

Machine Learning: Programming Exercise 8

Anomaly Detection and Recommender Systems

In this exercise, you will implement the anomaly detection algorithm and apply it to detect failing servers on a network. In the second part, you will use collaborative filtering to build a recommender system for movies.

Files needed for this exercise

- `ex8.mlx` - MATLAB Live Script that steps you through the exercise
- `ex8data1.mat` - First example Dataset for anomaly detection
- `ex8data2.mat` - Second example Dataset for anomaly detection
- `ex8_movies.mat` - Movie Review Dataset
- `ex8_movieParams.mat` - Parameters provided for debugging
- `multivariateGaussian.m` - Computes the probability density function for a Gaussian distribution
- `visualizeFit.m` - 2D plot of a Gaussian distribution and a dataset
- `checkCostFunction.m` - Gradient checking for collaborative filtering
- `computeNumericalGradient.m` - Numerically compute gradients
- `fmincg.m` - Function minimization routine (similar to `fminunc`)
- `loadMovieList.m` - Loads the list of movies into a cell-array
- `movie_ids.txt` - List of movies
- `normalizeRatings.m` - Mean normalization for collaborative filtering
- `submit.m` - Submission script that sends your solutions to our servers
- `*estimateGaussian.m` - Estimate the parameters of a Gaussian distribution with a diagonal covariance matrix
- `*selectThreshold.m` - Find a threshold for anomaly detection
- `*cofiCostFunc.m` - Implement the cost function for collaborative filtering

**** indicates files you will need to complete***

Confirm that your Current Folder is set correctly

Click into this section, then click the 'Run section' button above. This will execute the `dir` command below to list the files in your Current Folder. The output should contain all of the files listed above and the 'lib' folder. If it does not, right-click the 'ex8' folder and select 'Open' before proceeding or see the instructions in `README.mlx` for more details.

```
dir
```

Before you begin

The workflow for completing and submitting the programming exercises in MATLAB Online differs from the original course instructions. Before beginning this exercise, make sure you have read through the instructions in `README.mlx` which is included with the programming exercise files. `README` also contains solutions to the many common issues you may encounter while completing and submitting the exercises in MATLAB Online. Make sure you are following instructions in `README` and have checked for an existing solution before seeking help on the discussion forums.

Table of Contents

Anomaly Detection and Recommender Systems..... 1

Files needed for this exercise..... 1

Confirm that your Current Folder is set correctly..... 1

Before you begin..... 1

1. Anomaly Detection..... 2

1.1 Gaussian distribution..... 4

1.2 Estimating parameters for a Gaussian..... 4

1.3 Selecting the threshold, 6

1.4 High dimensional dataset..... 9

2. Recommender Systems..... 10

2.1 Movie ratings dataset..... 10

2.2 Collaborative filtering learning algorithm..... 11

2.2.1 Collaborative filtering cost function..... 12

2.2.2 Collaborative filtering gradient..... 13

2.2.3 Regularized cost function..... 15

2.2.4 Regularized gradient..... 15

2.3 Learning movie recommendations..... 16

2.3.1 Recommendations..... 17

Submission and Grading..... 21

1. Anomaly Detection

In this exercise, you will implement an anomaly detection algorithm to detect anomalous behavior in server computers. The features measure the throughput (mb/s) and latency (ms) of response of each server. While your servers were operating, you collected $m = 307$ examples of how they were behaving, and thus have an unlabeled dataset $\{x^{(1)}, \dots, x^{(m)}\}$. You suspect that the vast majority of these examples are 'normal' (non-anomalous) examples of the servers operating normally, but there might also be some examples of servers acting anomalously within this dataset.

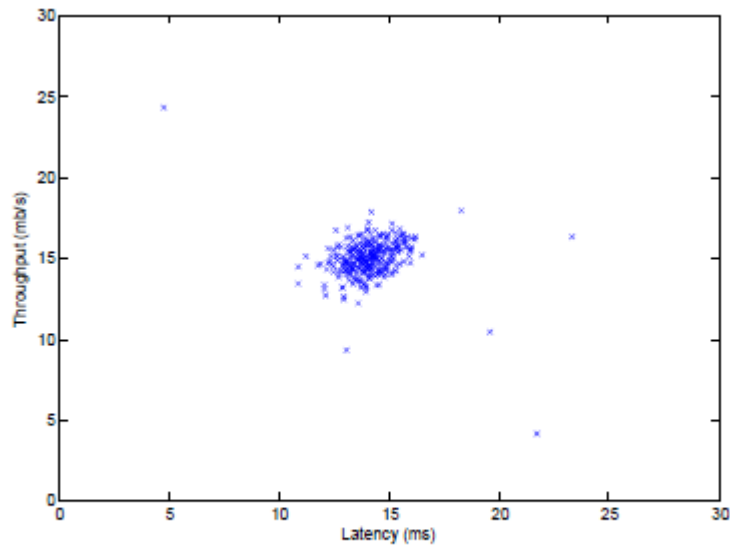
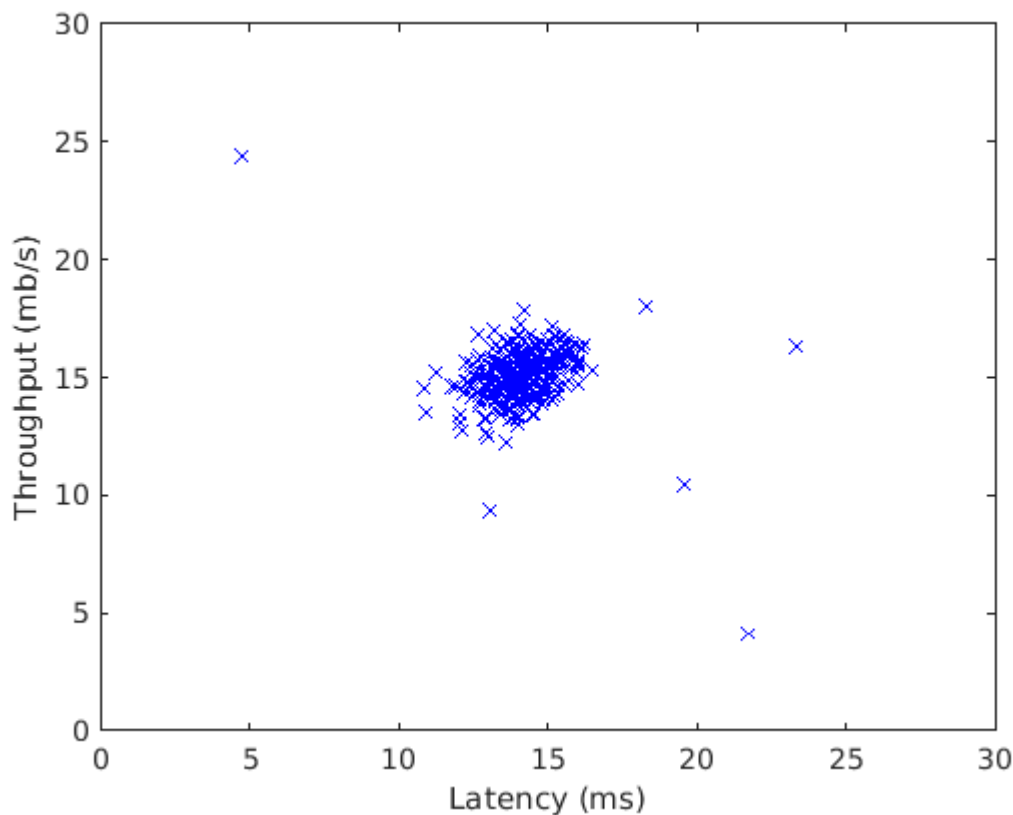


