

Homework Set 2, CPSC 8420, Fall 2023

Your Name

Due 10/26/2023, Thursday, 11:59PM EST

Problem 1

For PCA, from the perspective of maximizing variance, please show that the solution of ϕ to maximize $\|\mathbf{X}\phi\|_2^2$, s.t. $\|\phi\|_2 = 1$ is exactly the first column of \mathbf{U} , where $[\mathbf{U}, \mathbf{S}] = svd(\mathbf{X}^T \mathbf{X})$. (Note: you need prove why it is optimal than any other reasonable combinations of \mathbf{U}_i , say $\hat{\phi} = 0.8 * \mathbf{U}(:, 1) + 0.6 * \mathbf{U}(:, 2)$ which also satisfies $\|\hat{\phi}\|_2 = 1$.)

Problem 2

Given matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ (assume each column is centered already), where n denotes sample size while p feature size. To conduct PCA, we need find eigenvectors to the largest eigenvalues of $\mathbf{X}^T \mathbf{X}$, where usually the complexity is $\mathcal{O}(p^3)$. Apparently when $n \ll p$, this is not economic when p is large. Please consider conducting PCA based on $\mathbf{X} \mathbf{X}^T$ and obtain the eigenvectors for $\mathbf{X}^T \mathbf{X}$ accordingly and use experiment to demonstrate the acceleration.

Problem 3

Let's revisit Least Squares Problem: minimize $\frac{1}{2} \|\mathbf{y} - \mathbf{A}\boldsymbol{\beta}\|_2^2$, where $\mathbf{A} \in \mathbb{R}^{n \times p}$.

1. Please show that if $p > n$, then vanilla solution $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$ is not applicable any more.
2. Let's assume $\mathbf{A} = [1, 2, 4; 1, 3, 5; 1, 7, 7; 1, 8, 9], \mathbf{y} = [1; 2; 3; 4]$. Please show via experiment results that Gradient Descent method will obtain the optimal solution with Linear Convergence rate if the learning rate is fixed to be $\frac{1}{\sigma_{max}(\mathbf{A}^T \mathbf{A})}$, and $\boldsymbol{\beta}_0 = [0; 0; 0]$.
3. Now let's consider ridge regression: minimize $\frac{1}{2} \|\mathbf{y} - \mathbf{A}\boldsymbol{\beta}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2$, where $\mathbf{A}, \mathbf{y}, \boldsymbol{\beta}_0$ remains the same as above while learning rate is fixed to be $\frac{1}{\lambda + \sigma_{max}(\mathbf{A}^T \mathbf{A})}$ where λ varies from $0.1, 1, 10, 100, 200$, please show that Gradient Descent method with larger λ converges faster.

Problem 4

We consider matrix completion problem. As we discussed in class, the main issue of *softImpute (Matrix Completion via Iterative Soft-Thresholded SVD)* is when the matrix size is large, conducting *SVD* is computational demanding. Let's recall the original problem where $\mathbf{X}, \mathbf{Z} \in \mathbb{R}^{n \times d}$:

$$\min_{\mathbf{Z}} \frac{1}{2} \|P_{\Omega}(\mathbf{X}) - P_{\Omega}(\mathbf{Z})\|_F^2 + \lambda \|\mathbf{Z}\|_* \quad (1)$$

People have found that instead of finding optimal \mathbf{Z} , it might be better to make use of *Burer-Monteiro* method to optimize two matrices $\mathbf{A} \in \mathbb{R}^{n \times r}, \mathbf{B} \in \mathbb{R}^{d \times r} (r \geq \text{rank}(\mathbf{Z}^*))$ such that $\mathbf{AB}^T = \mathbf{Z}$. The new objective is:

$$\min_{\mathbf{A}, \mathbf{B}} \frac{1}{2} \|P_{\Omega}(\mathbf{X} - \mathbf{AB}^T)\|_F^2 + \frac{\lambda}{2} (\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2). \quad (2)$$

