

# K-means Clustering and Anomaly Detection

Adapted from the exercises of Andrew Ng. and Alex Ihler

## Introduction

In this exercise, you will implement the  $K$ -means clustering algorithm and apply it to compress an image. In the second part, you will learn about Anomaly Detection.

## 1 $K$ -means Clustering

In this exercise, you will implement the  $K$ -means algorithm and use it for image compression. You will first start on an example 2D dataset that will help you gain an intuition of how the  $K$ -means algorithm works. After that, you will use the  $K$ -means algorithm for image compression by reducing the number of colors that occur in an image to only those that are most common in that image. You will be using `ex7.m` for this part of the exercise.

## 1.1 Implementing $K$ -means

The  $K$ -means algorithm is a method to automatically cluster similar data examples together. Concretely, you are given a training set  $\{x^{(1)}, \dots, x^{(m)}\}$  (where  $x^{(i)} \in \mathbb{R}^n$ ), and want to group the data into a few cohesive “clusters”. The intuition behind  $K$ -means is an iterative procedure that starts by guessing the initial centroids, and then refines this guess by repeatedly assigning examples to their closest centroids and then recomputing the centroids based on the assignments.

The  $K$ -means algorithm is as follows:

```
% Initialize centroids
centroids = kMeansInitCentroids(X, K);
for iter = 1:iterations
    % Cluster assignment step: Assign each data point to the
    % closest centroid. idx(i) corresponds to  $\hat{c}^{(i)}$ , the index
    % of the centroid assigned to example i
    idx = findClosestCentroids(X, centroids);

    % Move centroid step: Compute means based on centroid
    % assignments
    centroids = computeMeans(X, idx, K);
end
```

The inner-loop of the algorithm repeatedly carries out two steps: (i) Assigning each training example  $x^{(i)}$  to its closest centroid, and (ii) Recomputing the mean of each centroid using the points assigned to it. The  $K$ -means algorithm will always converge to some final set of means for the centroids. Note that the converged solution may not always be ideal and depends on the initial setting of the centroids. Therefore, in practice the  $K$ -means algorithm is usually run a few times with different random initializations. One way to choose between these different solutions from different random initializations is to choose the one with the lowest cost function value (distortion).

You will implement the two phases of the  $K$ -means algorithm separately in the next sections.

### 1.1.1 Finding closest centroids

In the “cluster assignment” phase of the  $K$ -means algorithm, the algorithm assigns every training example  $x^{(i)}$  to its closest centroid, given the current positions of centroids. Specifically, for every example  $i$  we set

$$c^{(i)} := j \quad \text{that minimizes} \quad \|x^{(i)} - \mu_j\|^2,$$

where  $c^{(i)}$  is the index of the centroid that is closest to  $x^{(i)}$ , and  $\mu_j$  is the position (value) of the  $j$ 'th centroid. Note that  $c^{(i)}$  corresponds to `idx(i)` in the starter code.

Your task is to complete the code in `findClosestCentroids.m`. This function takes the data matrix `X` and the locations of all centroids inside `centroids` and should output a one-dimensional array `idx` that holds the index (a value in  $\{1, \dots, K\}$ , where  $K$  is total number of centroids) of the closest centroid to every training example.

You can implement this using a loop over every training example and every centroid.

Once you have completed the code in `findClosestCentroids.m`, the script `ex7.m` will run your code and you should see the output `[1 3 2]` corresponding to the centroid assignments for the first 3 examples.

### 1.1.2 Computing centroid means

Given assignments of every point to a centroid, the second phase of the algorithm recomputes, for each centroid, the mean of the points that were assigned to it. Specifically, for every centroid  $k$  we set

$$\mu_k := \frac{1}{|C_k|} \sum_{i \in C_k} x^{(i)}$$

where  $C_k$  is the set of examples that are assigned to centroid  $k$ . Concretely, if two examples say  $x^{(3)}$  and  $x^{(5)}$  are assigned to centroid  $k = 2$ , then you should update  $\mu_2 = \frac{1}{2}(x^{(3)} + x^{(5)})$ .

You should now complete the code in `computeCentroids.m`. You can implement this function using a loop over the centroids. You can also use a loop over the examples; but if you can use a vectorized implementation that does not use such a loop, your code may run faster.

Once you have completed the code in `computeCentroids.m`, the script `ex7.m` will run your code and output the centroids after the first step of  $K$ -means.

