original dimensionality of the data is much larger than the number of samples $d \gg m$ ($\mathbf{X} \in \mathbb{R}^{d \times m}$). What's the complexity of obtaining the optimal solution of PCA via Singular Value Decomposition? Please consider a more efficient solution by considering the relationships of eigenvalues/eigenvectors between $\mathbf{X}^T\mathbf{X}$ and $\mathbf{X}\mathbf{X}^T$.

  I found that the time complexity of SVD is $O(n^3)$. In a d x m matrix when d$\gg$m, we want to use $X^T X$ so that we are performing SVD on an m x m matrix.

  For matrix $X$ of size $mxn$, We can show that both $X^T X and X X^T$ are symmetric square matrices of size (n, n) and (m, m), respectively.

  $(X^T X)^T = X^T(X^T)^T = X^T X$
  $(X X^T)^T = (X^T)^T X^T = X X^T$

  It can also be shown that $X X^T$ and $X^T X$ share the same non-zero eigenvalues. If one has more eigenvalues than the other, all of the excess eigenvalues should be zero.
  If we let $v = $ the eigenvector of $X^T X$ and $\lambda$ be its eigenvalue:
  $X^T X b = \lambda v$ – if we multiply both sides by $X$ then:
  $X X^T (Xv) = \lambda(Xv)$.
  Since $\lambda$ does not $= 0$, it follows $X^T v$ does not $= 0$, therefore $X^T v$ is an eigenvector of $X^T X$ with eigenvalue of $\lambda$
  The eigenvectors to eigenvalue zero are the elements in the null space of $X : X^T v = 0$ implies $v^T X^T X v = ||Xv||^2 = 0$, which is $Xv=0$. You obtain the eigenvectors of $X^T X$ from $X^T v$ with $v$ eigenvector of $v$ of $X X^T$ to non-zero eigenvalue. Therefore you can obtain $V$ , whose each column is an eigenvector for $X X^T$, and D, a diagonal matrix holding $X X^T$'s eigenvalues.

## Problem 3 [10 pts]

Assume that in a community, there are 10% people suffer from COVID. Assume 80% of the patients come to breathing difficulty while 25% of those free from COVID also have symptoms of shortness of breath. Now please determine that if one has breathing difficulty, what's his/her probability to get COVID? (*hint*: you may consider Naive Bayes)

I want to calculate the probability of getting covid given that they have breathing difficulty. $P(C|B)$.

I say that the probability of getting covid, $P(C)$ is 0.1. The probability of having difficulty breathing given that they have covid is $P(B|C) = 0.8$.

To calculate the probability that one has breathing difficulty, I need to add the probability of breathing difficulty given covid and without covid.

$P(B) = P(B|C) * P(C) + P(B|\bar{C}) * P(\bar{C})$
$P(B) = 0.8 * 0.1 + 0.25 * 0.9$
$P(B) = .08 + .225 = .305$

$P(C|B) = \frac{P(B|C)*P(C)}{P(B)} = \frac{0.8*0.1}{0.305} = \frac{0.08}{0.305} = 0.2623$ or $26.23\%$

## Problem 4 [20 pts]

Recall the objective for RatioCut: $RatioCut(A_1, A_2, ...A_k) = \frac{1}{2} \sum\limits_{i=1}^{k} \frac{W(A_i, \overline{A}_i)}{|A_i|}$. If we introduce indicator vector: $h_j \in \{h_1, h_2, ..h_k\}, j \in [1, k]$, for any vector $h_j \in R^n$, we define: $h_{ij} = \begin{cases} 0 & v_i \notin A_j \\ \frac{1}{\sqrt{|A_j|}} & v_i \in A_j \end{cases}$,

we can prove: $h_i^T L h_i = \frac{cut(A_i, \overline{A}_i)}{|A_i|}$, and therefore:

$$RatioCut(A_1, A_2, ...A_k) = \sum_{i=1}^{k} h_i^T L h_i = \sum_{i=1}^{k} (H^T L H)_{ii} = tr(H^T L H), \tag{5}$$

thus we relax it as an optimization problem:

$$\underbrace{arg\ min}_{H}\ tr(H^T L H)\ \ s.t.\ H^T H = I. \tag{6}$$

Now let's explore Ncut, with objective: $NCut(A_1, A_2, ...A_k) = \frac{1}{2} \sum\limits_{i=1}^{k} \frac{W(A_i, \overline{A}_i)}{vol(A_i)}$, where $vol(A) :=$