- Assume $[\mathbf{U}, \boldsymbol{\Sigma}, \mathbf{V}] = svd(\mathbf{Z})$, show that if $\mathbf{A} = \mathbf{U} \boldsymbol{\Sigma}^{\frac{1}{2}}, \mathbf{B} = \mathbf{V} \boldsymbol{\Sigma}^{\frac{1}{2}}$, then Eq. (2) is equivalent to Eq. (1).
- The *Burer-Monteiro* method suggests if we can find $\mathbf{A}^*, \mathbf{B}^*$, then the optimal \mathbf{Z} to Eq. (1) can be recovered by $\mathbf{A}^* \mathbf{B}^{*T}$. It boils down to solve Eq. (2). Show that we can make use of least squares with ridge regression to update \mathbf{A}, \mathbf{B} row by row in an alternating minimization manner as below. Assume $n = d = 2000, r = 200$, please write program to find \mathbf{Z}^* .

```

 $T \leftarrow 100, i \leftarrow 1$  % you can also set T to be other number instead of 100
if  $i \leq T$  then
    update A row by row while fixing B
    update B row by row while fixing A
     $i \leftarrow i + 1$ 
end if

```

Problem 1:

For matrix with self centred data,
max variance :

$$\max \frac{\|X\phi\|_2^2}{\|\phi\|_2} \quad [\phi \rightarrow \text{vector}]$$

$$\|\phi\|_2 = 1$$

To prove: solution of ϕ for above eqn
is exactly the first column of V

where $[U, S] = \text{svd}(X^T X)$

→ we can modify $\|X\phi\|_2^2$ as

$$[(X\phi)]^T [X\phi] = \phi^T X^T X \phi$$

As $X^T X$ is van spd matrix,

for svd of $X^T X$, we have
 $V = U$ in $U \Sigma V^T$.

Also, singular values are absolute
values of eigen values.

Also, columns of $U_2 V$ are eigen vectors
of A .

Consider,

Case 1: when $\phi = aU_1 + bU_2$.

where a and b are real numbers

S.t. $a^2 + b^2 = 1$

$$\begin{aligned} \text{So, } \| \phi \|^2 &= \phi^T \cdot \phi \\ &= (aU_1 + bU_2)^T (aU_1 + bU_2) \\ &= a^2 U_1^T U_1 + b^2 U_2^T U_2 + 2ab U_1^T U_2 \end{aligned}$$

As columns of $U (U_1, U_2, \dots)$ are
orthonormal vectors, $U_1^T U_2 = 0$

and $U_1^T U_1 = I$

$$U_2^T U_2 = I$$

so, $\| \phi \|^2 = 1$.

i.e., $\phi = aU_1 + bU_2$, ($a, b \in \mathbb{R}$)

s.t. $a^2 + b^2 = 1$

satisfy the given constraint

9. If we substitute the value of ϕ in : $\phi^T X^T X \phi$

$$\phi^T [U \Sigma U^T] \phi$$

$$= \phi^T \left[\sum_{i=1}^p \sigma_i U_i U_i^T \right] \phi$$

$$(u = \phi^T [\sigma_1 U_1 U_1^T + \sigma_2 U_2 U_2^T$$

$$+ \sigma_3 U_3 U_3^T + \sigma_4 U_4 U_4^T]) \phi$$

$$= (au_1 + bu_2)^T (\sigma_1 u_1 u_1^T + \sum u_2 u_2^T + \dots) (au_1 + bu_2)$$

$$= (au_1^T \sigma_1 u_1 u_1^T + b u_2^T \sigma_2 u_2 u_2^T$$

$$+ a u_1^T \sigma_2 u_2 u_2^T + \dots) / (au_1 + bu_2)$$

$$= (a^2 \sigma_1 u_1^T u_1 u_1^T + b^2 \sum u_2^T u_2 u_2^T + \dots + a u_1^T \sigma_2 u_2 u_2^T u_1^T + \dots)$$

As U_1 & U_2 are orthonormal vectors,

all terms with $U_1^T \cdot U_2$ & $U_2^T \cdot U_1 = 0$.