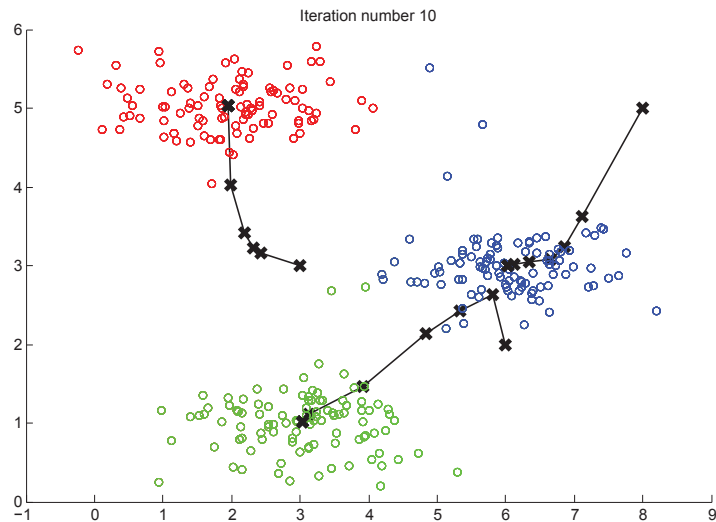


Figure 1: The expected output.

## 1.2 $K$ -means on example dataset

After you have completed the two functions (`findClosestCentroids` and `computeCentroids`), the next step in `ex7.m` will run the  $K$ -means algorithm on a toy 2D dataset to help you understand how  $K$ -means works. Your functions are called from inside the `runKmeans.m` script. We encourage you to take a look at the function to understand how it works. Notice that the code calls the two functions you implemented in a loop.

When you run the next step, the  $K$ -means code will produce a visualization that steps you through the progress of the algorithm at each iteration. Press *enter* multiple times to see how each step of the  $K$ -means algorithm changes the centroids and cluster assignments. At the end, your figure should look as the one displayed in Figure 1.

## 1.3 Random initialization

The initial assignments of centroids for the example dataset in `ex7.m` were designed so that you will see the same figure as in Figure 1. In practice, a good strategy for initializing the centroids is to select random examples from the training set.

In this part of the exercise, you should complete the function `kMeansInitCentroids.m` with the following code:

```
% Initialize the centroids to be random examples

% Randomly reorder the indices of examples
randidx = randperm(size(X, 1));
% Take the first K examples as centroids
centroids = X(randidx(1:K), :);
```

The code above first randomly permutes the indices of the examples (using `randperm`). Then, it selects the first  $K$  examples based on the random permutation of the indices. This allows the examples to be selected at random without the risk of selecting the same example twice.

## 1.4 Image compression with $K$ -means

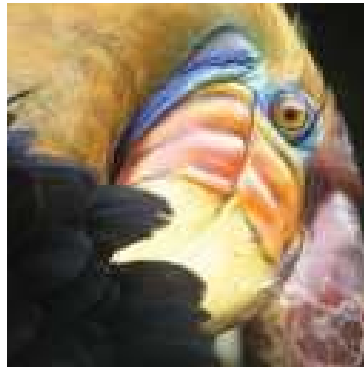


Figure 2: The original 128x128 image.

In this exercise, you will apply  $K$ -means to image compression. In a straightforward 24-bit color representation of an image,<sup>1</sup> each pixel is represented as three 8-bit unsigned integers (ranging from 0 to 255) that specify the red, green and blue intensity values. This encoding is often referred to as the RGB encoding. Our image contains thousands of colors, and in this part of the exercise, you will reduce the number of colors to 16 colors.

By making this reduction, it is possible to represent (compress) the photo in an efficient way. Specifically, you only need to store the RGB values of

---

<sup>1</sup>The provided photo used in this exercise belongs to Frank Wouters and is used with his permission.

the 16 selected colors, and for each pixel in the image you now need to only store the index of the color at that location (where only 4 bits are necessary to represent 16 possibilities).

In this exercise, you will use the  $K$ -means algorithm to select the 16 colors that will be used to represent the compressed image. Concretely, you will treat every pixel in the original image as a data example and use the  $K$ -means algorithm to find the 16 colors that best group (cluster) the pixels in the 3-dimensional RGB space. Once you have computed the cluster centroids on the image, you will then use the 16 colors to replace the pixels in the original image.