$\sum\limits_{i \in A} d_i, d_i := \sum\limits_{j=1}^{n} w_{ij}$.Similar to Ratiocut, we define: $h_{ij} = \begin{cases} 0 & v_i \notin A_j \\ \frac{1}{\sqrt{vol(A_j)}} & v_i \in A_j \end{cases}$. Now

1. Please show that $h_i^T L h_i = \frac{cut(A_i, \overline{A}_i)}{vol(A_i)}$.

$h_i^T L h_i = \frac{1}{2} \sum\limits_{m=1}^{} \sum\limits_{n=1}^{} w_{mn}(h_{im} - h_{in})^2$

$= \frac{1}{2} ( \sum\limits_{m \in A_i, n \notin A_i} w_{mn}(\frac{1}{\sqrt{vol(A_i)}} - 0)^2 + \sum\limits_{m \notin A_i, n \in A_i} w_{mn}(0 - \frac{1}{\sqrt{vol(A_i)}})^2)$

$= \frac{1}{2} ( \sum\limits_{m \in A_i, n \notin A_i} w_{mn}\frac{1}{vol(A_i)} + \sum\limits_{m \notin A_i, n \in A_i} w_{mn}\frac{1}{vol(A_i)})$

$= \frac{1}{2}(cut(A_i, \overline{A}_i)\frac{1}{vol(A_i)} + cut(\overline{A}_i, A_i)\frac{1}{vol(A_i)})$

$= \frac{cut(A_i, \overline{A}_i)}{vol(A_i)}$

2. Show that $NCut(A_1, A_2, ...A_k) = tr(H^T L H)$.
   Since we have proved $h_i^T L h_i = \frac{cut(A_i, \overline{A}_i)}{vol(A_i)}$, therefore

$$NCut(A_1, A_2, ...A_k) = \sum_{i=1}^{k} h_i^T L h_i = \sum_{i=1}^{k} (H^T L H)_{ii} = tr(H^T L H), \tag{7}$$

3. The constraint now is: $H^T D H = I$.

4. Find the solution to $\underbrace{arg\ min}_{H}\ tr(H^T L H)\ \ s.t.\ H^T D H = I$.
   In Spectural Clustering we need to construct the Laplacian for the obtained graph and calculate the first k eigenvectors of the Laplacian. The value of k is fixed based on the number of