so, we are left with:

$$a^2 \sigma_1 + b^2 \sigma_2$$

$$\sigma_1 (a^2 + b^2 \frac{\sigma_2}{\sigma_1}) - ①$$

As we know that $\sigma_1 > \sigma_2$ (in Σ matrix, singular values are arranged along the diagonal in value decreasing order)

So, from eqn ① :

$$a^2 + b^2 = 1$$

$$a^2 + b^2 \left(\frac{\sigma_2}{\sigma_1} \right) < 1$$

$$\Rightarrow \sigma_1 (a^2 + b^2 \frac{\sigma_2}{\sigma_1}) < \sigma_1 - ②$$

SUNDAY 12

APRIL 2023

Case 2: Consider if $\phi = U_1$ (i.e., first column of U)

$$\|\phi\|_2^2 = 1 \quad (\text{orthonormal vector})$$

so, putting $\phi = U_1$ in:

$$\phi^T [V \Sigma V^T] \phi$$

$$= \phi^T [\sigma_1 U_1 U_1^T + \sigma_2 U_2 U_2^T + \sigma_3 U_3 U_3^T + \dots] \phi$$

$$= \phi^T [\sigma_1 U_1 U_1^T + \sigma_2 U_2 U_2^T + \sigma_3 U_3 U_3^T + \dots] \phi$$

$$= \phi^T [U_1^T [\sigma_1 U_1 U_1^T + \sigma_2 U_2 U_2^T + \sigma_3 U_3 U_3^T + \dots]] \phi$$

$$= \sigma_1 U_1^T U_1 U_1^T U_1 + \sigma_2 U_1^T U_2 U_2^T U_1 +$$

$$\sigma_3 U_1^T U_3 U_3^T U_1 + \dots$$

As columns of V matrix are orthonormal,

$$U_{i(i+1)}^T \cdot U_i \text{ & } U_i^T \cdot U_{i(i+1)} = 0.$$

So, we get simplified solution

$$= \sigma_i U_i^T U_i U_i^T U_i$$

$$= \underline{\sigma_i} - (3)$$

Comparing $W = U$ with argument to a general combination of vectors from U ,

$$\arg \max_{\|\phi\|_2^2 = 1} \|X\phi\|_2^2 \Rightarrow W = aU_1 + bU_2$$

(st. $a, b \in \mathbb{R}$)

$$\text{gives } \sigma_i \left(\frac{a^2 + b^2}{\sigma_i^2} \right),$$

which is less than σ_i , when

$$W = U_1$$

so, max. variance is achieved when $\phi = \text{first column of } U$, where $U_2 [U, S] = \text{svd}(X^T X)$.

Problem 2

Let $A = X^T X$ — (1)
& $B = X X^T$.

Let U be eigen vectors of B .

$$X X^T U = \lambda U$$

Multiply both sides by X^T :

$$X^T X X^T U = \lambda X^T U$$

from (1), $X^T X = A$

$$A(X^T U) = \lambda(X^T U)$$

So, we can conclude that $X^T U$ are eigen vectors of A .

So, we can find eigen vectors of B , which have the computational

complexity of $O(n^3)$ instead of $O(p^3)$, where $X \in \mathbb{R}^{n \times p}$.

These eigen vectors can be used to find eigen vectors of $A(X^T X)$.

To convert eigen vectors of $X^T X$ to those of $X^T X$

→ Multiply eigen vectors of $B(X X^T)$ by X^T and then normalize.

i.e. eigen vectors of A : $X^T U$

But $X^T U$ is not normalized.
so, the matrix which gives us eigen vectors of A is $X^T U$ with all of its column vectors normalized.