Figure 1: The first dataset.

You will use a Gaussian model to detect anomalous examples in your dataset. You will first start on a 2D dataset that will allow you to visualize what the algorithm is doing. On that dataset you will fit a Gaussian distribution and then find values that have very low probability and hence can be considered anomalies. After that, you will apply the anomaly detection algorithm to a larger dataset with many dimensions. The code below will visualize the dataset as shown in Figure 1.

```
% The following command loads the dataset. You should now have the variables X, Xval, y
load('ex8data1.mat');

% Visualize the example dataset
plot(X(:, 1), X(:, 2), 'bx');
axis([0 30 0 30]);
xlabel('Latency (ms)');
ylabel('Throughput (mb/s)');
```



1.1 Gaussian distribution

To perform anomaly detection, you will first need to fit a model to the data's distribution. Given a training set $\{x^{(1)}, \dots, x^{(m)}\}$ (where $x^{(i)} \in \mathbb{R}^n$), you want to estimate the Gaussian distribution for each of the features x_i . For each feature $i = 1 \dots n$ you need to find parameters μ_i and σ_i^2 that fit the data in the i -th dimension $\{x_i^{(1)}, \dots, x_i^{(m)}\}$ (the i -th dimension of each example).

The Gaussian distribution is given by

$$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

where μ is the mean and σ^2 controls the variance.

1.2 Estimating parameters for a Gaussian

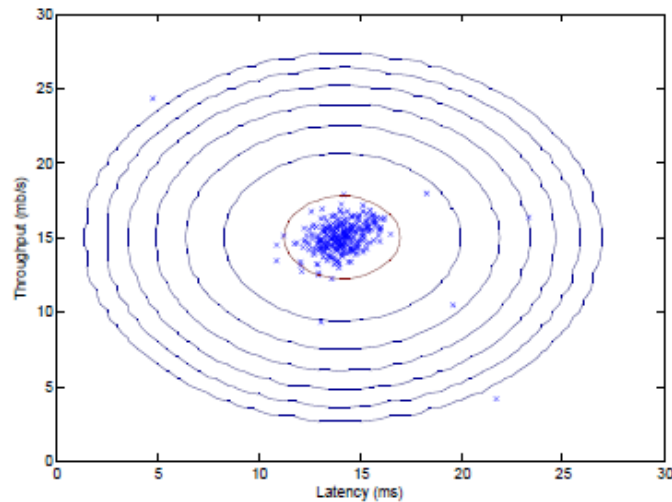
You can estimate the parameters, (μ_i, σ_i^2) , of the i -th feature by using the following equations. To estimate the mean, you will use:

$$\mu_i = \frac{1}{m} \sum_{j=1}^m x^{(j)},$$

and for the variance you will use:

$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x^{(j)} - \mu_j)^2.$$

Your task is to complete the code in `estimateGaussian.m`. This function takes as input the data matrix `X` and should output an n -dimension vector `mu` that holds the mean of all the n features and another n -dimension vector `sigma2` that holds the variances of all the features. You can implement this using a `for` loop over every feature and every training example (though a vectorized implementation might be more efficient; feel free to use a vectorized implementation if you prefer). Note that in MATLAB, the `var` function will (by default) use $\frac{1}{m-1}$, instead of $\frac{1}{m}$ when computing σ_i^2 .