### 1.4.1 $K$ -means on pixels

In Matlab and Octave, images can be read in as follows:

```
% Load 128x128 color image (bird_small.png)
A = imread('bird_small.png');

% You will need to have installed the image package to used
% imread. If you do not have the image package installed, you
% should instead change the following line to
%
%   load('bird_small.mat'); % Loads the image into the variable A
```

This creates a three-dimensional matrix  $A$  whose first two indices identify a pixel position and whose last index represents red, green, or blue. For example,  $A(50, 33, 3)$  gives the blue intensity of the pixel at row 50 and column 33.

The code inside `ex7.m` first loads the image, and then reshapes it to create an  $m \times 3$  matrix of pixel colors (where  $m = 16384 = 128 \times 128$ ), and calls your  $K$ -means function on it.

After finding the top  $K = 16$  colors to represent the image, you can now assign each pixel position to its closest centroid using the `findClosestCentroids` function. This allows you to represent the original image using the centroid assignments of each pixel. Notice that you have significantly reduced the number of bits that are required to describe the image. The original image required 24 bits for each one of the  $128 \times 128$  pixel locations, resulting in total size of  $128 \times 128 \times 24 = 393,216$  bits. The new representation requires some overhead storage in form of a dictionary of 16 colors, each of which require 24 bits, but the image itself then only requires 4 bits per pixel location. The final number of bits used is therefore  $16 \times 24 + 128 \times 128 \times 4 = 65,920$  bits, which corresponds to compressing the original image by about a factor of 6.

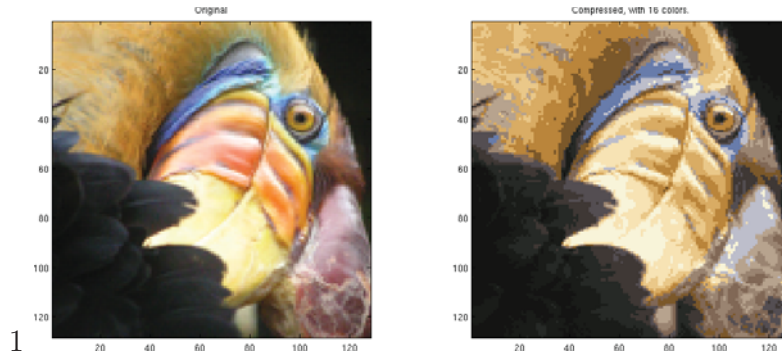


Figure 3: Original and reconstructed image (when using  $K$ -means to compress the image).

Finally, you can view the effects of the compression by reconstructing the image based only on the centroid assignments. Specifically, you can replace each pixel location with the mean of the centroid assigned to it. Figure 3 shows the reconstruction we obtained. Even though the resulting image retains most of the characteristics of the original, we also see some compression artifacts.

# 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.

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. You will be using `ex8.m` for this part of the exercise.

The first part of `ex8.m` will visualize the dataset as shown in Figure 1.

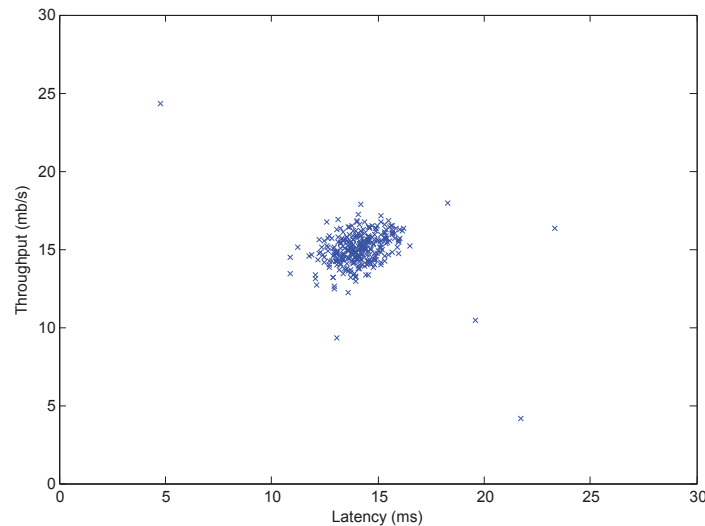


Figure 1: The first dataset.

## 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  $\mu_i$  and  $\sigma_i^2$  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  $\mu$  is the mean and  $\sigma^2$  controls the variance.

## 1.2 Estimating parameters for a Gaussian

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

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

and for the variance you will use:

$$\sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2. \quad (2)$$

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 Octave, the `var` function will (by default) use  $\frac{1}{m-1}$ , instead of  $\frac{1}{m}$ , when computing  $\sigma_i^2$ .