Problem 2

```
% Experimental Setup
n = 100; % number of samples
p = 10000; % number of features

% Generate random data for matrix X (centered)
X = randn(n, p);

% Start timer for PCA using X^T X
tic;

% Compute covariance matrix X^T X
C1 = X' * X;

% Compute eigenvectors and eigenvalues of X^T X
[V1, D1] = eig(C1);

% Sort eigenvectors based on eigenvalues in descending order
[~, idx1] = sort(diag(D1), 'descend');
V1 = V1(:, idx1);

time_XTX = toc;

% Start timer for PCA using XX^T
tic;

% Compute covariance matrix XX^T
C2 = X * X';

% Compute eigenvectors and eigenvalues of XX^T
[V2, D2] = eig(C2);

% Convert eigenvectors of XX^T to those of X^T X
V2 = X' * V2;

% Normalize the eigenvectors
for i = 1:n
    V2(:,i) = V2(:,i) / norm(V2(:,i));
end

% Sort eigenvectors based on eigenvalues in descending order
[~, idx2] = sort(diag(D2), 'descend');
V2 = V2(:, idx2);

time_XXT = toc;

% Display the results
fprintf('Time for PCA using X^T X: %.4f seconds\n', time_XTX);
```

```
Time for PCA using X^T X: 71.1933 seconds
```

```
fprintf('Time for PCA using XX^T: %.4f seconds\n', time_XXT);
```

```
Time for PCA using XX^T: 0.0941 seconds
```

Problem 3 (1)

Condition: $n < p$

```
n = 10; % number of samples
p = 10000; % number of features

% Generate random data for matrix X (centered)
A = randn(n, p);

M = A'*A;
disp("Size of the square matrix:");
```

Size of the square matrix:

```
disp(size(M));
```

10000 10000

```
disp("Rank of the square matrix");
```

Rank of the square matrix

```
disp(rank(M));
```

10

Condition $n > p$:

```
n = 1000; % number of samples
p = 100; % number of features

% Generate random data for matrix X (centered)
A = randn(n, p);

M = A'*A;
disp("Size of the square matrix:");
```

Size of the square matrix:

```
disp(size(M));
```

100 100

```
disp("Rank of the square matrix");
```

Rank of the square matrix

```
disp(rank(M));
```

100

As the rank of the matrix $<$ size of the matrix, $p > n$ condition leads to a non invertible matrix. So, the closed form solution fails.