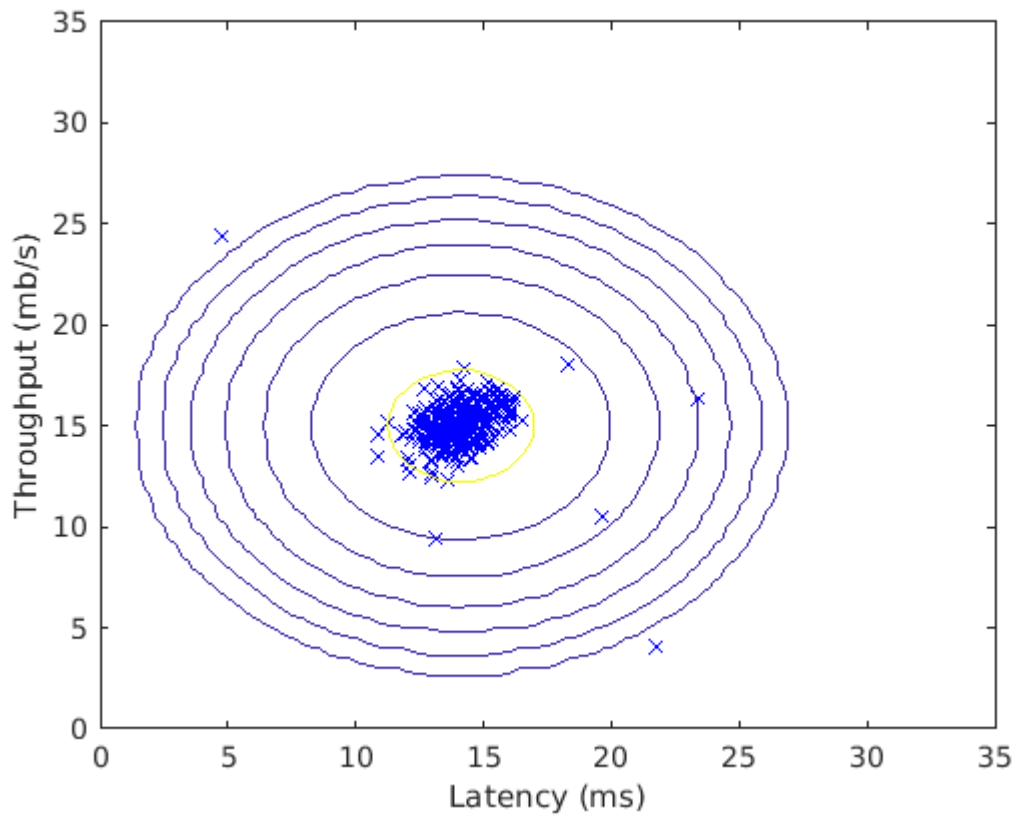


Once you have completed the code in `estimateGaussian.m`, the code below will visualize the contours of the fitted Gaussian distribution. You should get a plot similar to Figure 2. From your plot, you can see that most of the examples are in the region with the highest probability, while the anomalous examples are in the regions with lower probabilities.

```
% Estimate mu and sigma2
[mu, sigma2] = estimateGaussian(X);

% Returns the density of the multivariate normal at each data point (row) of X
p = multivariateGaussian(X, mu, sigma2);

% Visualize the fit
visualizeFit(X, mu, sigma2);
xlabel('Latency (ms)');
ylabel('Throughput (mb/s)');
```



*You should now submit your solutions. Enter **submit** at the command prompt, then and enter or confirm your login and token when prompted.*

1.3 Selecting the threshold, ϵ

Now that you have estimated the Gaussian parameters, you can investigate which examples have a very high probability given this distribution and which examples have a very low probability. The low probability examples are more likely to be the anomalies in our dataset. One way to determine which examples are anomalies is to select a threshold based on a cross validation set. In this part of the exercise, you will implement an algorithm to select the threshold ϵ using the F_1 score on a cross validation set.

You should now complete the code in `selectThreshold.m`. For this, we will use a cross validation set $\{(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m)}, y_{cv}^{(m)})\}$, where the label $y = 1$ corresponds to an anomalous example, and $y = 0$ corresponds to a normal example. For each cross validation example, we will compute $p(x_{cv}^{(i)})$. The vector of all of these probabilities $p(x_{cv}^{(1)}), \dots, p(x_{cv}^{(m)})$ is passed to `selectThreshold.m` in the vector `pval`. The corresponding labels $y_{cv}^{(1)}, \dots, y_{cv}^{(m)}$ is passed to the same function in the vector `yval`.

The function `selectThreshold.m` should return two values; the first is the selected threshold ϵ . If an example x has a low probability $p(x) < \epsilon$, then it is considered to be an anomaly. The function should also return the F_1 score, which tells you how well you're doing on finding the ground truth anomalies given a certain threshold.

For many different values of ϵ , you will compute the resulting F_1 score by computing how many examples the current threshold classifies correctly and incorrectly.

The F_1 score is computed using precision ($prec$) and recall (rec):

$$F_1 = \frac{2(prec)(rec)}{prec + rec},$$

You compute precision and recall by:

$$prec = \frac{tp}{tp + fp}$$

$$rec = \frac{tp}{tp + fn}$$

where

- tp is the number of true positives: the ground truth label says it's an anomaly and our algorithm correctly classified it as an anomaly.
- fp is the number of false positives: the ground truth label says it's not an anomaly, but our algorithm incorrectly classified it as an anomaly.
- fn is the number of false negatives: the ground truth label says it's an anomaly, but our algorithm incorrectly classified it as not being anomalous.

In the provided code `selectThreshold.m`, there is already a loop that will try many different values of ϵ and select the "best" based ϵ on the F_1 score. You should now complete the code in `selectThreshold.m`. You can implement the computation of the F_1 score using a `for` loop over all the cross validation examples (to compute the values tp , fp , and fn). You should see a value for epsilon of about $8.99e-05$.

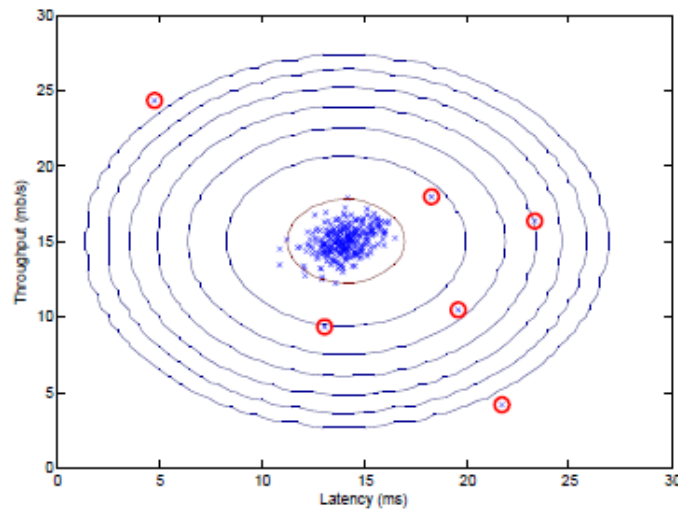


Figure 3: The classified anomalies.

Implementation Note: In order to compute tp , fp , and fn , you may be able to use a vectorized implementation rather than loop over all the examples. This can be implemented by MATLAB's equality test between a vector and a single number. If you have several binary values in an n -dimensional binary vector $v \in \{0, 1\}^n$, you can find out how many values in this vector are 0 by using: `sum(v == 0)`. You can also apply a logical and operator to such binary vectors. For instance, let `cvPredictions` be a binary vector of the size of your cross validation set, where the i -th element is 1 if your algorithm considers $x_{cv}^{(i)}$ an anomaly, and 0 otherwise. You can then, for example, compute the number of false positives using: `fp = sum((cvPredictions == 1) & (yval == 0))`.

