

# Machine Learning: Assessed Coursework 2

Yadu Bhageria

## 1 Question 1: Image Segmentation and Counting

The image provided is a  $640 \times 640$  pixel image thus there are  $N = 409600$  data points,  $\mathbf{P}$ . The data for each pixel,  $p_i$ , is stored in the Red-Green-Blue (RGB) colour format as  $p_i = [R_i, G_i, B_i]$ . The image is segmented by clustering pixels into a predefined number of groups in the 3 dimensional RGB space. Two method have been considered: the K-means algorithm and then the more general Gaussian mixture model (GMM).

### 1.1 K-Means

The basic premise of the K-means algorithm is to partition the data points into a given number of clusters,  $K$ , to minimise the within cluster distance from the cluster centre/mean. For cluster  $k$  this can be defined as

$$\sum_{p_i \in C_k} ||p_i - \mu_k|| = \sum_{i=1}^N z_{ki} ||p_i - \mu_k|| \quad (1)$$

$z_{ki} \in \{0, 1\}$  is a binary indicator that assigns each data point,  $p_i$ , to a single cluster,  $k$ .  $\mu_k$  is the cluster mean,  $\frac{1}{N_k} \sum_{p_i \in C_k} p_i$  i.e. the average values of the data points in the cluster.

The total measure of the K-means algorithm can then be written as

$$\mathcal{E}_K = \sum_{k=1}^K \sum_{i=1}^N z_{ki} ||p_i - \mu_k|| \quad (2)$$

which is the overall cluster goodness. The K-means algorithm minimizes this by a 2 step algorithm. Given a segmentation of the data, the cluster centres can be found using the formula above. And given cluster centres,  $\{\mu_k\}$ , it can be found which data points belong to that cluster, i.e. have the smallest distance function to a cluster. This updates the indicator variables  $z_{ki}$ . This process can be repeated until the cluster centres stop moving and the solution converges.

The MATLAB code for 1 iteration is as follows:

```
1 Mu = Mu_new; % Set Mu to previously computed Cluster Center
2
3 dist = zeros(N,K); % Compute L2 norm
4 for k = 1:K
5     dist(:,k) = sum((X(:, :) - Mu(:,k)).^2, 1);
6 end
7 [~,Z] = min(dist,[],2); % Get minimum for each N
8
9 for k = 1:K % Set new cluster centers to mean of their elements.
10     if isempty(X(:,Z == k)) == 1 % Incase a cluster has no members move it to a random point
11         Mu_new(:,k) = rand(M,1);
12     else
13         Mu_new(:,k) = mean(X(:,Z == k),2);
14     end
15 end
```

Thus for the pixels in the RGB space, each pixel intensity vector can be replaced by cluster mean value corresponding to the cluster of that pixel. Furthermore for storage purposes it means that each picture can be stored using only  $K$  cluster mean vectors and a single integer value assigning the data points to one of the cluster means.

## 1.2 Gaussian Mixture Models

A Gaussian mixture model differs from the K-means algorithm in that it assigns a distribution over the indicator variables rather than only assigning binary indicator variables for each cluster-data point pair. This means the model is more complex but allows better segmentation of the image. This is because the K-means algorithm uses the L2 norm and thus is biased towards spherical distributions whereas GMM can have gaussians with covariance between dimensions.

A GMM tries to fit  $K$  gaussians to the data. That means the parameters  $\mu_k \in \mathbb{R}^3$  and  $\Sigma_k \in \mathbb{R}^{3 \times 3}$  for each cluster must be found. Let  $\theta$  be the parameters of the entire model. Maximising the the log likelihood of the data given the parameters is a method to do so. In lectures the lower bound on the log likelihood was derived as

$$\mathcal{L}_B = \log p(\mathbf{P}|\theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{P}) \log p(\mathbf{P}, \mathbf{Z}|\theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{P}) \log p(\mathbf{Z}|\mathbf{P}) \quad (3)$$

where  $\mathbf{Z}$  is the matrix of indicator variables. This can be maximized using the Expectation-Minimization algorithm. Given estimates for the model parameters, the probability density of the latent variables,  $p(\mathbf{Z}|\mathbf{P})$ , can be calculated. And given the probabilities of the latent variables, the parameters values,  $\theta$ , can be estimated by maximizing the log likelihood. These two steps are repeated until convergence.

So the algorithm is set up to start at random initial conditions and then first calculate the posterior distribution over the mixture components

$$p(k|p_i) = \frac{p(p_i|\theta_k)p(k)}{\sum_{k'=1}^K p(p_i|\theta_{k'})p(k)} \quad (4)$$

where  $k$  is a cluster. Then estimate the parameters,  $\theta$ ,

$$\hat{\mu}_k = \frac{\sum_{i=1}^N p(p_i|\theta_k)p_i}{\sum_{j=1}^N p(p_j|\theta_k)} \quad \hat{\Sigma}_k = \frac{\sum_{i=1}^N p(p_i|\theta_k)(p_i - \hat{\mu}_k)(p_i - \hat{\mu}_k)^T}{\sum_{j=1}^N p(p_j|\theta_k)} \quad (5)$$

and finally repeat these two steps until convergence.

1 iteration of this algorithm in MATLAB is written as:

```

1 % E-Step: Compute expectations
2 for k = 1:K
3     sqrt_Sigma = 1/ sqrt( det( Sigma(:, :, k) ) );
4     Xminus = X - Mu(:, :, k);
5     invSX = -0.5 * (Sigma(:, :, k) \ Xminus);
6     PKX(:, k) = Pk(k) * sqrt_Sigma * exp( dot(Xminus, invSX, 1) );
7 end
8 PKX = PKX ./ sum(PKX, 2); % Normalize posterior
9
10 for k = 1:K
11     Pk(k) = sum(PKX(:, k)) / N;
12 end
13
14 % M-Step: Maximise parameters values based on the likelihood function

```

```

15 for k = 1:K
16     Mu_new(:,k) = 0; % Compute new values for Cluster centers
17     PKXX = X';
18     for m = 1:M
19         PKXX(m,:) = PKX(:,k) .* X(m,:);
20     end
21     Mu_new(:,k) = Mu_new(:,k) + sum( PKXX, 2) / sum( PKX(:,k) );
22
23     Sigma_new(:,k) = 0; % Compute new values for Cluster variances
24     Xminus = X - Mu_new(:,k);
25     PKXminus = Xminus;
26     for m = 1:M
27         PKXminus(m,:) = PKX(:,k) .* Xminus(m,:);
28     end
29     Sigma_new(:,k) = Sigma_new(:,k) + PKXminus * Xminus' / sum( PKX(:,k) );
30     % In case sigma is a non invertible matrix
31     if rcond(Sigma_new(:,k)) < 10^(-5*M)
32         for m = 1:M
33             Sigma_new(m,m,k) = 1;
34         end
35     end

```

Thus for the pixels in the RGB space, each pixel intensity vector can be replaced by a centroid vector. A centroid vector is formed using the probabilities of the indicator variables for each data point along with the  $K$  cluster means using the formula  $\sum_k P(k|x_i)\mu_k$ .

### 1.3 Results

Image segmentation results using the K-means algorithm for  $K = 2, 3, 5, 10$  are shown in the figures. The normalised within cluster sum of squares for  $K = 1, \dots, 10$  clearly has a sharp decrease until  $K = 3$  and then plateaus/flat-line a little<sup>1</sup>. Given that we would prefer a simpler model (Occam's Razor) and thus a smaller value of  $K$ , this sharp turn indicates that a value of  $K = 3$  is a good choice to segment the colour pixels of this image. This result is also reinforced by our intuition as looking at the image, there are clearly 3 dominant colours present: black, green, and blue.

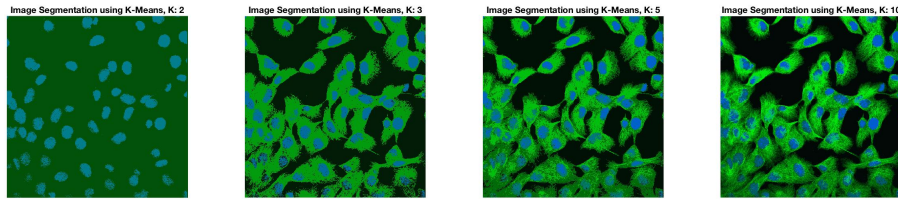


Figure 1: Image segmentation using the K-means algorithm for different values of  $K$

Image segmentation results using a Gaussian Mixture Model algorithm for  $K = 2, 3, 5, 10$  are as shown. There is large variation in the likelihood of the model fitted to the data in GMMs due to the random starting positions of the initial cluster centres. This can be seen in the likelihood figures below. There is flatlining of the negative LogLikelihood after a large drop in one of the cases at  $K = 4$ . Also later using cross validation it can also be seen that the negative LogLikelihood plateaus at  $K = 4$  after a large drop. Thus using the idea of the "Elbow method" and Occam's Razor it seems that  $K = 4$  is a reasonable choice for a Gaussian Mixture Model.

It should be noted that the inferences in both the cases are subjective.