Problem 3 (2)

```
% Given A, y, lr and initial weights
A = [1 2 4; 1 3 5; 1 7 7; 1 8 9];
y = [1; 2; 3; 4];

% Finding the maximum singular value of A
max_singular_value = max(svd(A'*A));

% Learning rate is fixed as 1/max singular value
lr = 1/max_singular_value;
weights = [0; 0; 0];

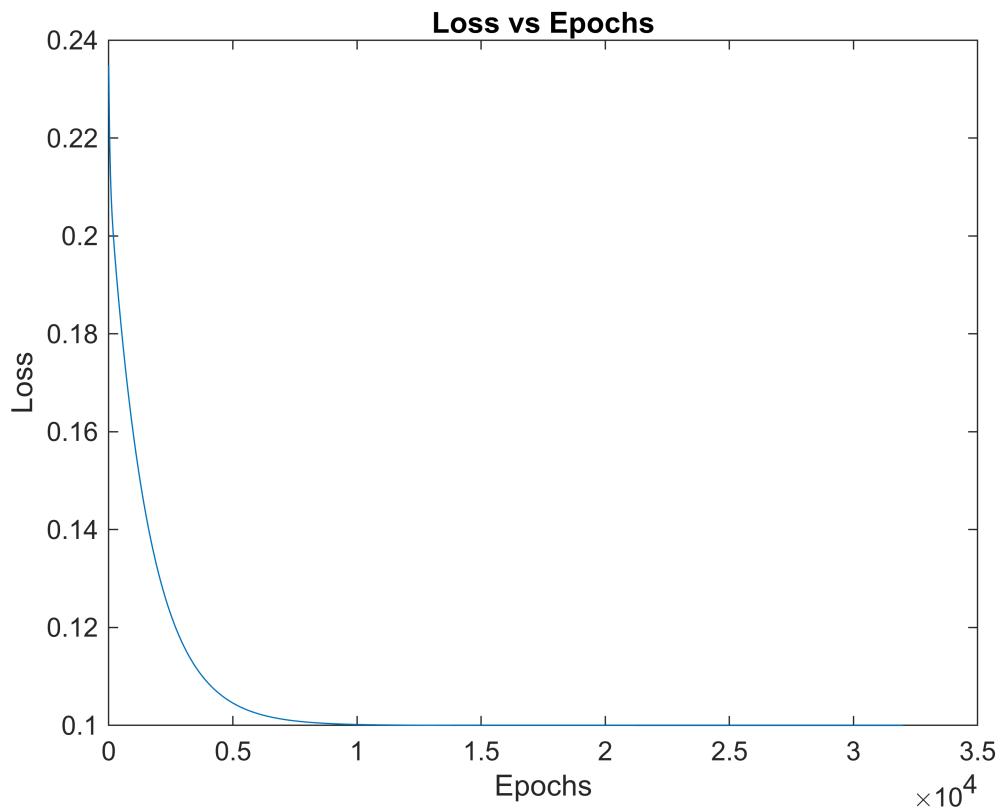
% Initializing a loss row vector to track progress
loss = [];

% Implementing Vanilla Gradient Descent (32k epochs)
for i = 1:32000
    grad = A'*(A*weights - y);
    weights = weights - grad*lr;
    loss(i) = (A*weights - y)' * (A*weights - y);
end

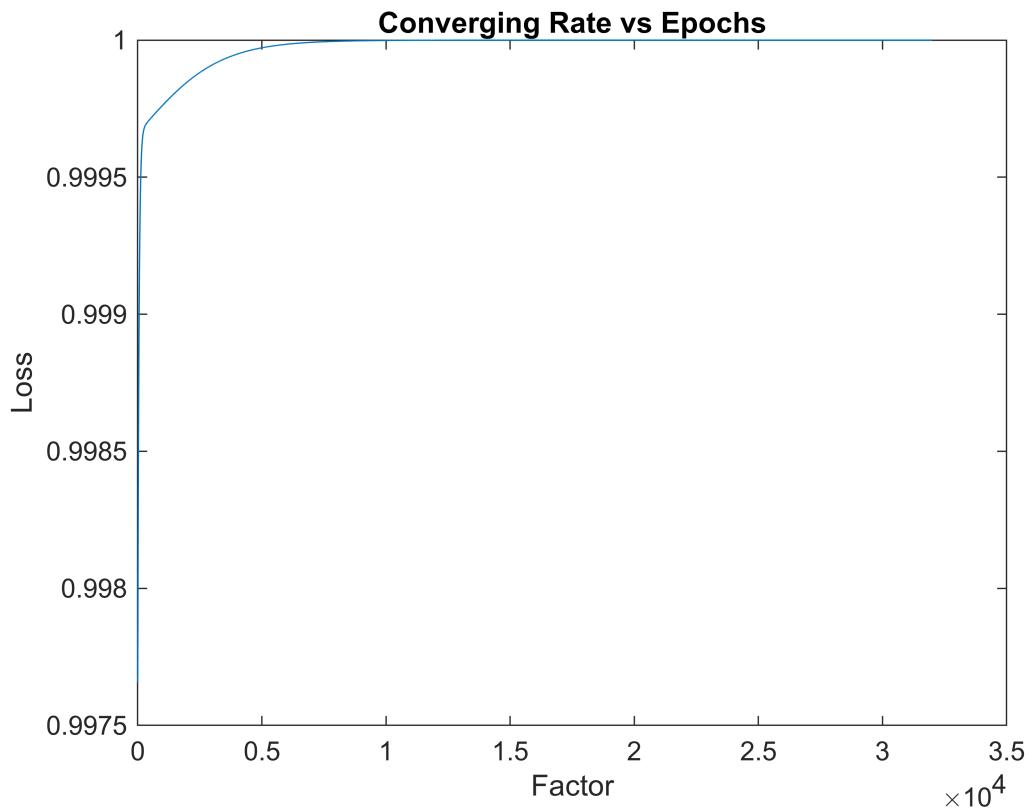
convergence_factor = loss(2:end) ./ loss(1:end-1);
display(convergence_factor);
```

```
convergence_factor = 1x31999
 0.9977    0.9977    0.9977    0.9978    0.9978    0.9978    0.9979    0.9979 ...
```

```
% Plotting the Loss obtained with time
plot(loss);
title("Loss vs Epochs");
xlabel("Epochs");
ylabel("Loss");
```



```
% Plotting the convergence factor for each epoch
plot(convergence_factor);
title("Converging Rate vs Epochs");
xlabel("Factor");
ylabel("Loss");
```



```
%Verifying the results obtained from closed form solution
weights2 = pinv(A'*A)*A'*y;
```

```
display(weights2);
```

```
weights2 = 3x1
-1.0000
 0.0333
 0.5333
```

```
display(weights);
```

```
weights = 3x1
-1.0000
 0.0333
 0.5333
```