Once you have completed the code in `selectThreshold.m`, the code below will detect and circle the anomalies in the plot (Figure 3).

```
pval = multivariateGaussian(Xval, mu, sigma2);

[epsilon, F1] = selectThreshold(yval, pval);
fprintf('Best epsilon found using cross-validation: %e\n', epsilon);
```

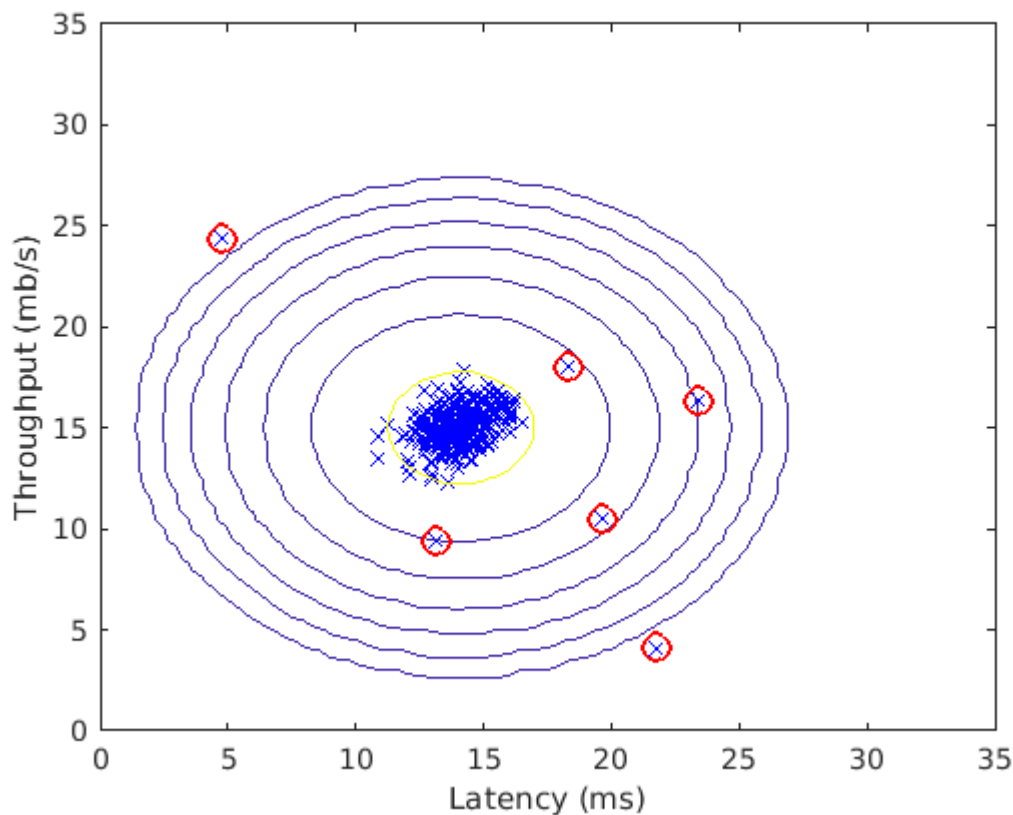
```
Best epsilon found using cross-validation: 8.990853e-05
```

```
fprintf('Best F1 on Cross Validation Set: %f\n', F1);
```

```
Best F1 on Cross Validation Set: 0.875000
```

```
% Find the outliers in the training set and plot the
outliers = find(p < epsilon);

% Visualize the fit
visualizeFit(X, mu, sigma2);
xlabel('Latency (ms)');
ylabel('Throughput (mb/s)');
% Draw a red circle around those outliers
hold on
plot(X(outliers, 1), X(outliers, 2), 'ro', 'LineWidth', 2, 'MarkerSize', 10);
hold off
```

*You should now submit your solutions. Enter **submit** at the command prompt, then and enter or confirm your login and token when prompted.*

1.4 High dimensional dataset

The code in this section will run the anomaly detection algorithm you implemented on a more realistic and much harder dataset. In this dataset, each example is described by 11 features, capturing many more properties of your compute servers. The code below will use your code to estimate the Gaussian parameters (μ_i and σ_i^2), evaluate the probabilities for both the training data x from which you estimated the Gaussian parameters, and do so for the the cross-validation set x_{val} . Finally, it will use `selectThreshold` to find the best threshold ϵ . You should see a value `epsilon` of about $1.38e-18$, and 117 anomalies found.

```
% Loads the second dataset. You should now have the variables X, Xval, yval in your environment
load('ex8data2.mat');

% Apply the same steps to the larger dataset
[mu, sigma2] = estimateGaussian(X);

% Training set
p = multivariateGaussian(X, mu, sigma2);

% Cross-validation set
pval = multivariateGaussian(Xval, mu, sigma2);
```

```
% Find the best threshold
[epsilon, F1] = selectThreshold(yval, pval);
fprintf('Best epsilon found using cross-validation: %e\n', epsilon);
```

```
Best epsilon found using cross-validation: 1.377229e-18
```

```
fprintf('Best F1 on Cross Validation Set: %f\n', F1);
```

```
Best F1 on Cross Validation Set: 0.615385
```

```
fprintf('# Outliers found: %d\n', sum(p < epsilon));
```

```
# Outliers found: 117
```

2. Recommender Systems

In this part of the exercise, you will implement the collaborative filtering learning algorithm and apply it to a dataset of movie ratings*. This dataset consists of ratings on a scale of 1 to 5. The dataset has $n_u = 943$ users, and $n_m = 1682$ movies. In the next parts of this exercise, you will implement the function `cofiCostFunc.m` that computes the collaborative filtering objective function and gradient. After implementing the cost function and gradient, you will use `fmincg.m` to learn the parameters for collaborative filtering.

*[MovieLens 100k Dataset](#) from GroupLens Research.

2.1 Movie ratings dataset

The code in this section will load the dataset `ex8_movies.mat`, providing the variables `Y` and `R` in your MATLAB environment. The matrix `Y` (a `num_movies` × `num_users` matrix) stores the ratings $y^{(i,j)}$ (from 1 to 5). The matrix `R` is a binary-valued indicator matrix, where $R(i, j) = 1$ if user j gave a rating to movie i , and $R(i, j) = 0$ otherwise. The objective of collaborative filtering is to predict movie ratings for the movies that users have not yet rated, that is, the entries with $R(i, j) = 0$. This will allow us to recommend the movies with the highest predicted ratings to the user.