<sup>1</sup>This method is often called the "Elbow Method"

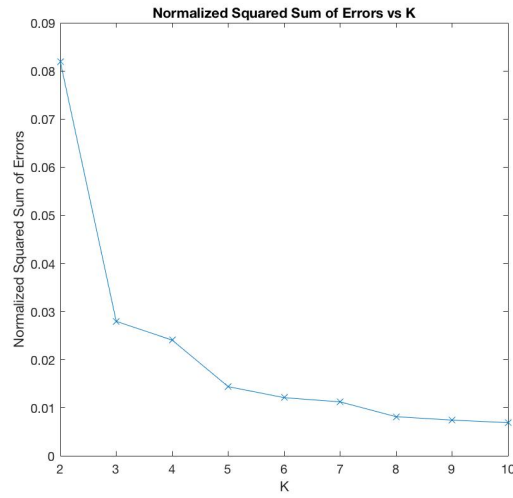


Figure 2: The normalised within cluster sum of squares vs K values using the K-means algorithm

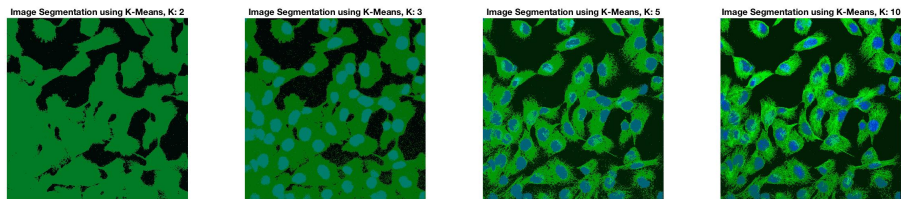


Figure 3: Image segmentation using a Gaussian Mixture Model for different values of K

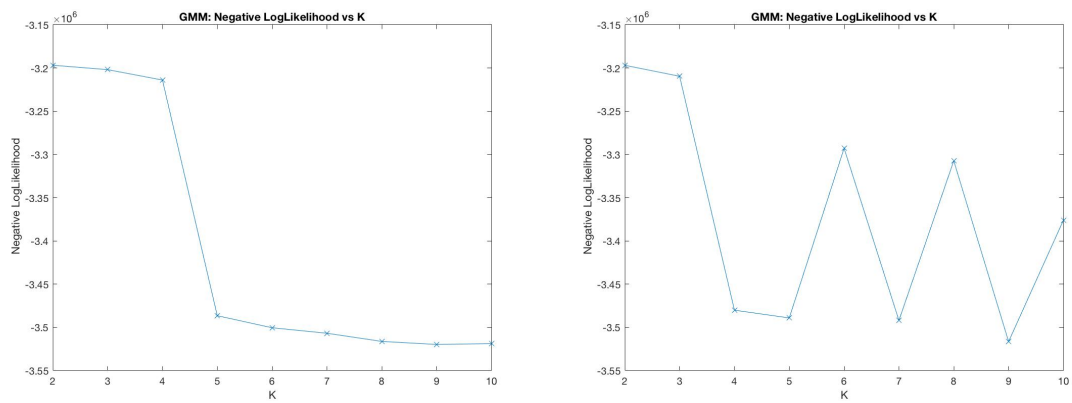


Figure 4: The normalised within cluster sum of squares vs K values using a GMM with K clusters

Another technique that was used to try and see if a better choice of  $K$  could be made was 2 fold cross validation. I took random permutation of the data and split it into test and training sets.

```
1 Xperm = X(:,randperm(length(X),length(X))); % X is the entire data here of size (M*N)
2 Xtest = Xperm(:,1:N/2);
3 Xtrain = Xperm(:, N/2+1:N);
```

Then fit the parameters of the models using the training data and compared the fit to the test data. But the results were mostly inconclusive except for reinforcing the idea that  $K = 4$  is a reasonable cluster number choice for a GMM. I believe this because many of the pixels lie very close together in the 3D RGB space which can be seen on the scatter figure produced.

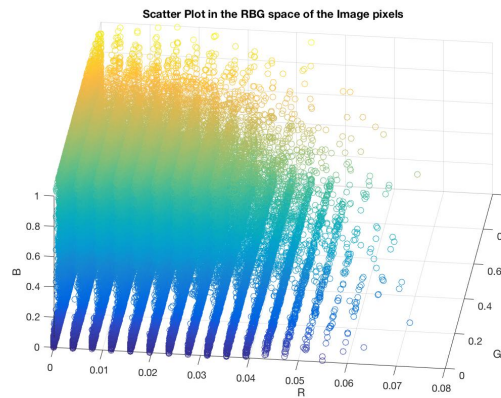


Figure 5: Scatter graph of the pixels in the normalized RGB 3D space. Shows how most of the pixels are concentrated in 1 corner

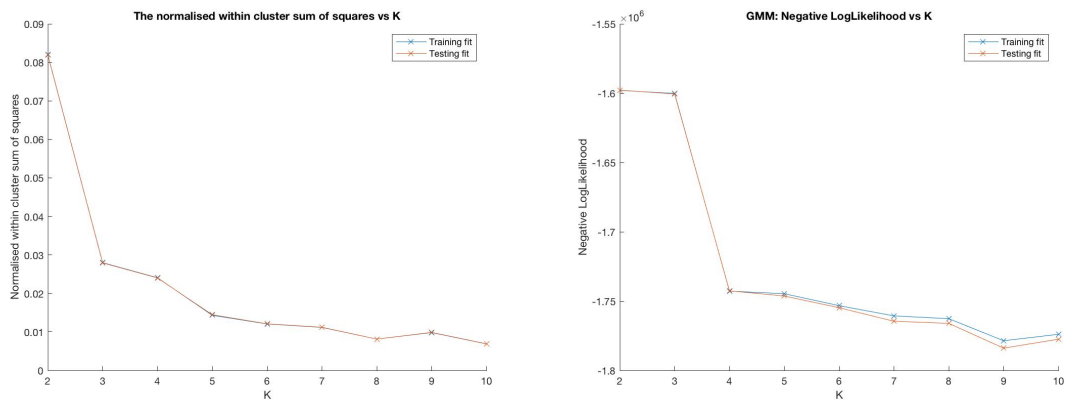


Figure 6: The results for cross validation using K means (left) and GMM (right) fitting. It can clearly be seen that the squared sum of errors and the likelihood for both the test and training data are very closely matched

## 1.4 Counting the Number of Cells

After an image has been segmented, there is still the issue of counting the number of cells present in the image. In order to do so, I extract the cluster of the images with its centre located closest to the blue colour vector.

This task has been attempted to be automated by using a Gaussian Mixture Model or K-Means algorithm to try and fit it with K clusters. Then model comparison can be utilised to find the value of K. First lets look at the negative LogLikelihood for various values of K fit onto the cell data.

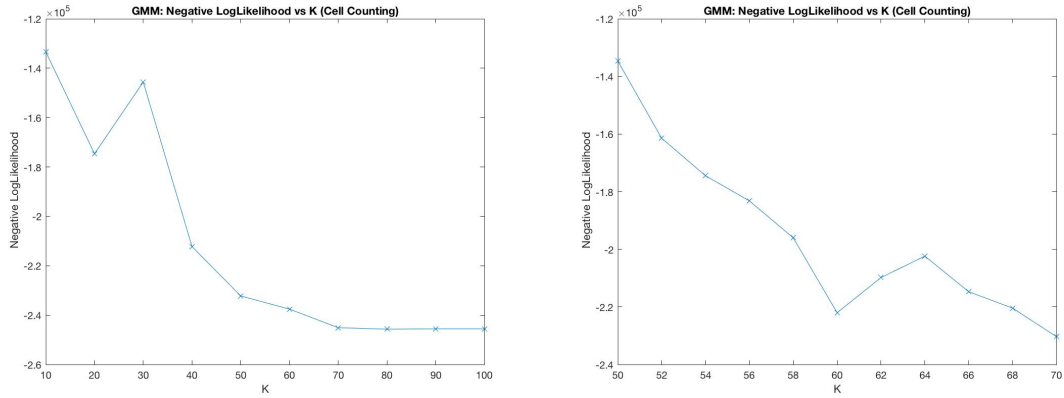


Figure 7: Negative LogLikelihood vs K for fitting a GMM to the cell data. The figure on the right is a higher resolution version of the one on the left

The negative LogLikelihood for the GMM fits for various values of K do not have as sharp a drop to plateau change. Thus it is hard to find a good cut-off point by inspecting the graph. But the value flatlines after  $K = 80$ . Running the same test for values between  $K = 50$  and  $K = 70$  is useful and done on the image on the right. It indicates that  $K = 62$  might be a reasonable choice.

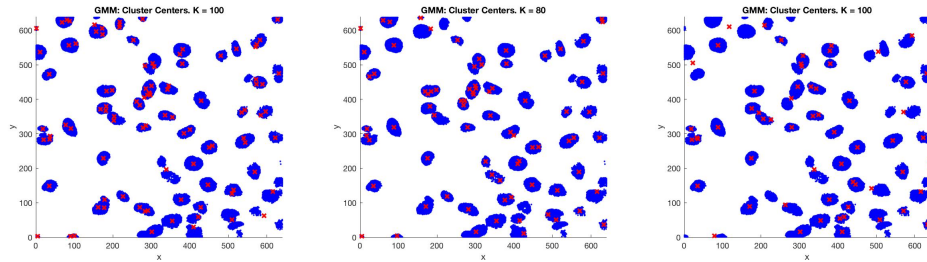


Figure 8: The cluster centres for fitting a GMM to the cell data for various values of K. The effect of some cluster centres getting 'stuck' between other clusters can be observed here.