From the above weight matrices, we can deduce that Gradient Descent method obtains the optimal solution with almost a Linear Convergence rate with the given learning rate.

Problem 3 (3)

```
% Define the matrix A and vector y
```

```

A = [1, 2, 4; 1, 3, 5; 1, 7, 7; 1, 8, 9];
y = [1; 2; 3; 4];

% Define a range of λ values
lambda_values = [0.1, 1, 10, 100, 200];

% Define the number of iterations
num_iterations = 20;

% Initialize variables to store losses
losses = zeros(length(lambda_values), num_iterations);

% Perform gradient descent for each λ
for k = 1:length(lambda_values)
    lambda = lambda_values(k);
    alpha = 1 / (max(max(A' * A)) + lambda_values(k));

    % Initialize θ to [0; 0; 0]
    theta = [0; 0; 0];
    for iteration = 1:num_iterations
        % Compute the gradient for ridge regression
        gradient = -A' * (y - A * theta) + lambda * theta;

        % Update θ using the gradient and learning rate
        theta = theta - alpha * gradient;

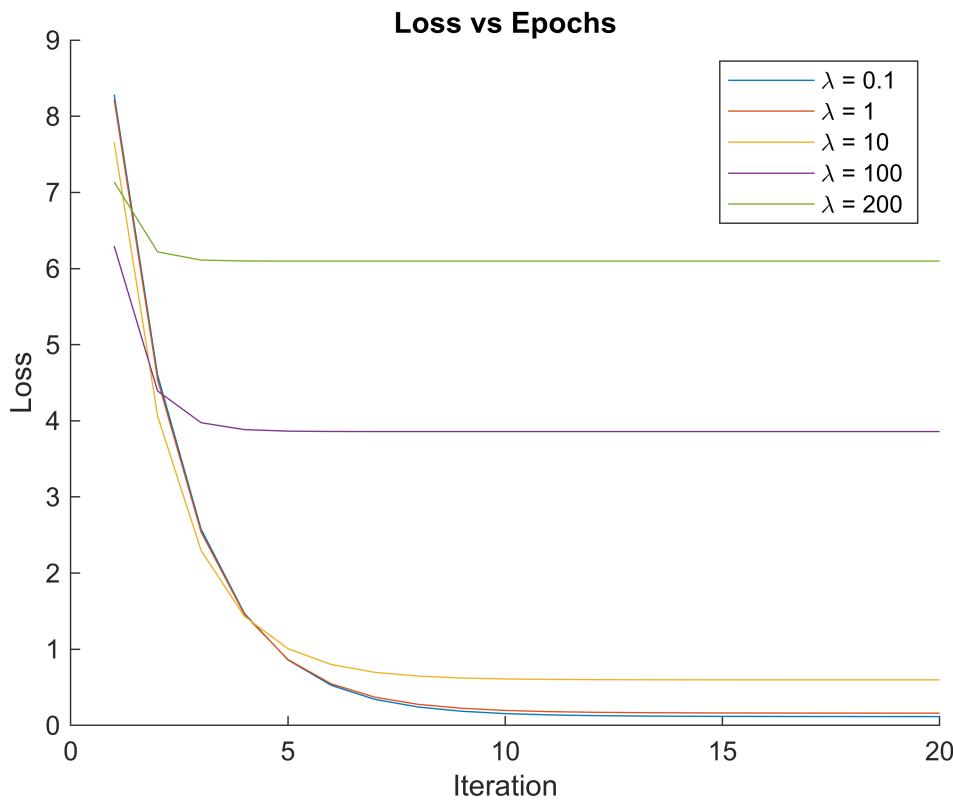
        % Calculate the cost (ridge regression loss)
        loss = (1 / (2 )) * norm(A * theta - y)^2 + (theta' * theta) * (lambda / 2);

        % Store the loss
        losses(k, iteration) = loss;
    end
end

% Plot the loss curves for different λ values
figure;
hold on;
for k = 1:length(lambda_values)
    % plot(1:num_iterations, losses(k, :), 'DisplayName', ['λ = ', num2str(lambda_values(k))]);
    semilogy(1:num_iterations, losses(k, :), 'DisplayName', ['\lambda = ', num2str(lambda_values(k))]);
end

hold off;
xlabel('Iteration');
ylabel('Loss');
title("Loss vs Epochs");
legend;