clusters required. We can consider the normalized Laplacian where $L = I - D^{-1}W$. Therefor for $k$ clusters, the solution can be achieved by taking the first $k$ eigenvectors of $I - D^{-1}W$ as the columns of $H$.
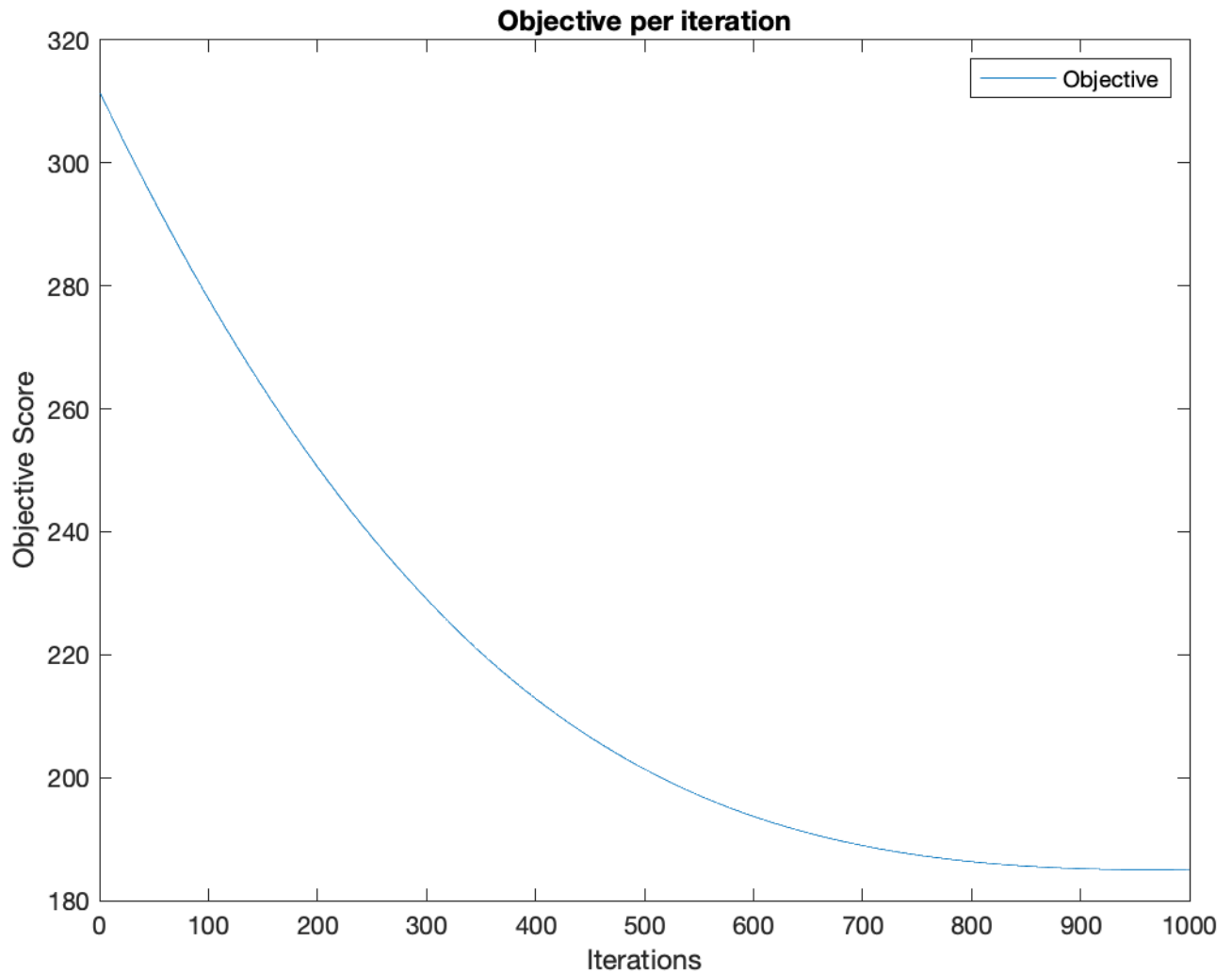
# Problem 5 [10 pts]

We consider the following optimization problem ($\mathbf{Y}$ is given and generated randomly):

$$\min_{\mathbf{X}} \frac{1}{2}\|\mathbf{X} - \mathbf{Y}\|_F^2 + \|\mathbf{X}\|_*  \tag{8}$$

where $\mathbf{Y}, \mathbf{X} \in \mathbb{R}^{100 \times 100}$ and $\|\cdot\|_*$ denotes the nuclear norm (sum of singular values). Now please use gradient descent method to update $\mathbf{X}$. ($\frac{\partial \|\mathbf{X}\|_*}{\partial \mathbf{X}} = \mathbf{U}\mathbf{V}^T$, where $\mathbf{U}, \mathbf{V}$ is obtained from reduced SVD, namely $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = svd(\mathbf{X}, 0)$). Plot the objective changes with 1000 iteration.