Due to the random starting points of the GMM algorithm, it does not always produce optimal segmentation. i.e. the result converges but not necessarily to the best solution. Thus looking at the plots of the cells with clusters centres marked it can be seen that some clusters are, in a sense, stuck between two other clusters and thus cannot move towards a visually unoccupied cluster of points

that do not have a cluster near them. Hence simply looking at the likelihood is not a good idea. Although the above plots give an indication that the answer might be  $K = 62$ .

## 2 Question 2: Bayesian Linear Regression and Gaussian Processes

1. Given data <sup>2</sup>,  $\mathcal{D} = \{x_t, y_t\}_{t=1}^N$ , linear regression models the data as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} = f(\mathbf{X}) + \boldsymbol{\epsilon} \quad (6)$$

where  $\mathbf{y}$  is the vector of regressand values,  $\mathbf{X}$  is a matrix whose rows contain the covariates for each data point, and  $\boldsymbol{\epsilon}$  is a vector of the noise.  $f(\mathbf{X})$  is then the hidden function of time.

This model describes a family of functions and a single realisation is generated by choosing values for the parameter vector  $\boldsymbol{\beta}$  along with appropriate parameters for the noise vector. In this question it is assumed that  $\epsilon_t \sim N(0, \sigma_\epsilon^2)$ .

The likelihood function of the data given the parameter values is

$$P(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}, \sigma_\epsilon) = \prod_{t=1}^N p(y_t|x_t, \boldsymbol{\beta}, \sigma_\epsilon) = \prod_{t=1}^N \mathcal{N}_{y_t}(\mathbf{x}_t\boldsymbol{\beta}, \sigma_\epsilon) \quad (7)$$

Maximizing this likelihood function gives the most likely set of parameters. But in a Bayesian framework, a probability distribution over the parameters is of interest,  $p(\boldsymbol{\beta}, \sigma_\epsilon|\mathbf{y}, \mathbf{X})$  i.e. posterior distribution for the parameters given the data. This posterior can be found using Bayes Rule.

Bayes Rule for a model with parameters  $\Theta$  and given data  $\mathcal{D}$  is

$$P(\Theta|\mathcal{D}) = \frac{P(\mathcal{D}|\Theta)P(\Theta)}{P(\mathcal{D})} \quad (8)$$

where  $P(\Theta|\mathcal{D})$  is the posterior probability (the probability of a certain model parameters given the data),  $P(\mathcal{D}|\Theta)$  is the probability of observing the data given a model with parameters  $\Theta$ ,  $P(\Theta)$  is the probability of certain model parameters, i.e. our prior beliefs about the model, and  $P(\mathcal{D})$  is the probability of obtaining the given data over all values of  $\Theta$ .

Going back to our case of Bayesian Linear regression and using the simplification that  $\sigma_\epsilon$  is known this gives that the posterior distribution can be written as

$$\begin{aligned} p(\boldsymbol{\beta}|\sigma_\epsilon, \mathbf{y}, \mathbf{X}) &= \frac{p(\sigma_\epsilon, \mathbf{y}, \mathbf{X}|\boldsymbol{\beta})p(\boldsymbol{\beta})}{p(\sigma_\epsilon, \mathbf{y}, \mathbf{X})} \\ &\propto p(\sigma_\epsilon, \mathbf{y}, \mathbf{X}|\boldsymbol{\beta})p(\boldsymbol{\beta}) \end{aligned} \quad (9)$$

Now assuming an independent Gaussian prior over each of the parameters in the weight vector,  $\beta_i$ , gives

---

<sup>2</sup>In this question, I constrain my independent variable to be 1 dimensional given the CO2 dataset being analyzed

$$p(\boldsymbol{\beta}) = \mathcal{N}(0, \Sigma_{\boldsymbol{\beta}}) = \frac{1}{\sqrt{2\pi}|\Sigma_{\boldsymbol{\beta}}|^{1/2}} \exp -\frac{1}{2}\boldsymbol{\beta}^T \Sigma_{\boldsymbol{\beta}}^{-1} \boldsymbol{\beta} \quad (10)$$

$$p(\sigma_{\epsilon}, \mathbf{y}, \mathbf{X} | \boldsymbol{\beta}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma_{\epsilon}} \exp -\frac{1}{2\sigma_{\epsilon}^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \quad (11)$$

Using the fact that the product of two gaussians is also gaussian this gives another gaussian distribution for the posterior probability. This can be written as

$$p(\boldsymbol{\beta} | \sigma_{\epsilon}, \mathbf{y}, \mathbf{X}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (12)$$

where  $\boldsymbol{\mu} = \frac{1}{\sigma_{\epsilon}^2} (\boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{X}^T \mathbf{X})^{-1}$  and  $\boldsymbol{\Sigma} = \frac{1}{\sigma_{\epsilon}^2} (\boldsymbol{\Sigma}_{\boldsymbol{\beta}}^{-1} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

So finally, the underlying hidden function of time  $f(\mathbf{X})$  in the form  $\mathbf{X}\boldsymbol{\beta}$  can be considered as a linear combination between the covariate matrix  $\mathbf{X}$  defined by the independent variables and the parameter vector,  $\boldsymbol{\beta}$ , for which a posterior distribution can be found using the given data.

2. Consider two points  $x_i$  and  $x_j$ . For linear regression each point can be represented as a linear combination of its basis function (the basis functions form the covariates from the independent variables) as  $f(\mathbf{X}_i) = \mathbf{X}_i \boldsymbol{\beta} = \sum_k \beta_k \phi_k(\mathbf{x}_i) = \phi(\mathbf{x}_i) \boldsymbol{\beta}$ .

So the covariance of the function  $f$  for the two time points will be

$$\text{Cov}(f(\mathbf{X}_i), f(\mathbf{X}_j)) = \mathbb{E}[f(\mathbf{X}_i) f(\mathbf{X}_j)^T] - \mathbb{E}[f(\mathbf{X}_i)] \mathbb{E}[f(\mathbf{X}_j)] \quad (13)$$

$$\mathbb{E}[f(\mathbf{X}_i)] = \mathbb{E}[\phi(\mathbf{x}_i) \boldsymbol{\beta}] = \phi(\mathbf{x}_i) \mathbb{E}[\boldsymbol{\beta}] = 0 \quad (14)$$

as  $\boldsymbol{\beta} \sim \mathcal{N}(0, \Sigma_{\boldsymbol{\beta}})$ . Therefore

$$\begin{aligned} \text{Cov}(f(\mathbf{X}_i), f(\mathbf{X}_j)) &= \mathbb{E}[f(\mathbf{X}_i) f(\mathbf{X}_j)^T] = \mathbb{E}[\phi(\mathbf{x}_i) \boldsymbol{\beta} \boldsymbol{\beta}^T \phi(\mathbf{x}_j)^T] = \phi(\mathbf{x}_i) \mathbb{E}[\boldsymbol{\beta} \boldsymbol{\beta}^T] \phi(\mathbf{x}_j)^T \\ &= \phi(\mathbf{x}_i) \Sigma_{\boldsymbol{\beta}} \phi(\mathbf{x}_j)^T \end{aligned}$$

Note  $\mathbf{x}_i$  is a row of values and thus  $\mathbf{x}_i \mathbf{x}_j^T$  defines an inner-product.

Thus the covariance is a linear combination of Gaussian variables. A Gaussian process is defined by its mean function  $\mu(\mathbf{x})$  and Covariance matrix,  $K(\mathbf{x}, \mathbf{x}')$ . In this case defining  $K(\mathbf{x}, \mathbf{x}') = \text{Cov}(\mathbf{x}, \mathbf{x}')$  gives that  $f$  is then a gaussian process.

3. In the following section are results produced by fitting a Gaussian Process to the CO2 data recorded in Hawaii. This is done by first finding the hyperparameters that maximise the log likelihood of the data. This is done in GP.m. The basic process is as follows [1].

Firstly assumptions must be made by looking at the data set on what covariances function to utilize. In the figures below I have used 3 different types of Kernels.

The first is a linear kernel function with added white noise. This has the form.

$$K(x_i, x_j) = \sigma_f(x_i - c)(x_j - c) + \sigma_n \delta_{x_i, x_j} \quad (15)$$



The second is a squared exponential function with added white noise

$$K(x_i, x_j) = \sigma_f e^{-\frac{1}{2l_1^2}(x_i - x_j)^2} + \sigma_n \delta_{x_i, x_j} \quad (16)$$

And finally the third and the one with most hyper parameters was a squared exponential plus squared sin exponential

$$K(x_i, x_j) = \sigma_f e^{-\frac{1}{2l_1^2}(x_i - x_j)^2} + e^{-\frac{1}{2l_2^2} \sin(f * (x_i - x_j))^2} + \sigma_n \delta_{x_i, x_j} \quad (17)$$

Each of these Kernels has a different hyper-parameters and for each case these must be optimized. This can be done by using the `fminunc()` function in MATLAB where an input of the Kernel generating function that itself takes the hyper parameters and training data as input and gives the negative LogLikelihood of the Kernel and its derivatives with respect to the various hyper parameters. This is done in `OptimHP.m`.

```
1 [theta, fval] = fminunc(@(hyperparameters) optimHP(hyperparameters, constants), startVals(s,:), options);
```

An example of a gradient of with respect to the hyper-parameters can be seen here for the squared exponential [2]:

```
1 % Covariance
2 K = sigma.f^2 * exp(-(x1-x2)^2/(2*l^2));
3 % Derivatives
4 d_l = K * (1^-3) * (x1-x2)^2;
5 dsigma.f = 2*sigma.f * exp(-(x1-x2)^2/(2*l^2));
6
```