To help you understand the matrix `Y`, the code below will compute the average movie rating for the first movie (Toy Story) and output the average rating to the screen. Throughout this part of the exercise, you will also be working with the matrices, `X` and `Theta`:

$$X = \begin{bmatrix} -(x^{(1)})^T & - \\ -(x^{(2)})^T & - \\ \vdots & \\ -(x^{(n_m)})^T & - \end{bmatrix}, \quad \text{Theta} = \begin{bmatrix} -(\theta^{(1)})^T & - \\ -(\theta^{(2)})^T & - \\ \vdots & \\ -(\theta^{(n_u)})^T & - \end{bmatrix}.$$

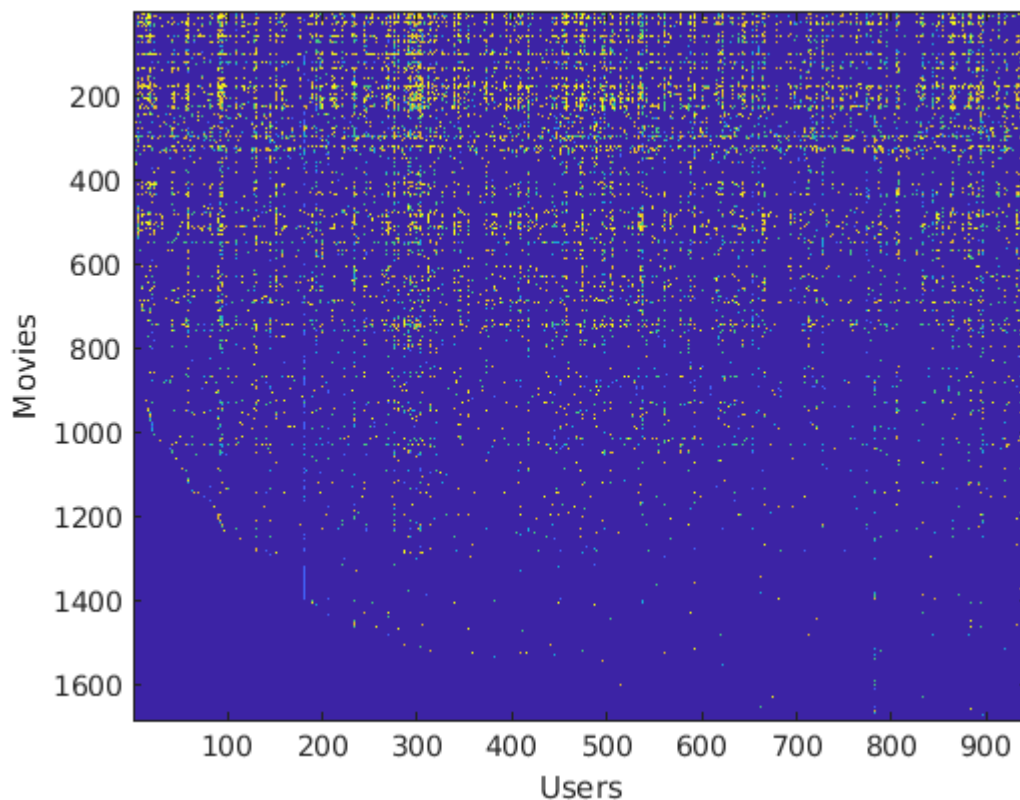
```
% Load data
load('ex8_movies.mat');
```

- `Y` is a 1682 x 943 matrix, containing ratings (1 - 5) of 1682 movies on 943 users
- `R` is a 1682 x 943 matrix, where $R(i, j) = 1$ if and only if user j gave a rating to movie i

```
% From the matrix, we can compute statistics like average rating.
fprintf('Average rating for movie 1 (Toy Story): %f / 5\n\n', mean(Y(1, R(1, :))));
```

```
Average rating for movie 1 (Toy Story): 3.878319 / 5
```

```
% We can "visualize" the ratings matrix by plotting it with imagesc
imagesc(Y);
ylabel('Movies');
xlabel('Users');
```



The i -th row of X corresponds to the feature vector $x^{(i)}$ for the i -th movie, and the j -th row of Θ corresponds to one parameter vector $\theta^{(j)}$, for the j -th user. Both $x^{(i)}$ and $\theta^{(j)}$ are n -dimensional vectors. For the purposes of this exercise, you will use $n = 100$, and therefore, $x^{(i)} \in \mathbb{R}^{100}$ and $\theta^{(j)} \in \mathbb{R}^{100}$. Correspondingly, X is a $n_m \times 100$ matrix and Θ is a $n_u \times 100$ matrix.

2.2 Collaborative filtering learning algorithm

Now, you will start implementing the collaborative filtering learning algorithm. You will start by implementing the cost function (without regularization). The collaborative filtering algorithm in the setting of movie recommendations considers a set of n -dimensional parameter vectors $x^{(1)}, \dots, x^{(n_m)}$ and $\theta^{(1)}, \dots, \theta^{(n_u)}$, where the model predicts the rating for movie i by user j as $y^{(i,j)} = (\theta^{(j)})^T x^{(i)}$. Given a dataset that consists

of a set of ratings produced by some users on some movies, you wish to learn the parameter vectors

$x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ that produce the best fit (minimizes the squared error).

You will complete the code in `cofiCostFunc.m` to compute the cost function and gradient for collaborative filtering. Note that the parameters to the function (i.e., the values that you are trying to learn) are `X` and `Theta`. In order to use an off-the-shelf minimizer such as `fmincg`, the cost function has been set up to unroll the parameters into a single vector `params`. You had previously used the same vector unrolling method in the neural networks programming exercise.

2.2.1 Collaborative filtering cost function

The collaborative filtering cost function (without regularization) is given by

$$J(x^{(i)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2$$

You should now modify `cofiCostFunc.m` to return this cost in the variable `J`. Note that you should be accumulating the cost for user j and movie i only if $R(i, j) = 1$. After you have completed the function, the code below will run your cost function. You should expect to see an output of 22.22.

Implementation Note: We strongly encourage you to use a vectorized implementation to compute J , since it will later be called many times by the optimization package `fmincg`. As usual, it might be easiest to first write a non-vectorized implementation (to make sure you have the right answer), and then modify it to become a vectorized implementation (checking that the vectorization steps don't change your algorithm's output). To come up with a vectorized implementation, the following tip might be helpful: You can use the R matrix to set selected entries to 0. For example, `R.*M` will do an element-wise multiplication between M and R ; since R only has elements with values either 0 or 1, this has the effect of setting the elements of M to 0 only when the corresponding value in R is 0. Hence, `sum(sum(R.*M))` is the sum of all the elements of M for which the corresponding element in R equals 1.

```
% Load pre-trained weights (X, Theta, num_users, num_movies, num_features)
load('ex8_movieParams.mat');

% Reduce the data set size so that this runs faster
num_users = 4; num_movies = 5; num_features = 3;
X = X(1:num_movies, 1:num_features);
Theta = Theta(1:num_users, 1:num_features);
Y = Y(1:num_movies, 1:num_users);
R = R(1:num_movies, 1:num_users);

% Evaluate cost function
J = cofiCostFunc([X(:); Theta(:)], Y, R, num_users, num_movies, num_features, 0);
fprintf('Cost at loaded parameters: %f ', J);
```

```
Cost at loaded parameters: 22.224604
```

You should now submit your solutions. Enter `submit` at the command prompt, then and enter or confirm your login and token when prompted.