I first noticed that the optimization problem can be rewritten as $tr((X - Y)^T(X - Y)) + \|\mathbf{X}\|_*$ In which the gradient $w.r.t.X$ can be seen as just being the gradient $(X - Y)$ plus the nuclear norm gradient which is given as $\frac{\partial \|\mathbf{X}\|_*}{\partial \mathbf{X}} = \mathbf{U}\mathbf{V}^T$. Through a web search on Matlab, I found that calculating the nuclear norm is as follows, $nucnorm = norm(svd(X), 1)$; In my code I used a lambda of .001.
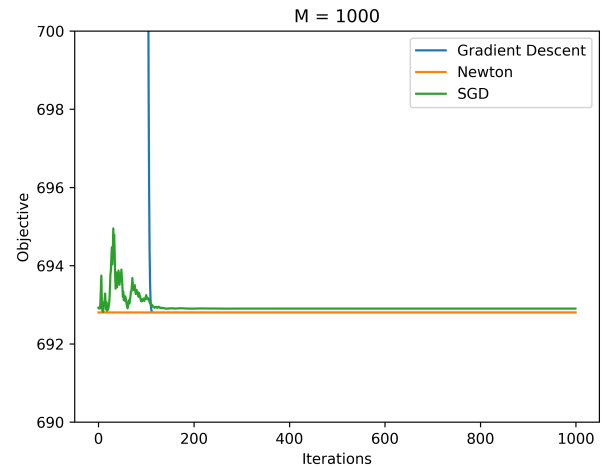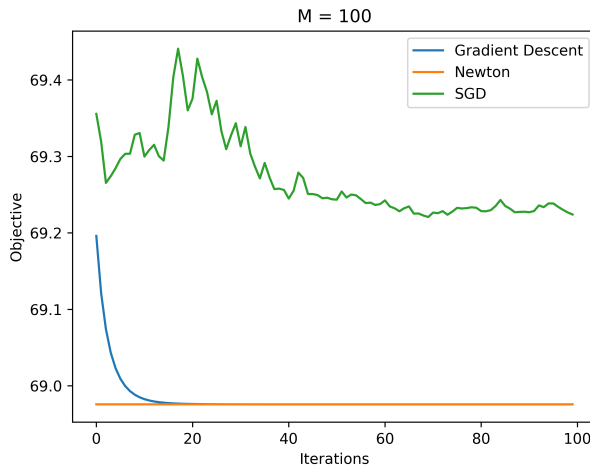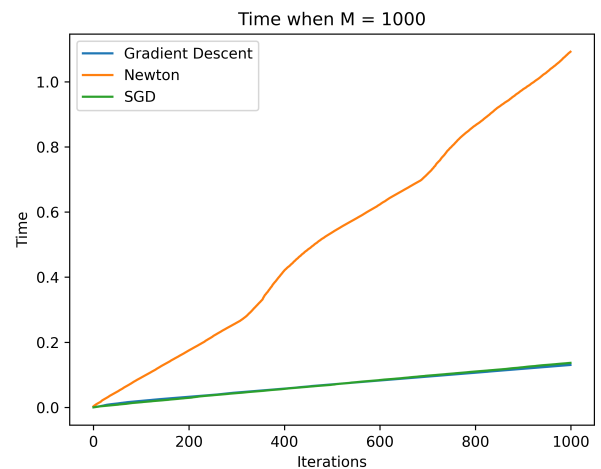


7

# Problem 6 [20 pts]
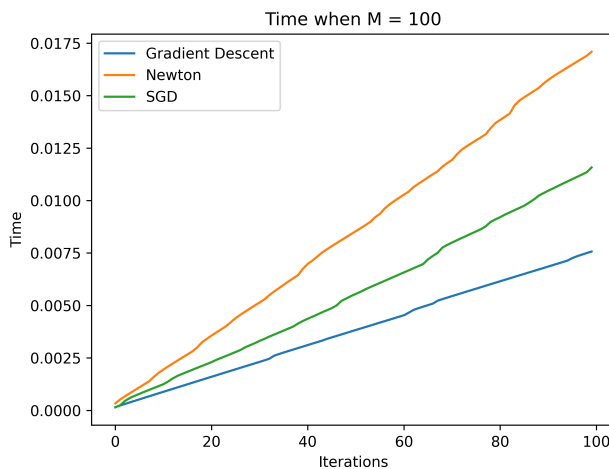
We turn to Logistic Regression:

$$\min_{\beta} \sum_{i=1}^{m} ln(1 + e^{\langle \beta, \hat{x_i} \rangle}) - y_i \langle \beta, \hat{x_i} \rangle, \tag{9}$$