Once the hyper-parameters have been optimized, then the realized Kernel matrix is computed from the training data. Using this the mean for the test data along with 95% confidence intervals bounds can be found.

```
1 % Fit test data to the GP with the found hyperparameters
2 [ytest, bounds, K, kstar] = computeGP (@K.SE, X, y, theta, xtest);
```

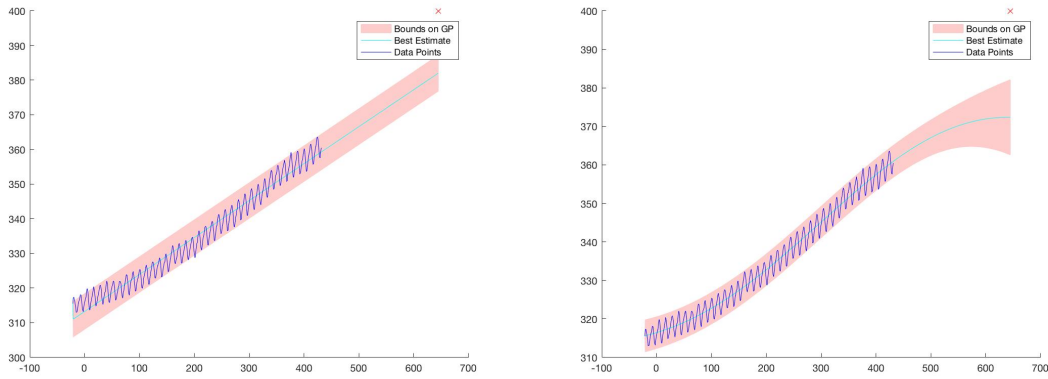


Figure 9: Predictions made by my Gaussian Process models in the 95% confidence interval. The observation of 400pm in 2013 is marked with a red cross. The x-axis here is Months after 1960. The y-axis is the CO2 levels.

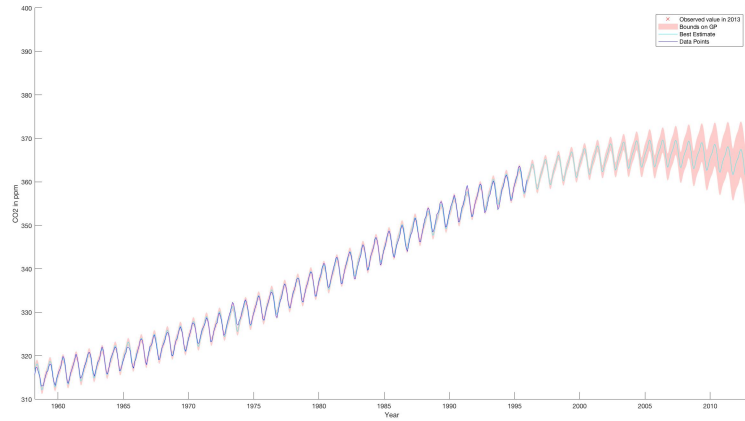


Figure 10: Predictions made by my Gaussian Process models in the 95% confidence interval. The observation of 400ppm in 2013 is marked with a red cross

The graphs above have 95% confidence intervals plotted for the data points until 2013. Looking at the various fits of the Gaussian Process for the different Kernel functions, it can be seen none of them predict the dramatic increase in CO2 level ppm recorded in 2013.

The linear Kernel is unable to capture the nature of the exponential rise in the data.

And the squared exponential and periodic kernels are limited due to the training data stopping about 20 years before the prediction point. This point at 2013 is far out from the actual training data and thus the non-noise components start to diminish and the kernel becomes dominated by the noise element leading to the flattening of the curve.

An exponential increase occurred in the CO2 emission level between 1993 and 2013. This is hard to model using a Gaussian Process as seen from the previous part that the underlying generating function is a linear combination of basis functions. And in the case of a Gaussian Process, having an exponentially increasing or decreasing basis function is hard to model.

### 3 Question 3: Social Consequences of AI and Machine Learning

Machine learning, and artificial intelligence (AI<sup>[3]</sup>), have gained popularity in recent years. The heart of these techniques lies in trying to find optimisations such that a computer can 'learn' how to handle a specific set of problems given enough training data rather than finding explicit solutions. The wide scope of these techniques means that no single researcher can hope to assess all of their possible impacts. In light of this, I shall look at an example that is near the horizon and will quite possibly have a large impact.

My chosen example is that of autonomous cars. Many leading automobile companies, such as Tesla [3], BMW, and more have invested significantly in this field along with technology companies such as Google and Intel as well [4]. Driverless cars use GPS systems and map databases to locate themselves in their environments. Moreover, they utilise machine vision systems, often using lasers, to map out their immediate surroundings. The systems then need to classify various elements on the road such as cars, busses, cyclists and pedestrians. This is done through combinations of PCA basis analysis and classification algorithms. Furthermore, they also need to predict where each element on the road will be moving in the future and utilize algorithms such as Hidden Markov Monte Carlo techniques, as well as classical IF-THEN rules [5] [6].

It is estimated that over 1 hour per day, per person, is spent driving in many European countries [7]. Consider a reality where people no longer need to drive cars at all. This time could instead be utilized for other activities. People could sleep en-route and thus live further away from work leading to less densely populated urban areas. Rented cars might immerse us in advertisements that are personalized using machine learning algorithms, just like for webpage adverts nowadays [8]. Furthermore, there will be a revolution in delivery services that could lead to a large number of jobs becoming automated. A large change in one industry will undoubtedly impact entire economies. Only time will tell whether these impacts will be, on the whole, positive or negative.

A key consequence of autonomous cars will be our requirement of having to deal with new ethical grey-areas. Suppose a car with a passenger runs over a pedestrian, killing them instantly. If the situation is deemed avoidable, is the car company at fault or is it the passenger? Does it matter if the passenger owned the car or rented it? New laws will have to be written and entire systems, such as driving licenses, redone. Already many automobile owners in China, Germany, Korea, and USA have started to use driver-assistance systems [9]. I do not believe the transition to autonomous cars will be sudden, but as their prevalence rises I am certain they will impact us all regardless of whether we choose to adopt them. I only hope that their impact is not physical.

---

<sup>3</sup>For the sake of clarity I'd like to distinguish between AI, artificial intelligence, and AGI, artificial general intelligence. The former is able to tackle a set of given tasks, e.g. categorising images based on the objects present in them. The latter is able to tackle any sort of problems/tasks and is comparable to creating intelligence that is, in a sense, similar to our own intelligence.

## 4 Question 4: Bitcoin Transactions

Before choosing a model, I believe it is useful to understand and visualise the data. A key issue with the data set provided is that each transaction does not have a time stamp but only a date stamp along with a transaction ID (assumed to be chronological in time <sup>4</sup>). Hence the data is not linearly spaced time. For example, it is not possible to know whether 1 second passed between 2 transactions or 1 hour. This notion is further reenforced by noting the variation in the number of transactions per day.

```
1 NumTransactionsPerDay
2 = 1261 3261 3754 1485 1447 1993 2928 1231 2473 115]
```

The data can then be visualized by looking at the spread of transactions per day. This gives

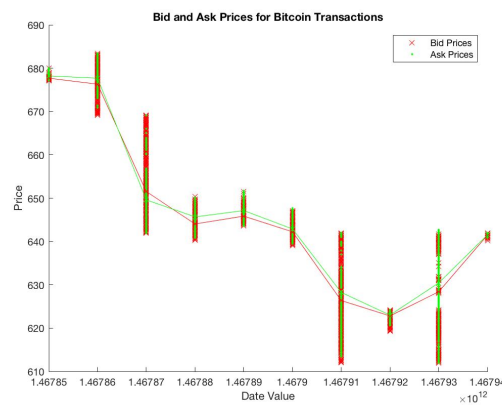


Figure 11: Bid and Ask prices for the data set. Also the lines represent the average bid and ask price traded each day.

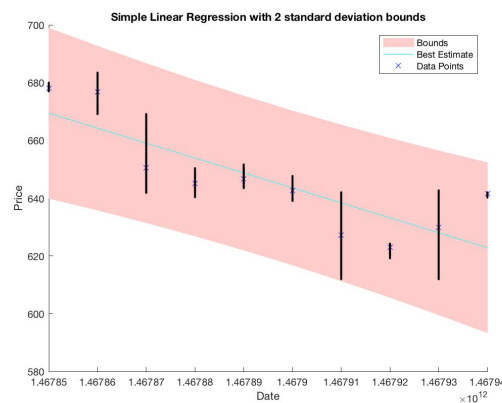


Figure 12: Simple linear regression on the data points that represent the average price of the day

<sup>4</sup>The provided data does have the transaction IDs in order either and thus the data has been sorted to fix this. See Figure 14 in the Appendix for more details

Simple linear regression gives extremely large prediction bounds (Figure 12) that would not be very useful in order to predict the price and make transactions. An example of this can be seen in the figure produced.

Even though the data is not linearly spaced in time. It is helpful to see the impact the bid and ask price have on the trend of the data. It can clearly be seen that often when transactions are repeatedly occurring at bid price, the price is falling and often when the transactions are repeatedly occurring at ask price, the price is increasing. This agrees with our basic market knowledge that if people are willing to pay some incremental amount more than the bid to match the ask price then there is larger demand than supply and thus the price moves up (and vice versa).