2.2.2 Collaborative filtering gradient

Now, you should implement the gradient (without regularization). Specically, you should complete the code in `cofiCostFunc.m` to return the variables `X_grad` and `Theta_grad`. Note that `X_grad` should be a matrix of the same size as `X` and similarly, `Theta_grad` is a matrix of the same size as `Theta`. The gradients of the cost function is given by:

$$\frac{\partial J}{\partial x_k^{(i)}} = \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)}$$

$$\frac{\partial J}{\partial \theta_k^{(j)}} = \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)}$$

Note that the function returns the gradient for both sets of variables by unrolling them into a single vector. After you have completed the code to compute the gradients, the code below will run a gradient check (`checkCostFunction`) to numerically check the implementation of your gradients. (This is similar to the numerical check that you used in the neural networks exercise.) If your implementation is correct, you should find that the analytical and numerical gradients match up closely.

Implementation Note: You can get full credit for this assignment without using a vectorized implementation, but your code will run much more slowly (a small number of hours), and so we recommend that you try to vectorize your implementation. To get started, you can implement the gradient with a `for` loop over movies (for computing $\frac{\partial J}{\partial x_k^{(i)}}$) and a `for` loop over users (for computing $\frac{\partial J}{\partial \theta_k^{(j)}}$). When you first implement the gradient, you might start with an unvectorized version, by implementing another inner `for` loop that computes each element in the summation. After you have completed the gradient computation this way, you should try to vectorize your implementation (vectorize the inner `for` loops), so that you're left with only two `for` loops (one for looping over movies to compute $\frac{\partial J}{\partial x_k^{(i)}}$ for each movie, and one for looping over users to compute $\frac{\partial J}{\partial \theta_k^{(j)}}$ for each user).

To perform the vectorization, you might find this helpful: You should come up with a way to compute all the derivatives associated with $x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}$, (i.e., the derivative terms associated with the feature vector $x^{(i)}$) at the same time. Let us define the derivatives for the feature vector of the i -th movie as:

$$(X_{\text{grad}}(i, :))^T = \begin{bmatrix} \frac{\partial J}{\partial x_1^{(i)}} \\ \frac{\partial J}{\partial x_2^{(i)}} \\ \vdots \\ \frac{\partial J}{\partial x_n^{(i)}} \end{bmatrix} = \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta^{(j)}$$

To vectorize the above expression, you can start by indexing into `Theta` and `Y` to select only the elements of interest (that is, those with $r(i, j) = 1$). Intuitively, when you consider the features for the i -th movie, you only

need to be concerned about the users who had given ratings to the movie, and this allows you to remove all the other users from Θ and Y .

Concretely, you can set `idx = find(R(i,:) == 1)` to be a list of all the users that have rated movie i . This will allow you to create the temporary matrices `Theta_temp = Theta(idx,:)` and `Y_temp = Y(i,idx)` that index into Θ and Y to give you only the set of users which have rated the i -th movie. This will allow you to write the derivatives as:

$$X_{\text{grad}}(i,:) = (X(i,:) * \Theta_{\text{temp}}^T - Y_{\text{temp}}) * \Theta_{\text{temp}}.$$

(Note: The vectorized computation above returns a row-vector instead.) After you have vectorized the computations of the derivatives with respect to $x^{(i)}$, you should use a similar method to vectorize the derivatives with respect to $\theta^{(j)}$ as well.

```
% Check gradients by running checkNNGradients
checkCostFunction;
```

0.8700	0.8700
0	0
2.0479	2.0479
-2.7341	-2.7341
2.3462	2.3462
0	0
-3.1380	-3.1380
4.7809	4.7809
1.9846	1.9846
0	0
1.2999	1.2999
-0.9485	-0.9485
0.4700	0.4700
-0.2653	-0.2653
4.1396	4.1396
-0.7753	-0.7753
-0.9104	-0.9104
-0.9614	-0.9614
0.6241	0.6241
-5.6192	-5.6192
0.8609	0.8609
0.6332	0.6332
-0.5509	-0.5509
-0.3591	-0.3591
0.6803	0.6803
-0.6341	-0.6341
-0.5817	-0.5817

The above two columns you get should be very similar.
(Left-Your Numerical Gradient, Right-Analytical Gradient)

If your cost function implementation is correct, then
the relative difference will be small (less than $1e-9$).

Relative Difference: 1.21911e-12

*You should now submit your solutions. Enter **submit** at the command prompt, then and enter or confirm your login and token when prompted.*

2.2.3 Regularized cost function

The cost function for collaborative filtering with regularization is given by

$$J(x^{(i)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \left(\frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2 \right) + \left(\frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 \right)$$

You should now add regularization to your original computations of the cost function, J . After you are done, the code below will run your regularized cost function, and you should expect to see a cost of about 31.34.

```
% Evaluate cost function
J = cofiCostFunc([X(:); Theta(:)], Y, R, num_users, num_movies, num_features, 1.5);
fprintf('Cost at loaded parameters (lambda = 1.5): %f', J);
```

*You should now submit your solutions. Enter **submit** at the command prompt, then and enter or confirm your login and token when prompted.*

2.2.4 Regularized gradient

Now that you have implemented the regularized cost function, you should proceed to implement regularization for the gradient. You should add to your implementation in `cofiCostFunc.m` to return the regularized gradient by adding the contributions from the regularization terms. Note that the gradients for the regularized cost function is given by:

$$\frac{\partial J}{\partial x_k^{(i)}} = \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)}$$
$$\frac{\partial J}{\partial \theta_k^{(j)}} = \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)}$$

This means that you just need to add $\lambda x^{(i)}$ to the `X_grad(i, :)` variable described earlier, and add $\lambda \theta^{(j)}$ to the `Theta_grad(j, :)` variable described earlier. After you have completed the code to compute the gradients, the code below will run another gradient check (`checkCostFunction`) to numerically check the implementation of your gradients.

```
% Check gradients by running checkNNGradients
checkCostFunction(1.5);
```

```
-0.2471    -0.2471
 1.9414     1.9414
-1.9787    -1.9787
 0.7006     0.7006
 3.5101     3.5101
 0.2543     0.2543
 1.5416     1.5416
 2.3990     2.3990
-3.2561    -3.2561
-4.1617    -4.1617
 4.7051     4.7051
-1.1998    -1.1998
 0.3539     0.3539
 1.3291     1.3291
```

1.2457	1.2457
-1.5213	-1.5213
0.5858	0.5858
-1.5169	-1.5169
2.1835	2.1835
-3.5664	-3.5664
1.4779	1.4779
1.0621	1.0621
-1.8107	-1.8107
-4.9811	-4.9811
-2.8814	-2.8814
0.5047	0.5047
-4.5334	-4.5334

The above two columns you get should be very similar.
(Left-Your Numerical Gradient, Right-Analytical Gradient)

If your cost function implementation is correct, then
the relative difference will be small (less than $1e-9$).