Once you have completed the code in `estimateGaussian.m`, the next part of `ex8.m` 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.

## 1.3 Selecting the threshold, $\varepsilon$

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

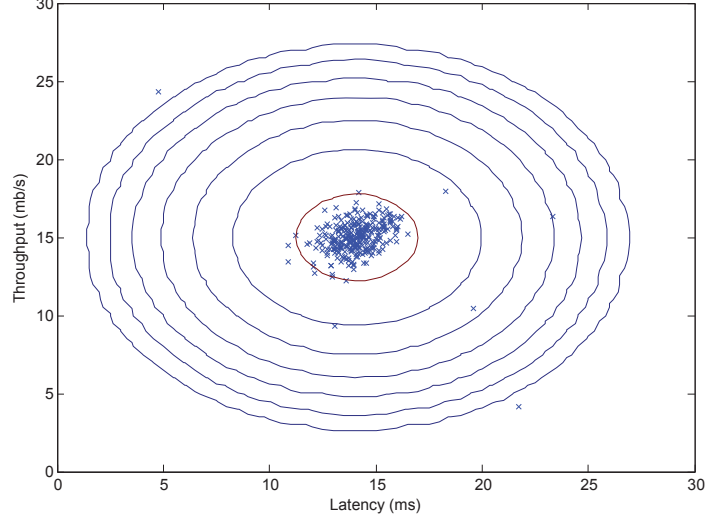


Figure 2: The Gaussian distribution contours of the distribution fit to the dataset.

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  $\varepsilon$  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  $\{(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})\}$ , 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  $p(x_{cv}^{(1)}), \dots, p(x_{cv}^{(m_{cv})})$  is passed to `selectThreshold.m` in the vector `pval`. The corresponding labels  $y_{cv}^{(1)}, \dots, y_{cv}^{(m_{cv})}$  is passed to the same function in the vector `yval`.

The function `selectThreshold.m` should return two values; the first is the selected threshold  $\varepsilon$ . If an example  $x$  has a low probability  $p(x) < \varepsilon$ , 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  $\varepsilon$ , 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*):

$$F_1 = \frac{2 \cdot \text{prec} \cdot \text{rec}}{\text{prec} + \text{rec}}, \quad (3)$$

You compute precision and recall by:

$$prec = \frac{tp}{tp + fp} \quad (4)$$

$$rec = \frac{tp}{tp + fn}, \quad (5)$$

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  $\varepsilon$  and select the best  $\varepsilon$  based 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  $tp$ ,  $fp$ ,  $fn$ ). You should see a value for epsilon of about 8.99e-05.

**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 Octave'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 number of 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 next step in `ex8.m` will run your anomaly detection code and circle the anomalies in the plot (Figure 3).

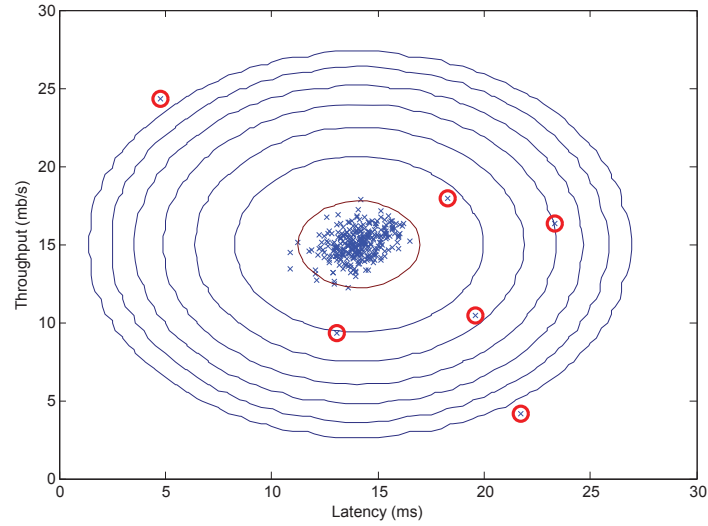


Figure 3: The classified anomalies.

## 1.4 High dimensional dataset

The last part of the script `ex8.m` 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 script will use your code to estimate the Gaussian parameters ( $\mu_i$  and  $\sigma_i^2$ ), 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 `Xval`. Finally, it will use `selectThreshold` to find the best threshold  $\epsilon$ . You should see a value epsilon of about  $1.38\text{e-}18$ , and 117 anomalies found.