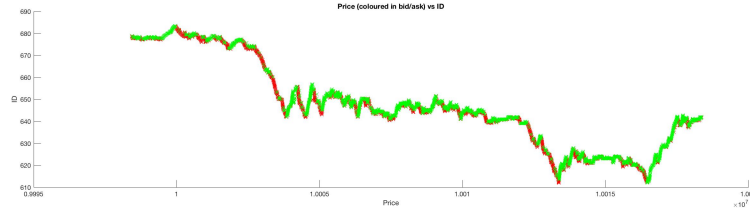


Figure 13: A plot of the prices of Bitcoin in Price vs ID. Coloured red for bid price and green for ask price.

Now as the assumption is made that the transaction IDs are chronological in time, this Bitcoin data must then be treated as a times-series. The price at transaction  $t$  appears to be dependent on the prices of the transactions before along with the bid-ask binary information for each transaction. This is an ideal problem for a neural network to solve!

The Neural Network Architecture that I have considered has 3 layers with a sigmoid activation function. The input layer is specified by the size of number of prices plus bid/ask binaries inputted to the system. The hidden layer of a chosen size (mostly through trial and error - more on this later). And finally the output layer, in my chosen model, has 2 nodes representing a binary probability of the price going up or down in the future.

A key element for this model is determining what it means for the price to go up or down in the future. Simple taking the difference between the price between two immediate transactions is not reasonable. Visualising the bit-coin data for small intervals shows that it behaves very much like any other financial instrument in that there is huge amount of fluctuations in the price at small time intervals. But nonetheless there are clear trends in the direction the average price is heading. So to test the data a moving average is considered. It is set up in a simple way

$$\text{movingAvg}_t = \frac{1}{\tau} \sum_{i=1}^{\tau} \text{movingAvg}_{t-i}$$

where  $\tau$  is the size of the moving average window. So now in our model we wish to input some known values up to the point  $t$  and get an indication of whether the average price in the next  $\tau$  points will be higher or lower than the current average price (also over the same  $\tau$  time period).

$N_I$  is the number of input layers.  $N_H$  is the number of hidden layers. And  $N_O$  is the number of output layers.

Then for a 3 layer neural network with input vector  $X \in \mathbb{R}^{N_I}$ ,

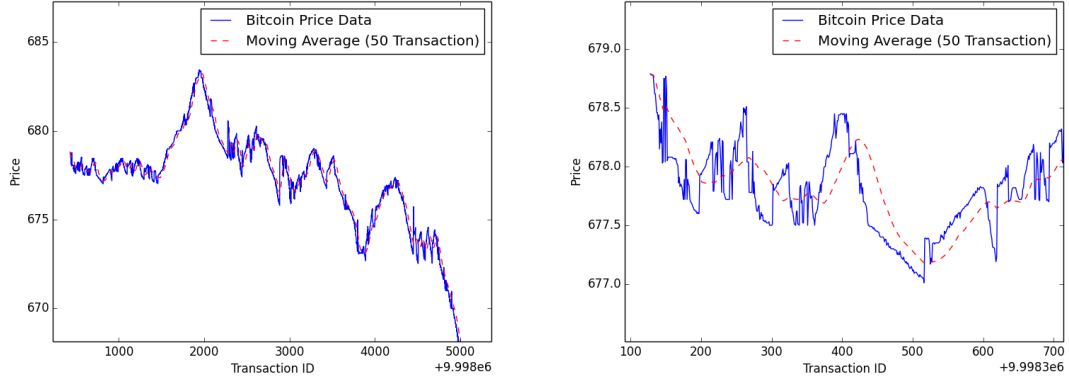


Figure 14: Moving Average over 50 transactions for the Bitcoin Data. The right figure is a zoomed in version of the left. The moving in red is for the previous 50 data points and comparing the moving average between 2 points is more indicative than simply comparing the prices between them.

$$\begin{aligned}
u_1 &= W_1 X + c_1 \\
v_1 &= \tanh(u_1) \\
u_2 &= W_2 v_1 + c_2 \\
v_2 &= \frac{\exp(u_2)}{\sum_i \exp(u_{2i})}
\end{aligned} \tag{18}$$

where  $W_1 \in \mathbb{R}^{N_H \times N_I}$ ,  $c_1 \in \mathbb{R}^{N_H}$ ,  $W_2 \in \mathbb{R}^{N_O \times N_H}$ , and  $c_2 \in \mathbb{R}^{N_O}$ . Also  $v_1 \in \mathbb{R}^{N_H}$  and  $v_2 \in \mathbb{R}^{N_O}$ . The model learns but propagating the error backwards by using a test solution vector  $y \in \mathbb{R}^{N_O}$  such that

$$\begin{aligned}
\Delta_1 &= v_2 - y \\
\delta W_2 &= \Delta v_1^T \\
\delta c_2 &= \sum_i \Delta_{1i} \\
\Delta_2 &= W_2^T \Delta_1 (1 - v_1^2) \\
\delta W_1 &= \Delta_2 X^T \\
\delta c_1 &= \sum_i \Delta_{2i} \\
\delta W_2 &= \delta W_2 + \epsilon_2 W_2 \\
\delta W_1 &= \delta W_1 \epsilon_2 W_1 \\
W_1 &= W_1 - \epsilon_1 \delta W_1 \\
c_1 &= c_1 - \epsilon_1 \delta c_1 \\
W_2 &= W_2 - \epsilon_1 \delta W_2 \\
c_2 &= c_2 - \epsilon_1 \delta c_2
\end{aligned} \tag{19}$$

## 4.1 Results

Given that the data is sequential, I believe it is not possible to shuffle the data around, else the intrinsic behaviour of the impact between bid/ask price and the movement of price is lost. Furthermore as the neural network components are initialized to random values it is important to model test multiple times to ensure the results are not generated simply due to chance and rather due to the network learning from the training set.

One method is to use 10-fold cross validation where a new random model is initialized for each set of data to be tested. Using this I see get an averaged prediction performance over all testing sets of 0.547%. The values of individual sections can be seen in Figure 15. Some aspects of interest are that the sections with relatively flat prices and/or that have less autocorrelated data (more up/down spikes) are closer to 50% accuracy as we would reasonably suspect. Sections with strong trends in either direction and/or large autocorrelations perform significantly better.

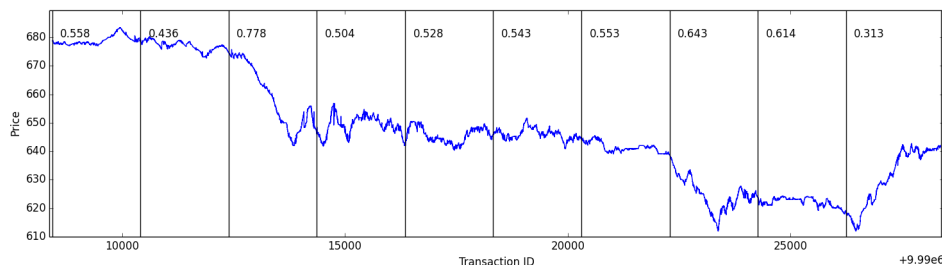


Figure 15: 10-fold Cross Validation. Lines display where the sets are split and the percentage of values correct that specific set was tested using all the other sets for training.

Another method of testing used given the sequential nature of the data is iteratively testing the data using the first X points as training data and then testing on the next X followed by training on the currently tested points and repeating this process multiple times for different random starting conditions of the neural network. It can be seen that this method perhaps with an average performance of 52.9% for batches of size 2000 and of 53.4% for batches of size 4000. The batch size is the value of X defined before. Once again the model predicts better when there are strong underlying trends and worse when the data is mostly flat.

Overall it seems that the neural network is able to get a 52-55% rate of predicting the price correctly over the entire data set. In terms of trading an asset this is quite significant as over time a lot of money could be made with such an advantage. But there are also several assumptions made in the implementation that compromise this result.

1. Firstly, this model does not have an accurate time-stamp and thus it is not possible to know whether all these transactions occurred in a small time interval. Perhaps one even too small for a neural network to train and predict based on past values.
2. Secondly, it is assumed that transactions IDs are chronological in time but this must be verified.
3. Most importantly, our training set is quite small - only 10 days. This means that even though the results have been cross validated, it is hard to assess whether the entire data set is rep-

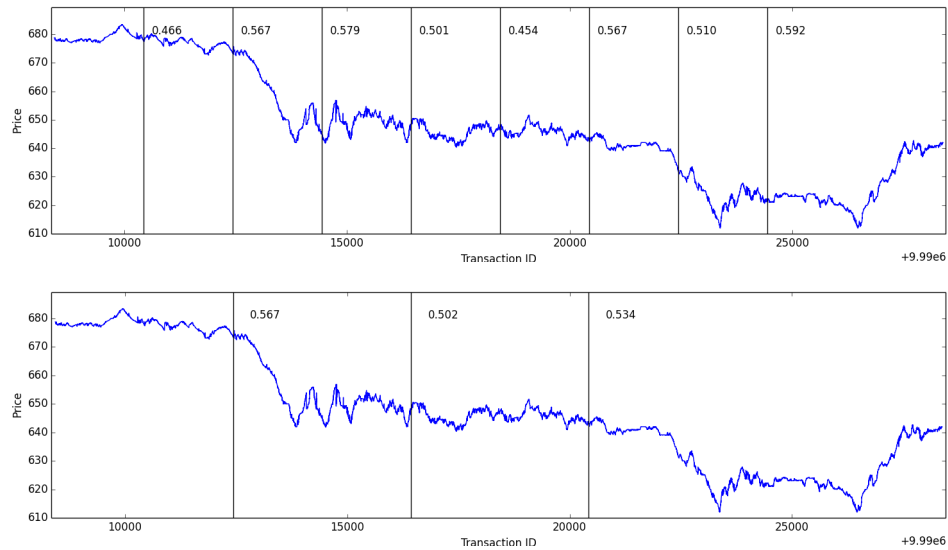


Figure 16: Iterative Batch Cross Validation with 10 Repeats. In the upper figure, the first 2000 values are used to training and then the next 2000 are tested after which they are used for training and this process is iteratively repeated. In the lower figure the same has been done but for a batch size of 4000. The values indicate the performance (percentage of times the model correctly predicted the price trend) in each section

representative of the general market of Bitcoin transactions. Furthermore, this ties into the idea that neural networks require tremendous amounts of data to train properly.