where $\beta = (w; b), \hat{x} = (x; 1)$. Assume $m = 100, x \in \mathbb{R}^{99}$. Please randomly generate $x, y$ and find the optimal $\beta$ via 1) gradient descent; 2) Newton's method and 3) stochastic gradient descent (SGD) where the batch-size is 1. (need consider choosing appropriate step-size if necessary). Change $m = 1000, x \in \mathbb{R}^{999}$, observe which algorithm will decrease the objective faster in terms of iteration ($X$-axis denotes number of iteration) and CPU time. [You will receive another 5 bonus points if you implement backtracking line search]

### Norm of the matrix per iteration



### Time per iteration



There are several observations I made from these graphs.

8

First I see that SGD does not converge, while is this unusual it is not expected to converge since the batch size is 1 and I am taking the gradient of 1 random training sample. This wouldn't matter as much however x and y are randomly generated with no specific pattern. Therefore the gradient of 1 point may be very different from the gradient of another. When M = 1000 Gradient descent appears from the top and quickly converges to 0. I found that in the beginning the Norm is in the 3000s and then it starts to converge. I believe this is because with a higher training set there is room for higher error. I also noticed that for all random data, the best $\beta$ was normally not far from [0, 0],from my understanding this is because none of the data has any real world pattern since it is all randomly generated. Because of the close intial guess, Newton's method converge within < 5 iterations which does not appear on the graphs.

When M = 100 Newton's method takes the longest unsurprisingly, however SGD takes 5 milliseconds longer than Gradient Descent. I believe that in this case the extra line of code in the SGD gradient calculation creating a random number accounts for this difference as I am calling an external function. This hypothesis is further justified as SGD is narrowly faster than Gradient Descent when M = 1000. Therefore I believe when M is Larger SGD will be faster than Gradient Descent, and at that time Newton's method would outlive the entire human race with how long it would take to calculate.

My code is at the bottom of exam, I didn't include the plot code as it is just plotting the data.

# Problem 7 [10 pts]

Please design an (either toy or real-world) experiment to demonstrate that PCA can be helpful for denoising.

In my experiment I took the Lenna image and an image from the MNIST dataset. I then added noise to it and performed PCA by SVD of the first 50 columns of each image. Code will be included at end of exam.

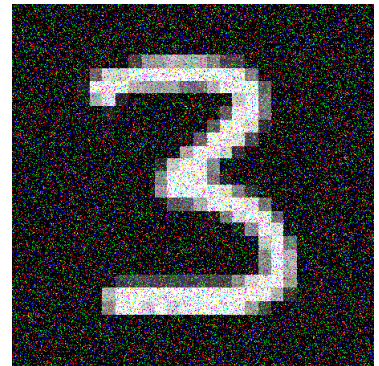Original Image



Image with noise



Image with reduced noise by PCA

# Bonus Problem 8 [10 pts]

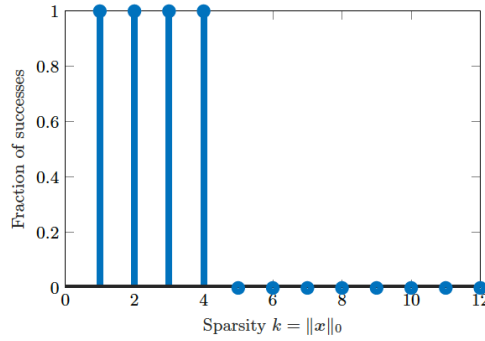$$\textbf{Solve:} \quad \min \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{y}. \tag{10}$$

We have proved that if $\mathbf{y} = \mathbf{Ax}_o$ with

$$\|\mathbf{x}_o\|_0 \leq \tfrac{1}{2}\,\mathrm{krank}(\mathbf{A}). \tag{11}$$

Then $\mathbf{x}_o$ is the unique optimal solution to the $\ell^0$ minimization problem

$$\min \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{y}. \tag{12}$$

However, when $\mathbf{A}$ is of size $5 \times 12$, the following figure illustrates the fraction of success across 100 trials. Apparently $krank(\mathbf{A}) \leq rank(\mathbf{A}) \leq 5$, therefore, when sparsity $k = 1, 2$ satisfying Eq. (11)



it has 100% recovery success rate is not surprising. However, the above experiment also shows even $k = 3, 4$ which violates Eq. (11), still it can be recovered at 100%. Please explain this phenomenon.