Relative Difference: 2.28773e-12

*You should now submit your solutions. Enter **submit** at the command prompt, then and enter or confirm your login and token when prompted.*

2.3 Learning movie recommendations

After you have finished implementing the collaborative filtering cost function and gradient, you can now start training your algorithm to make movie recommendations for yourself. In the code below, you can enter your own movie preferences, so that later when the algorithm runs, you can get your own movie recommendations! We have filled out some values according to our own preferences, but you should change this according to your own tastes. The list of all movies and their number in the dataset can be found listed in the file movie idx.txt.

```
% Load movie list
movieList = loadMovieList();

% Initialize my ratings
my_ratings = zeros(1682, 1);

% Check the file movie_idx.txt for id of each movie in our dataset
% For example, Toy Story (1995) has ID 1, so to rate it "4", you can set
my_ratings(1) = 4;
% Or suppose did not enjoy Silence of the Lambs (1991), you can set
my_ratings(98) = 2;

% We have selected a few movies we liked / did not like and the ratings we gave are as
my_ratings(7) = 3;
my_ratings(12)= 5;
my_ratings(54) = 4;
my_ratings(64)= 5;
my_ratings(66)= 3;
my_ratings(69) = 5;
my_ratings(183) = 4;
my_ratings(226) = 5;
my_ratings(355)= 5;
```



```
fprintf('\n\nNew user ratings:\n');
```

New user ratings:

```
for i = 1:length(my_ratings)
    if my_ratings(i) > 0
        fprintf('Rated %d for %s\n', my_ratings(i), movieList{i});
    end
end
```

```
Rated 4 for Toy Story (1995)
Rated 3 for Twelve Monkeys (1995)
Rated 5 for Usual Suspects, The (1995)
Rated 4 for Outbreak (1995)
Rated 5 for Shawshank Redemption, The (1994)
Rated 3 for While You Were Sleeping (1995)
Rated 5 for Forrest Gump (1994)
Rated 2 for Silence of the Lambs, The (1991)
Rated 4 for Alien (1979)
Rated 5 for Die Hard 2 (1990)
Rated 5 for Sphere (1998)
```

2.3.1 Recommendations

After the additional ratings have been added to the dataset, the code below will proceed to train the collaborative filtering model. This will learn the parameters X and Θ .

```
% Load data
load('ex8_movies.mat');

% Y is a 1682x943 matrix, containing ratings (1-5) of 1682 movies by 943 users
% R is a 1682x943 matrix, where R(i,j) = 1 if and only if user j gave a rating to movie i
% Add our own ratings to the data matrix
Y = [my_ratings Y];
R = [(my_ratings ~= 0) R];

% Normalize Ratings
[Ynorm, Ymean] = normalizeRatings(Y, R);

% Useful Values
num_users = size(Y, 2);
num_movies = size(Y, 1);
num_features = 10;

% Set Initial Parameters (Theta, X)
X = randn(num_movies, num_features);
Theta = randn(num_users, num_features);
initial_parameters = [X(:); Theta(:)];

% Set options for fmincg
options = optimset('GradObj', 'on', 'MaxIter', 100);

% Set Regularization
lambda = 10;
```

```
theta = fmincg(@(t)(cofiCostFunc(t, Ynorm, R, num_users, num_movies, num_features, lambda
```

Iteration	1	Cost: 3.147179e+05
Iteration	2	Cost: 1.258353e+05
Iteration	3	Cost: 1.003162e+05
Iteration	4	Cost: 8.098216e+04
Iteration	5	Cost: 6.174263e+04
Iteration	6	Cost: 4.990692e+04
Iteration	7	Cost: 4.574697e+04
Iteration	8	Cost: 4.333935e+04
Iteration	9	Cost: 4.264168e+04
Iteration	10	Cost: 4.188885e+04
Iteration	11	Cost: 4.120568e+04
Iteration	12	Cost: 4.082442e+04
Iteration	13	Cost: 4.067975e+04
Iteration	14	Cost: 4.031180e+04
Iteration	15	Cost: 4.009608e+04
Iteration	16	Cost: 3.991522e+04
Iteration	17	Cost: 3.976256e+04
Iteration	18	Cost: 3.960765e+04
Iteration	19	Cost: 3.951726e+04
Iteration	20	Cost: 3.945804e+04
Iteration	21	Cost: 3.941524e+04
Iteration	22	Cost: 3.935598e+04
Iteration	23	Cost: 3.930324e+04
Iteration	24	Cost: 3.928257e+04
Iteration	25	Cost: 3.925262e+04
Iteration	26	Cost: 3.924258e+04
Iteration	27	Cost: 3.920573e+04
Iteration	28	Cost: 3.917665e+04
Iteration	29	Cost: 3.913786e+04
Iteration	30	Cost: 3.911741e+04
Iteration	31	Cost: 3.910725e+04
Iteration	32	Cost: 3.909607e+04
Iteration	33	Cost: 3.908635e+04
Iteration	34	Cost: 3.907298e+04
Iteration	35	Cost: 3.906507e+04
Iteration	36	Cost: 3.906113e+04
Iteration	37	Cost: 3.905604e+04
Iteration	38	Cost: 3.905421e+04
Iteration	39	Cost: 3.905268e+04
Iteration	40	Cost: 3.904750e+04
Iteration	41	Cost: 3.904235e+04
Iteration	42	Cost: 3.903743e+04
Iteration	43	Cost: 3.903417e+04
Iteration	44	Cost: 3.903007e+04
Iteration	45	Cost: 3.902828e+04
Iteration	46	Cost: 3.902688e+04
Iteration	47	Cost: 3.902591e+04
Iteration	48	Cost: 3.902381e+04
Iteration	49	Cost: 3.902131e+04
Iteration	50	Cost: 3.901993e+04
Iteration	51	Cost: 3.901909e+04
Iteration	52	Cost: 3.901821e+04
Iteration	53	Cost: 3.901748e+04
Iteration	54	Cost: 3.901690e+04
Iteration	55	Cost: 3.901566e+04
Iteration	56	Cost: 3.901368e+04
Iteration	57	Cost: 3.901262e+04
Iteration	58	Cost: 3.901214e+04
Iteration	59	Cost: 3.901016e+04
Iteration	60	Cost: 3.900879e+04
Iteration	61	Cost: 3.900806e+04
Iteration	62	Cost: 3.900702e+04
Iteration	63	Cost: 3.900597e+04

```

Iteration 64 | Cost: 3.900493e+04
Iteration 65 | Cost: 3.900187e+04
Iteration 66 | Cost: 3.899706e+04
Iteration 67 | Cost: 3.899633e+04
Iteration 68 | Cost: 3.899439e+04
Iteration 69 | Cost: 3.899294e+04
Iteration 70 | Cost: 3.899236e+04
Iteration 71 | Cost: 3.899176e+04
Iteration 72 | Cost: 3.899129e+04
Iteration 73 | Cost: 3.899046e+04
Iteration 74 | Cost: 3.898975e+04
Iteration 75 | Cost: 3.898909e+04
Iteration 76 | Cost: 3.898877e+04
Iteration 77 | Cost: 3.898865e+04
Iteration 78 | Cost: 3.898782e+04
Iteration 79 | Cost: 3.898689e+04
Iteration 80 | Cost: 3.898655e+04
Iteration 81 | Cost: 3.898582e+04
Iteration 82 | Cost: 3.898414e+04
Iteration 83 | Cost: 3.898199e+04
Iteration 84 | Cost: 3.897960e+04
Iteration 85 | Cost: 3.897868e+04
Iteration 86 | Cost: 3.897771e+04
Iteration 87 | Cost: 3.897719e+04
Iteration 88 | Cost: 3.897651e+04
Iteration 89 | Cost: 3.897624e+04
Iteration 90 | Cost: 3.897616e+04
Iteration 91 | Cost: 3.897539e+04
Iteration 92 | Cost: 3.897517e+04
Iteration 93 | Cost: 3.897486e+04
Iteration 94 | Cost: 3.897454e+04
Iteration 95 | Cost: 3.897429e+04
Iteration 96 | Cost: 3.897341e+04
Iteration 97 | Cost: 3.897182e+04
Iteration 98 | Cost: 3.897043e+04
Iteration 99 | Cost: 3.896884e+04
Iteration 100 | Cost: 3.896780e+04