4. Lastly, the analysis of neural networks is still in its infancy. There are no exact methods for determining what sort of architecture will lead to the best result. This also makes it extremely hard to choose and then justify parameter choices, such as the learning rate.

A similar implementation could be attempted on properly formatted data. Properly formatted data is received regularly, i.e. in set time steps, and sequentially in time. Without regularity, many forms of inference do not make sense such as moving averages and volatility which would be extremely beneficial here. Also, then other simpler and better understood models, for example AR regression models, could be implemented for comparison given the lack of rigour involved in neural networks.

Ultimately, I would not trade based on this model but I think it is a good proof of concept.



# Appendix

## 4.2 Q4: Figure

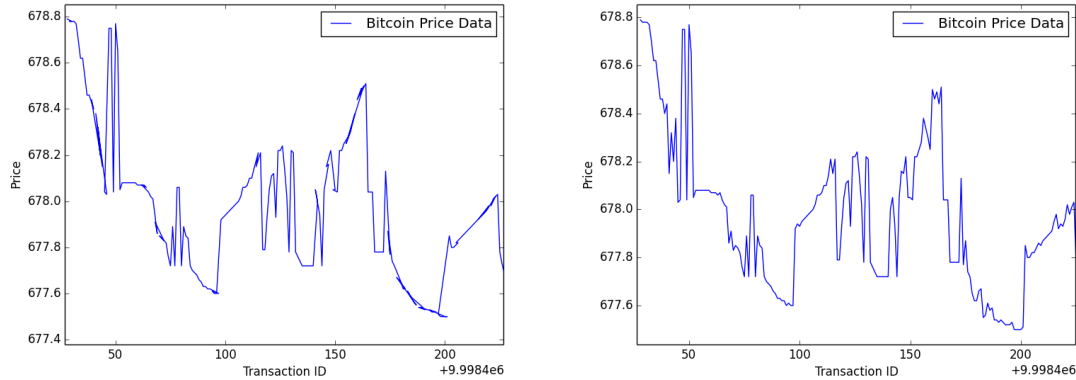


Figure 17: A zoomed in view of small intervals in Transaction IDs with the figure on the left plotted without transaction IDs sorted and the one of the right plotted with the transaction IDs sorted. The locally random ordering of transaction ID data is not visible on a global scale but clearly impacts any implementation that is dependent on assuming the Bitcoin data is time-series data and the underlying function for the price depends on the previous prices

## 4.3 Code

Note that the bound plotting function that I have used is heavily influenced from [1]. Also the Gaussian Process code is implemented in line with the same tutorial. Furthermore I believe it is worth mentioning that I have used a similar code for investigations in my on-going M4R course as well.

### 4.3.1 Q1

Kmeans.m

```
1 function [Z, Mu, EK] = Kmeans( X, K, tol, display )
2 % Uses the K Means algorithm to segment the inputted data into K clusters
3 % INPUTS: X: M x N points of data to be segmented using a gaussians
4 %         K: Number of clusts. Default set to 3.
5 % OUTPUTS: Z: Contains the centroid vector for each data point
6 %         C: Defines the cluster means for each centroid vector
7 %         Ek: Normalized Within Cluster Sum Squared of Errors
8
9 if (nargin < 2) K = 3; end % Set default cluster value to 3 if not given
10 if (nargin < 3) tol = 1e-10; end % Set default value for the break tolerance
11 if (nargin < 4) display = false; end % Display iteration count/means
12
13 [M,N] = size(X); % Extract the size of K
14
15 Z = zeros(N,1); % Array to store the indicator variables
16 dist = zeros(K,1); % Array to store the norm between a point and the cluster means
17 size_k = zeros(K,1); % Array to store the size of each cluster in each iteration
18 Mu = zeros(M,K); % Array to store the cluster means
19 Mu_new = rand(M,K); % Array to store the new cluster means
20 iters = 0; % iteration count
21
```

```

22 while norm(Mu - Mu_new) > tol
23     if (display) tic; end
24
25     Mu = Mu_new; % Set Mu to previously computed Cluster Center
26
27     dist = zeros(N,K); % Compute L2 norm
28     for k = 1:K
29         dist(:,k) = sum((X(:, :) - Mu(:,k)).^2,1);
30     end
31     [I,Z] = min(dist,[],2); % Get minimum for each N
32
33     for k = 1:K % Set new cluster centers to mean of their elements.
34         if isempty(X(:,Z == k)) == 1 % Incase a cluster has no members move it to a random point
35             Mu_new(:,k) = rand(M,1);
36         else
37             Mu_new(:,k) = mean(X(:,Z == k),2);
38         end
39     end
40     if display
41         iters = iters + 1
42         Mu_new = Mu_new
43         diff = norm(Mu - Mu_new)
44         toc
45     end
46 end
47
48 EK = 0; % Compute Normalized Within Cluster Sum of Squares
49 for k = 1:K
50     EK = EK + sum( sum( (X(:,Z == k) - Mu(:,k)).^2 ));
51 end
52 EK = EK / N;
53
54 end

```

## GaussianMixtureModel.m

```

1 function [Z, Mu, LK, PKX, Sigma] = GaussianMixtureModel( X, K, tol, display, Mu, Sigma)
2 % Uses the EM algorithm to use the Gaussian Mixture model to fit to the
3 % data
4 % INPUTS: X: M x N points of data to be segmented using a gaussians
5 %         K: Number of clusts. Default set to 3.
6 %         Mu: Optional initial conditions for the starting cluster means
7 %         Sigma: Optional initial starting cluster variance
8 % OUTPUTS: Z: Contains the centroid vector for each data point
9 %         Mu: Defines the cluster means for each centroid vector
10 %         Lk: The Likelihood
11 %         PKX: Matrix (N X K) The probability of point X_n to be in
12 %              cluster k
13 %         Sigma: Variance matrices of the K gaussians
14
15 [M,N] = size(X); % Get dimensions of data and gaussians to be fit
16
17 if (nargin < 2) K = 3; end % If K is not given
18 if (nargin < 3) tol = 1e-4; end % If tol is not given
19 if (nargin < 4) display = false; end % If display is not specified
20 if (nargin < 5) Mu = zeros(M,K); iMu = true; end % If Mu not given
21 if (nargin < 6) Sigma = zeros(M,M,K); iSigma = true; end % If Sigma not given
22
23 PKX = zeros(N,K); % Centroid vector for each X point
24 Pk = ones(K,1); % To store prior probabilities
25 Z = zeros(N,1); % Indicator for each X point
26
27 % STARTING CONDITIONS: Equally spaced cluster means
28 if iMu == true
29     for m = 1:M
30         mmax = max(X(m,:));
31         mmin = min(X(m,:));
32         mdiff = mmax - mmin;
33         Mu(m,:) = mdiff*rand(1,K);
34     end
35 end
36 if iSigma == true
37     for m = 1:M
38         Sigma(m,m,:) = 10^(-M)*ones(1,1,K);
39     end
40 end
41
42 % Variables to compare if clusters have stopped moving
43 Mu_new = Mu;
44 Sigma_new = Sigma;
45 iters = 0;
46 n1 = 0;
47 while (true)
48     if (display) tic; end % Start the clock for each iteration if display is on
49
50     % E-Step: Compute expectations
51     for k = 1:K
52         sqrt_Sigma = 1/ sqrt( det( Sigma(:, :, k) ) );
53         Xminus = X - Mu(:,k);
54         invSX = -0.5 * (Sigma(:, :, k) \ Xminus);
55         PKX(:,k) = Pk(k) * sqrt_Sigma * exp( dot(Xminus, invSX, 1) );
56     end
57     PKX = PKX ./ sum(PKX,2); % Normalize posterior
58