**Proof:**

$$A\hat{x} = y \rightarrow A(\hat{x} - x_0) = A\hat{x} - Ax_0 = y - y = 0 \tag{13}$$

 

This phenomenon can be explain because while the data theoretically cannot be recovered in the $\ell^0$ norm, it can be recovered in the $\ell^1$ norm. While $\ell^0$ is not convex, $\ell^1$ is and can be used as a convex surrogate for $\ell^0$. $|x|$ is the convex envelope for $||x||_0$. We know from experimental results like the one above that $\ell^1$ minimization succeeds with very high rate whenever the constants of n/m and k/m are relatively small. Here we introduce the Restricted Isometry Property (RIP) of order $k$ with constant $\delta \in [0, 1)$ if

$$\forall \mathbf{x}\; k - sparse, \quad (1 - \delta)||x||_2^2 \leq ||\mathbf{Ax}||_2^2 \leq (1 = \delta)||x||_2^2 \tag{14}$$

For $\ell^0$ norm, we say that $y = Ax_o$, with $k = ||x_o||_0$. If $\delta_{2k}(A) < 1$, then $x_o$ is the unique optimal solution to

$$min||x||_0 \quad s.t.\; Ax = y \tag{15}$$

We can prove this by contradiction, let's say there exist an $x' \neq x$ with $||x'||_o \leq k$ Then $x_o - x' \in null(A)$, and $||x_o - x'||_0 \leq 2k$. This implies that $\delta_{2k}(A) \geq 1$, contradicting our statement and therefore we can conclude the original statement is correct.

For $\ell^1$ norm recovery under RIP, we say that $y = Ax_o$, with $k = ||x_o||_0$. If $\delta_{2k}(A) < \sqrt{2} - 1$, then $x_o$ is the unique optimal solution to:

$$min||x||_1 \quad s.t. \ Ax = y \tag{16}$$

Candes and Tao proved in 2006 originally that $\delta_{2k} < \sqrt{2} - 1$

Computing the RIP constant $\delta_k(A)$ is an np-hard problem, however we know that there exists a numerical constant C > 0 such that if $A \in \mathbb{R}^{mxn}$ is a random matrix with entries independent $\mathcal{N}(0, \frac{1}{m})$ random variables, with high probability, $delta_k(A) < \delta$, provided that

$$m \geq Ck \, log(n/k)/\delta^2 \tag{17}$$

With incoherence, we need $m \geq \Omega(k^2)$. Here this result allows (k,m,n) to scale proportionally: $m \geq \Omega(k)$.

12

# 1 Code

## 1.1 Minmize Elastic Net Objective Problem 1

```
1   clc; clear;
2
3   rng('default');
4
5   X = rand(100, 9);
6   y = rand(100, 1);
7   beta = rand(9, 1);
8   learning_rate = 0.0001;
9   ob_scores = [];
10  t = y - X * beta;
11  for i=1:100
12      grad = beta - 2 * X' * t + 0.5 * sign(beta);
13      beta = beta - learning_rate * grad;
14      t = y - X * beta;
15      ob = (norm(t)^2) + (0.5 * norm(beta)^2) + (0.5 * norm(beta, 1));
16      ob_scores = [ob_scores ob];
17  end
18
19  h = figure();
20  plot(ob_scores)
21  xlabel("Iterations")
22  ylabel("Objective Score")
23  title("Elastic Net Objective per Iteration")
24  legend(["Objective"])
25  saveas(gcf, "elastic_net_ob.png")
26  waitfor(h)
```