```

```

% Unfold the returned theta back into U and W
X = reshape(theta(1:num_movies*num_features), num_movies, num_features);
Theta = reshape(theta(num_movies*num_features+1:end), num_users, num_features);

```

To predict the rating of movie i for user j , you need to compute $(\theta^{(j)})^T x^{(i)}$. The code below computes the ratings for all the movies and users and displays the movies that it recommends (Figure 4), according to ratings that were entered earlier in the script. Note that you might obtain a different set of the predictions due to different random initializations.

```

Top recommendations for you:
Predicting rating 9.0 for movie Titanic (1997)
Predicting rating 8.9 for movie Star Wars (1977)
Predicting rating 8.8 for movie Shawshank Redemption, The (1994)
Predicting rating 8.5 for movie As Good As It Gets (1997)
Predicting rating 8.5 for movie Good Will Hunting (1997)
Predicting rating 8.5 for movie Usual Suspects, The (1995)
Predicting rating 8.5 for movie Schindler's List (1993)
Predicting rating 8.4 for movie Raiders of the Lost Ark (1981)
Predicting rating 8.4 for movie Empire Strikes Back, The (1980)
Predicting rating 8.4 for movie Braveheart (1995)

Original ratings provided:
Rated 4 for Toy Story (1995)
Rated 3 for Twelve Monkeys (1995)
Rated 5 for Usual Suspects, The (1995)
Rated 4 for Outbreak (1995)
Rated 5 for Shawshank Redemption, The (1994)
Rated 3 for While You Were Sleeping (1995)
Rated 5 for Forrest Gump (1994)
Rated 2 for Silence of the Lambs, The (1991)
Rated 4 for Alien (1979)
Rated 5 for Die Hard 2 (1990)
Rated 5 for Sphere (1998)

```

Figure 4: Movie recommendations

```

p = X * Theta';
my_predictions = p(:,1) + Ymean;

movieList = loadMovieList();

[r, ix] = sort(my_predictions, 'descend');
for i=1:10
    j = ix(i);
    if i == 1
        fprintf('\nTop recommendations for you:\n');
    end
    fprintf('Predicting rating %.1f for movie %s\n', my_predictions(j), movieList{j});
end

```

```

Top recommendations for you:
Predicting rating 5.0 for movie Someone Else's America (1995)
Predicting rating 5.0 for movie Entertaining Angels: The Dorothy Day Story (1996)
Predicting rating 5.0 for movie They Made Me a Criminal (1939)
Predicting rating 5.0 for movie Saint of Fort Washington, The (1993)
Predicting rating 5.0 for movie Prefontaine (1997)
Predicting rating 5.0 for movie Marlene Dietrich: Shadow and Light (1996)
Predicting rating 5.0 for movie Aiqing wansui (1994)
Predicting rating 5.0 for movie Great Day in Harlem, A (1994)
Predicting rating 5.0 for movie Santa with Muscles (1996)
Predicting rating 5.0 for movie Star Kid (1997)

```

```

for i = 1:length(my_ratings)
    if i == 1
        fprintf('\n\nOriginal ratings provided:\n');
    end
    if my_ratings(i) > 0

```

```

        fprintf('Rated %d for %s\n', my_ratings(i), movieList{i});
    end
end

```

```

Original ratings provided:
Rated 4 for Toy Story (1995)
Rated 3 for Twelve Monkeys (1995)
Rated 5 for Usual Suspects, The (1995)
Rated 4 for Outbreak (1995)
Rated 5 for Shawshank Redemption, The (1994)
Rated 3 for While You Were Sleeping (1995)
Rated 5 for Forrest Gump (1994)
Rated 2 for Silence of the Lambs, The (1991)
Rated 4 for Alien (1979)
Rated 5 for Die Hard 2 (1990)
Rated 5 for Sphere (1998)

```

Submission and Grading

After completing various parts of the assignment, be sure to use the submit function system to submit your solutions to our servers. The following is a breakdown of how each part of this exercise is scored.

Part	Submitted File	Points
Estimate Gaussian Parameters	estimateGuassian.m	15 points
Select Threshold	selectThreshold.m	15 points
Collaborative Filtering Cost	cofiCostFunc.m	20 points
Collaborative Filtering Gradient	cofiCostFunc.m	30 points
Regularized Cost	cofiCostFunc.m	10 points
Gradient with regularization	cofiCostFunc.m	10 points
Total Points		100 points

You are allowed to submit your solutions multiple times, and we will take only the highest score into consideration.