```

```

59     for k = 1:K
60         Pk(k) = sum(PKX(:,k)) / N;
61     end
62
63     % M-Step: Maximise parameters values based on the likelihood function
64     for k = 1:K
65         Mu.new(:,k) = 0; % Compute new values for Cluster centers
66         PKXX = X';
67         for m = 1:M
68             PKXX(m,:) = PKX(:,k) .* X(m,:)' ;
69         end
70         Mu.new(:,k) = Mu.new(:,k) + sum( PKXX, 2) / sum( PKX(:,k) );
71
72         Sigma.new(:,k) = 0; % Compute new values for Cluster variances
73         Xminus = X - Mu.new(:,k);
74         PKXminus = Xminus;
75         for m = 1:M
76             PKXminus(m,:) = PKX(:,k) .* Xminus(m,:)' ;
77         end
78         Sigma.new(:,k) = Sigma.new(:,k) + PKXminus * Xminus' / sum( PKX(:,k) );
79         % In case sigma is a non invertible matrix
80         if rcond(Sigma.new(:,k)) < 10-5 * M
81             for m = 1:M
82                 Sigma.new(m,m,k) = 1;
83             end
84         end
85     end
86
87     n2 = norm(Mu.new - Mu); % norm to check if the solution has converged
88     if (n2 < tol) || (abs(n1 - n2) < tol) % Stop once converged
89         Mu = Mu.new;
90         Sigma = Sigma.new;
91         break
92     else
93         n1 = n2;
94         Mu = Mu.new;
95         Sigma = Sigma.new;
96     end
97
98     if (display) % Display variables if asked for
99         toc
100         display(n2)
101         iters = iters + 1;
102         display(iters)
103     end
104 end
105
106 % Get indicator variables
107 for n = 1:N
108     [~,index] = max(PKX(n,:));
109     Z(n) = index;
110 end
111
112 % Compute LogLikelihood
113 LK = 0;
114 for k = 1:K
115     LK = LK - Pk(k) * N * log( det( Sigma(:,k) ) ) / 2;
116     Xminus = X(:,k) - Mu(:,k);
117     invSX = (Sigma(:,k) \ Xminus);
118     LK = LK - (1/2) * sum( PKX(:,k)' .* dot(Xminus, invSX, 1) );
119     LK = LK + 0.5 * Pk(k) * N * log(Pk(k));
120     % LPKX = -log(PKX(:,k));
121     % LPKX(LPXX == inf) = 0;
122     % LK = LK + sum( PKX(:,k) .* (LPKX - (M/2) * log(2*pi) ) );
123     % LK = LK + sum( PKX(:,k) .* (LPKX) );
124 end
125
126 end

```

test\_Kmeans\_imageSegmentation.m.

Note test\_GMM\_imageSegmentation.m is very similar and just calls the GMM function instead of the K means one.

```

1  img = imread('FluorescentCells.jpg'); % Load the image and format it
2  img = double(img);
3  [Nx,Ny,M] = size(img);
4  X = reshape(img, Nx*Ny,M)';
5  X = X/255;
6  % Initialize arrays needed to store data
7  K_vals = 2:10;
8  Nk = length(K_vals);
9  Z = cell(Nk,1);
10 Mu = cell(Nk,1);
11 EK = zeros(Nk,1);
12 TimeTaken = zeros(Nk,1);
13 % Compute K Means clustering for each value of K
14 for K = K_vals
15     tic;
16     [Z{K-1},Mu{K-1},EK(K-1)] = Kmeans(X,K);
17     TimeTaken(K-1) = toc;
18 end

```

```

19 % Plot results
20 figure();
21 subplot(1,4,1);
22 plotImage(Z{2-1},Mu{2-1},2,Nx,Ny);
23 title(['Image Segmentation using K-Means, K: ',num2str(2)]);
24 subplot(1,4,2);
25 plotImage(Z{3-1},Mu{3-1},3,Nx,Ny);
26 title(['Image Segmentation using K-Means, K: ',num2str(3)]);
27 subplot(1,4,3);
28 plotImage(Z{5-1},Mu{5-1},5,Nx,Ny);
29 title(['Image Segmentation using K-Means, K: ',num2str(5)]);
30 subplot(1,4,4);
31 plotImage(Z{10-1},Mu{10-1},10,Nx,Ny);
32 title(['Image Segmentation using K-Means, K: ',num2str(10)]);
33
34 figure();
35 plot(K_vals,EK,'x-');
36 title('Normalized Squared Sum of Errors vs K'); xlabel('K'); ylabel('Normalized Squared Sum of Errors')

```

## test\_GMM\_segmentationCrossValidation.m

```

1 % Load the image
2 img = imread('FluorescentCells.jpg');
3 % imshow(img);
4 img = double(img);
5 [Nx,Ny,M] = size(img);
6 N = Nx*Ny;
7 X = reshape(img, N,M)';
8 X = X/255;
9
10 Xperm = X(:,randperm(length(X),length(X)));
11 Xtest = Xperm(:,1:N/2);
12 Xtrain = Xperm(:, N/2+1:N);
13
14 K_vals = 2:10;
15 Nk = length(K_vals);
16 Z = cell(Nk,1);
17 Ztest = cell(Nk,1);
18 Mu = cell(Nk,1);
19 PKX = cell(Nk,1);
20 Sigma = cell(Nk,1);
21 LK = zeros(Nk,1);
22 LKtest = zeros(Nk,1);
23 TimeTaken = zeros(Nk,1);
24
25 for K = K_vals
26     tic;
27     [Z{K-1},Mu{K-1},LK{K-1},PKX{K-1}, Sigma{K-1}] = GaussianMixtureModel(Xtrain,K);
28     TimeTaken(K-1) = toc;
29     pkx = PKX{K-1};
30     Pk = zeros(K,1);
31     for k = 1:K
32         Pk(k) = sum(pkx(:,k)) / length(pkx(:,k));
33     end
34     [ Ztest{K-1} , LKtest(K-1)] = GMMTest( Xtest , K, Mu{K-1}, Sigma{K-1}, Pk);
35
36 end
37
38 figure(); hold on;
39 plot(K_vals,-LK,'x-','DisplayName','Training fit');
40 plot(K_vals,-LKtest,'x-','DisplayName','Testing fit');
41 title('GMM: Negative LogLikelihood vs K'); xlabel('K'); ylabel('Negative LogLikelihood')
42 legend('show')

```

## test\_CellCounting.m

```

1 % Load the image
2 img = imread('FluorescentCells.jpg');
3 % imshow(img);
4 img = double(img);
5 [Nx,Ny,M] = size(img);
6 N = Nx * Ny;
7 X = reshape(img, N ,M)';
8 X = X/255;
9 K = 3;
10
11 load('Kmeans.counting.mat'); % Has the data stored for the segmentation with K = 3
12 tic;
13 % [Z1,C1] = GaussianMixtureModel(X,K,1e-3); % Data can also be found using
14 toc
15 Z = Z1;
16 C = C1;
17
18 Z = reshape(Z, Nx, Ny);
19
20 S = [];
21 Xp = []; Yp = [];
22 for x = 1:Nx
23     for y = 1:Ny
24         if Z(x,y) == 2
25             S = [S, [x;y]];
26             Xp = [Xp, x];

```

```

27     Yp = [Yp, y];
28     end
29 end
30 end
31 S(1,:) = S(1,:) / Nx;
32 S(2,:) = S(2,:) / Ny;
33
34 K_vals = 10:10:100;
35 Nk = length(K_vals);
36 Z = cell(Nk,1);
37 Mu = cell(Nk,1);
38 LK = zeros(Nk,1);
39 TimeTaken = zeros(Nk,1);
40
41 for i = 1:Nk
42     K = K_vals(i);
43     tic;
44     [Z{i},Mu{i},LK{i}] = GaussianMixtureModel(S,K,1e-3);
45     timet = toc;
46     TimeTaken(i) = timet;
47 end
48
49 BIC = N * log(LK / N) + K_vals.* log(N);
50
51 figure();
52 subplot(1,3,1);
53 plotCellSegmentation(Xp,Yp,Z{i},K,640*Mu{i}, false)
54 title('GMM Cluster Centers. K = 100'); xlabel('x'); ylabel('y')
55 axis([0 Nx 0 Ny])
56 subplot(1,3,2);
57 plotCellSegmentation(Xp,Yp,Z{i-2},80,640*Mu{i-2}, false)
58 title('GMM Cluster Centers. K = 80'); xlabel('x'); ylabel('y')
59 axis([0 Nx 0 Ny])
60 subplot(1,3,3);
61 plotCellSegmentation(Xp,Yp,Z{i-4},60,640*Mu{i-4}, false)
62 title('GMM Cluster Centers. K = 100'); xlabel('x'); ylabel('y')
63 axis([0 Nx 0 Ny])
64
65 figure();
66 plot(K_vals,-LK,'x-');
67 title('GMM Negative LogLikelihood vs K (Cell Counting)'); xlabel('K'); ylabel('Negative LogLikelihood')

```

### KmeansTest.m

for segmenting the testing data with given clusters.

```

1 function [ Z , EK] = KmeansTest( X, K, Mu )
2 % Segments the inputted test data into K clusters with means Mu
3 % INPUTS:      X: M x N points of data to be segmented using a gaussians
4 %             K: Number of clusts. Default set to 3.
5 %             Mu: Cluster centers
6 % OUTPUTS:     Z: Contains the centroid vector for each data point
7 %             EK: Goodness of Clustering
8
9 [~,N] = size(X); % Extract the size of K
10
11 dist = zeros(N,K); % Assign to cluster with smallest L2 norm
12 for k = 1:K
13     dist(:,k) = sum((X(:, :) - Mu(:,k)).^2,1);
14 end
15 [~,Z] = min(dist,[],2);
16
17 EK = 0; % Compute Goodness of fit
18 for k = 1:K
19     EK = EK + sum( sum( (X(:,Z == k) - Mu(:,k)).^2));
20 end
21 EK = EK/N;
22
23 end

```

### GMMTest.m

for segmenting the testing data with given clusters means and variances.

```

1 function [ Z , LK, PKX] = GMMTest( X, K, Mu, Sigma, Pk )
2 % Segments the inputted test data into K clusters with means Mu
3 % INPUTS:      X: M x N points of data to be segmented using a gaussians
4 %             K: Number of clusts. Default set to 3.
5 %             Mu: Cluster means
6 %             Sigma: Cluster variances
7 %             Pk: Prior probabilities of each cluster
8 % OUTPUTS:     Z: Contains the centroid vector for each data point
9 %             LK: Negative Log Likelihood
10 %             PKX: Probability vectors
11
12 [M,N] = size(X); % Get dimensions of data and gaussians to be fit
13
14 PKX = zeros(N,K); % Centroid vector for each X point
15 Z = zeros(N,1); % Indicator for each X point
16