## 1.2 Minimie Objective Problem 5

```
1   clc; clear;
2
3   %want to have the same random samples to replicate results
4   rng('default');
5
6   X = rand(100, 100);
7   Y = rand(100, 100);
8   lambda = 0.001;
9   ob = objective_p5(X, Y);
10  ob_scores = [];
11  for i=1:1000
12      [U,S,V] = svd(X, 0);
13      gradient = (X - Y) + U*V';
```

```matlab
14      X = X − lambda ∗ gradient;
15      ob = objective_p5(X, Y);
16      %add objective to array for graph
17      ob_scores = [ob_scores ob];
18  end
19
20  h = figure();
21  plot(ob_scores);
22  ylabel("Objective Score")
23  xlabel("Iterations")
24  title("Objective per iteration")
25  legend(["Objective"]);
26  saveas(gcf, 'ob_min.png')
27  waitfor(h);
```

```matlab
1  function [cost] = objective_p5(x, y)
2      cost = 0.5 ∗ norm(x − y, "fro") + norm(svd(x), 1);
3  end
```

## 1.3   Logistic Regression Problem 6

```python
1  import numpy as np
2  import time
3
4  #objective
5  def f(beta, X, Y):
6      return np.ravel(np.ones(len(Y))∗(np.log(1+np.exp(X∗beta)))−Y.T∗X∗
           beta)[0]
7
8  #first derivative of objective
9  def derivative_grad_f(beta, X, Y):
10     return X.T∗(1/(1+1/np.exp(X∗beta))−Y)
11
12 #SGD is different in vanilla gradient descent in that it randomly
       chooses a batch of x rather than use the whole dataset.
13 #The final said that the batch_size = 1
14 def derivative_SGD_f(beta, X, Y, m):
15     i = np.random.randint(0, m)
16     return X[i, :].T∗(1/(1+1/np.exp(X[i, :]∗beta))−Y[i, :])
17
18 #second derivative of objective
19 def second_derivative_f(beta, X, Y):
20     return X.T∗(np.diag(np.ravel(np.exp(X∗beta)/np.power(1+np.exp(X∗
           beta),2))))∗X)
21
22 if __name__ == "__main__":
```

```python
23
24      #create data arrays
25      objective_gd_100 = []
26      objective_new_100 = []
27      objective_sgd_100 = []
28      time_gd_100 = []
29      time_new_100 = []
30      time_sgd_100 = []
31
32      objective_gd_1000 = []
33      objective_new_1000 = []
34      objective_sgd_1000 = []
35      time_gd_1000 = []
36      time_new_1000 = []
37      time_sgd_1000 = []
38
39      size_array = [100, 1000]
40      for i in range(len(size_array)):
41
42          m = size_array[i]
43          X = np.matrix(np.random.rand(m, 2) * 10)
44          Y = np.matrix(np.random.randint(0, 2, (m, 1)))
45          beta = np.array([[0.], [0.]])
46          learning_rate = 0.001
47
48          t = time.time()
49          counter = 0
50          #gradient descent
51          for i in range(m):
52              counter += 1
53              beta = beta - learning_rate*derivative_grad_f(beta, X, Y)
54              learning_rate = learning_rate/1.02
55              if m == 100:
56                  objective_gd_100.append(f(beta, X, Y))
57                  time_gd_100.append(time.time() - t)
58              else:
59                  objective_gd_1000.append(f(beta, X, Y))
60                  time_gd_1000.append(time.time() - t)
61
62          print('iter =',counter)
63          print(beta)
64          print('norm =',np.linalg.norm(derivative_grad_f(beta, X, Y)))
65          print("\n")
66
67          beta = np.array([[0.], [0.]])
```

```python
68              counter = 0
69              t = time.time()
70             #newton's method
71             for i in range(m):
72                  counter += 1
73                  beta = beta -  np.linalg.inv(second_derivative_f(beta, X, Y
                        ))*derivative_grad_f(beta, X, Y)
74                  if m == 100:
75                      objective_new_100.append(f(beta, X, Y))
76                      time_new_100.append(time.time() - t)
77                  else:
78                      objective_new_1000.append(f(beta, X, Y))
79                      time_new_1000.append(time.time() - t)
80
81             print('iter =',counter)
82             print(beta)
83             print('norm =',np.linalg.norm(derivative_grad_f(beta, X, Y)))
84             print("\n")
85
86             counter = 0
87             beta = np.array([[0.], [0.]])
88             t = time.time()
89             learning_rate = 0.001
90             #Stochastic Gradient Descent
91             for i in range(m):
92                  counter += 1
93                  beta = beta - learning_rate*derivative_SGD_f(beta, X, Y, m)
94                  #lower the learning rate over time to prevent jumping
95                  learning_rate = learning_rate/1.02
96                  if m == 100:
97                      objective_sgd_100.append(f(beta, X, Y))
98                      time_sgd_100.append(time.time() - t)
99                  else:
100                     objective_sgd_1000.append(f(beta, X, Y))
101                     time_sgd_1000.append(time.time() - t)
102
103            print('iter =',counter)
104            print(beta)
105            print('norm =',np.linalg.norm(derivative_SGD_f(beta, X, Y, m)))
106            print("\n")
107
108            if m == 100:
109                np.save("data/objective_gd_100.npy", np.array(
                       objective_gd_100))
110                np.save("data/objective_new_100.npy", np.array(
```

```
                    objective_new_100))
111             np.save("data/objective_sgd_100.npy", np.array(
                    objective_sgd_100))
112             np.save("data/time_gd_100.npy", np.array(time_gd_100))
113             np.save("data/time_new_100.npy", np.array(time_new_100))
114             np.save("data/time_sgd_100.npy", np.array(time_sgd_100))
115         else:
116             np.save("data/objective_gd_1000.npy", np.array(
                    objective_gd_1000))
117             np.save("data/objective_new_1000.npy", np.array(
                    objective_new_1000))
118             np.save("data/objective_sgd_1000.npy", np.array(
                    objective_sgd_1000))
119             np.save("data/time_gd_1000.npy", np.array(time_gd_1000))
120             np.save("data/time_new_1000.npy", np.array(time_new_1000))
121             np.save("data/time_sgd_1000.npy", np.array(time_sgd_1000))
```