```



```
%Verifying the results obtained from closed form solution
```

```
weights2 = pinv(A'*A + lambda_values(5)*eye(size(A, 2)))*A'*y;
```

To verify, let us compare weights2 and theta

```
display(weights2);
```

```
weights2 = 3x1
0.0195
0.1231
0.1423
```

```
display(theta);
```

```
theta = 3x1
0.0195
0.1231
0.1423
```

Problem 4



$$9 \quad \min_{A,B} \frac{1}{2} \|P_2(X - AB^T)\|_F^2 + \frac{\lambda}{2} (\|A\|_F^2 + \|B\|_F^2) \quad - \textcircled{1}$$

$$10 \quad \text{Given: } A = U\Sigma^{1/2}; \quad B = V\Sigma^{1/2}$$

$$11 \quad [U, \Sigma, V] = \text{svd}(Z)$$

$$12 \quad P_2(X - AB^T) = P_2(X) - P_2(AB^T)$$

$$13 \quad = P_2(X) - P_2(Z) \quad - \textcircled{2}$$

$$14 \quad \|A\|_F^2 = \text{tr}(A^T A) = \text{tr}[(U\Sigma^{1/2})^T (U\Sigma^{1/2})] \\ = \text{tr}[\Sigma^{1/2} U^T U \Sigma^{1/2}]$$

$$15 \quad \text{As } U^T U = I,$$

$$16 \quad \|A\|_F^2 = \text{tr}[\Sigma^{1/2} \Sigma^{1/2}] = \text{tr}(\Sigma), \quad - \textcircled{3}$$

$$17 \quad \text{Similarly, } B^T B = (V\Sigma^{1/2})(V\Sigma^{1/2})^T$$

$$18 \quad \|B\|_F^2 = \text{tr}(B^T B) = \text{tr}(\Sigma), \quad - \textcircled{4}$$

Substituting these values from 2, 3 & 4.
in $\textcircled{1}$,

$$\min_{A, B} \frac{1}{2} \|P_n(X) - P_n(Z)\|_F^2 + \frac{\lambda}{2} (\text{tr}(\Sigma) + \text{tr}(\Sigma))$$

$$= \min_{A, B} \frac{1}{2} \|P_n(X) - P_n(Z)\|_F^2 + \lambda \cdot \text{tr}(\Sigma)$$

$\text{tr}(\Sigma)$ is the trace of Σ matrix, i.e.
sum of singular values

This is also = nuclear norm of Z .

So, above eqⁿ can be rewritten as:

$$\min_{A, B} \frac{1}{2} \|P_n(X) - P_n(Z)\|_F^2 + \lambda \|Z\|_*$$

Problem 4 (2)

```
% Given A, y, lambda, learning rate and initial weights
% Parameters
n = 2000;
d = 2000;
r = 200;

% Generate a random X matrix for testing
X = randn(n, d);

% Set the regularization parameter
lambda = 0.1;

% Initialize A and B
A = randn(n, r);
B = randn(d, r);

% Maximum number of iterations
T = 100;

num_iterations = 100;

% LEARNING RATE FORMULA TBD
lr = 1 / (max(max(A' * A)) + lambda);%1/(max_singular_value + lambda(k));
for i = 1:T
    A = A + lr * ((X - A*B')*B - A.*lambda);
    B = B + lr * ((X - A*B')*A - B.*lambda);
end
```