## 1.4  Denoise Problem 7

```matlab
1   clear; clc;
2   for i=1:2
3       if i == 1
4           A = imread("3.png");
5       else
6           A = imread("Lenna.png");
7       end
8       A = imnoise(A, 'salt & pepper', .2);
9       h = figure();
10      imshow(A)
11      waitfor(h);
12      if i == 1
13          imwrite(A,"3_noise.png");
14      else
15          imwrite(A,"Lenna_noise.png");
16      end
17
18      k=50;
19      R = A(:, :, 1);
20      G = A(:, :, 2);
21      B = A(:, :, 3);
22
23      IR = im2double(R);
24      IG = im2double(G);
25      IB = im2double(B);
26
27      [u1,s1,v1] = svd(IR);
```

```matlab
28        [u2,s2,v2] = svd(IG);
29        [u3,s3,v3] = svd(IB);
30
31        c1 = zeros(size(IR));
32        c2 = zeros(size(IG));
33        c3 = zeros(size(IB));
34
35        c1 = u1(:, 1:k) * s1(1:k, 1:k) * v1(:, 1:k)';
36        c2 = u2(:, 1:k) * s2(1:k, 1:k) * v2(:, 1:k)';
37        c3 = u3(:, 1:k) * s3(1:k, 1:k) * v3(:, 1:k)';
38        siz1 = size(IR(:, 1));
39        siz1 = siz1(1);
40        siz2 = size(IR(1, :));
41        siz2 = siz2(2);
42        q = zeros(siz1, siz2, 3);
43        q(:, :, 1) = c1;
44        q(:, :, 2) = c2;
45        q(:, :, 3) = c3;
46        h = figure();
47        imshow(q)
48        waitfor(h);
49        if i == 1
50            imwrite(q, "3_reduced_noise.png");
51        else
52            imwrite(q,"Lenna_reduced_noise.png");
53        end
54  end
```