```

```

17 for k = 1:K % Compute Posterior
18     sqrt.Sigma = 1/ sqrt( det(Sigma(:, :, k) ) );
19     Xminus = X - Mu(:, k);
20     invSX = -0.5 * (Sigma(:, :, k) \ Xminus);
21     PKX(:, k) = Pk(k) * sqrt.Sigma * exp( dot(Xminus, invSX, 1) );
22 end
23 PKX = PKX ./ sum(PKX, 2); % Normalize
24
25 % Get indicator variables
26 for n = 1:N
27     [~, index] = max(PKX(n, :));
28     Z(n) = index;
29 end
30
31 % Compute LogLikelihood
32 LK = 0;
33 for k = 1:K
34     LK = LK - Pk(k) * N * log( det( Sigma(:, :, k) ) ) / 2;
35     Xminus = X(:, :) - Mu(:, k);
36     invSX = (Sigma(:, :, k) \ Xminus);
37     LK = LK - (1/2) * sum( PKX(:, k)' .* dot(Xminus, invSX, 1) );
38     LK = LK + 0.5 * Pk(k) * N * log(Pk(k));
39     LPKX = -log(PKX(:, k));
40     LPKX(LPKX == inf) = 0;
41     LK = LK + sum( PKX(:, k) .* (LPKX - (M/2) * log(2*pi) ) );
42 end
43
44 end

```

### 4.3.2 Q2

testGP.m

Produces the plot for the periodic gaussian process fit using the functions defined below.

```

1 data = csvread('.../CO2.Mauna.Loa.Data.csv', 1); % load the data in the file
2
3 % For ease of reading
4 X = data(:, 1);
5 X = X/12 + 1960; % Convert to Years
6 y = data(:, 2);
7 sigma.n = sqrt(var(y));
8
9 nx = length(X); % Extract length of data
10 n_xtest = 1e3; % Length of GP extrapolation points over the domain
11 Nstarts = 1; % Num starting points sampled from a hyperbolic distribution to achieve fmin
12
13 ox = 2013; % Year of observation
14 % ox = 645; % Num months between 01/01/1960 and 01/01/2013
15 oy = 400; % ppm level
16
17 % Create a vector fom the smallest x value to the observation point
18 xtest = linspace( min(X), ox, n_xtest);
19 % Compute Hyperparameters
20 theta = computeHP(X, y, Nstarts);
21 % Fit test data to the GP with the found hyperparameters
22 [ytest, bounds, K, kstar] = computeGP (@K.SE, X, y, theta, xtest);
23 % Plot Data
24 display(theta)
25 fig.main = figure(); hold on;
26 plot(ox, oy, 'rx', 'DisplayName', 'Observed value in 2013');
27 plotGP( X, y, xtest, ytest, bounds, 'b-' );
28 ylabel('CO2 in ppm'); xlabel('Year')
29 % axis([min(X) ox 310 400])

```

Computes the covariance matrix K as well as all the derivatives wrt the hyper-parameters

K.SE.m

```

1 function [K, d_l, dsigma.n, dsigma.f, d_f, d_l2] = K.SE(x1, x2, theta)
2 % Computes the covariance matrix K using a squared exponential +
3 % exponential sin squared and white noise kernels as well as all the
4 % derivatives with respect the inputted theta
5
6 if exist('theta', 'var') == 0 % make sure all parameters have been inputted
7     theta = [];
8 end
9 if length(theta) < 5
10     l2 = 1;
11 else
12     l2 = theta(5);
13 end
14 if length(theta) < 4
15     f = 0;
16 else

```

```

17     f = theta(4);
18 end
19 if length(theta) < 3
20     sigma.f = 1;
21 else
22     sigma.f = theta(3);
23 end
24 if length(theta) < 2
25     sigma.n = 0;
26 else
27     sigma.n = theta(2);
28 end
29 if length(theta) < 1
30     l = 1;
31 else
32     l = theta(1);
33 end
34
35 % Covariance
36 K = sigma.f^2 * exp( - (x1-x2)^2 / (2*l^2) );
37 % Derivatives
38 d.l = K * (1^-3) * (x1-x2)^2;
39 dsigma.f = 2 * sigma.f * exp( - (x1-x2)^2 / (2*l^2) );
40 % If the Kernel has any periodic behaviour
41 if f > 0
42     K = K + exp( - 2 * sin( (x1-x2)*f)^2 / l2^2);
43     d.f = exp( - 2 * sin( (x1-x2)*f)^2 / l2^2) * -(4 / l2^2) * sin((x1-x2)*f) * cos((x1-x2)*f) * (x1-x2);
44     d.l2 = exp( - 2 * sin( (x1-x2)*f)^2 / l2^2) * (4*sin((x1-x2)*f)^2 / l2^3);
45 else
46     d.f = 0;
47     d.l2 = 0;
48 end
49 % The white noise terms in the Kernel
50 if x1==x2
51     K = K + sigma.n^2;
52     dsigma.n = 2*sigma.n;
53 else
54     dsigma.n = 0;
55 end
56
57 end

```

## optimHP.m

Based on the implementation in [1] but altered to handle more hyperparameters

```

1 function [nLogLikelihood, nLLDerivatives] = optimHP (logtheta, constants)
2 % Computes the Negative LogLikelihood for the Kernel Function and its
3 % various derivatives
4 theta = exp(logtheta); backupConstants = exp(constants);
5 toBeComputed = constants(6:10); % indicators
6 Ntheta = 1;
7
8 if toBeComputed(1) == 1
9     l1 = theta(1);
10    Ntheta = Ntheta + 1;
11 end
12 if toBeComputed(2) == 1
13    sigma.n = theta(Ntheta);
14    Ntheta = Ntheta + 1;
15 end
16 if toBeComputed(3) == 1
17    sigma.f = theta(Ntheta);
18    Ntheta = Ntheta + 1;
19 end
20 if toBeComputed(4) == 1
21    f = theta(Ntheta);
22    Ntheta = Ntheta + 1;
23 end
24 if toBeComputed(5) == 1
25    l2 = theta(Ntheta);
26    Ntheta = Ntheta + 1;
27 end
28
29 % Constants -- whatever's left
30 if exist('l1') == 0
31    l1 = backupConstants(1);
32 end
33 if exist('sigma.n') == 0
34    sigma.n = backupConstants(2);
35 end
36 if exist('sigma.f') == 0
37    sigma.f = backupConstants(3);
38 end
39 if exist('f') == 0
40    f = backupConstants(4);
41 end
42 if exist('l2') == 0
43    f = backupConstants(5);
44 end
45
46 n = constants(11);
47 X = constants(12:n+11); y = constants(n+12:end); y = y - mean(y);

```

```

48
49 % Compute the Kernel matrix and its derivative matrix wrt to the various
50 % hyper parameters
51 K = zeros(n); dKdl = zeros(n); dKds = zeros(n); dKdf = zeros(n); dKdw = zeros(n); dKdl2 = zeros(n);
52 for i = 1:n
53     for j = 1:n
54         [K(i,j), dKdl(i,j), dKds(i,j), dKdf(i,j), dKdw(i,j), dKdl2(i,j)] = K.SE (X(i),X(j),[l1 sigma_n sigma_f l2]);
55     end
56 end
57
58 % Factorize K using the fact that its positive semi definite
59 L = chol (K,'lower');
60 alpha = L\'(L\y);
61 invK = inv(K);
62 alpha2 = K\y; % (inv(K)*y) alpha from page 114, not page 19, of Rasmussen and Williams (2006)
63
64 % Log marginal likelihood and its gradient
65 logpyX = -y'*alpha/2 - sum(log(diag(L))) - n*log(2*pi)/2;
66 dlogp_dl = 11*trace((alpha2*alpha2' - invK)*dKdl)/2;
67 dlogp_ds = sigma_n*trace((alpha2*alpha2' - invK)*dKds)/2;
68 dlogp_df = sigma_f*trace((alpha2*alpha2' - invK)*dKdf)/2;
69 dlogp_dw = f*trace((alpha2*alpha2' - invK)*dKdw)/2;
70 dlogp_dl2 = 12*trace((alpha2*alpha2' - invK)*dKdl2)/2;
71
72 % Output Negative Log Likelihood and any derivatives of interest
73 nLogLikelihood = -logpyX; nLLDerivatives = [];
74 if toBeComputed(1) == 1
75     nLLDerivatives = [nLLDerivatives -dlogp_dl];
76 end
77 if toBeComputed(2) == 1
78     nLLDerivatives = [nLLDerivatives -dlogp_ds];
79 end
80 if toBeComputed(3) == 1
81     nLLDerivatives = [nLLDerivatives -dlogp_df];
82 end
83 if toBeComputed(4) == 1
84     nLLDerivatives = [nLLDerivatives -dlogp_dw];
85 end
86 if toBeComputed(5) == 1
87     nLLDerivatives = [nLLDerivatives -dlogp_dl2];
88 end
89
90 end

```

## computeGP.m

Computes the GP given training data, hyper-parameters and testing x values.

```

1 function [ytest, bounds, K, ytestvar] = computeGP (k, X, y, theta, xstar)
2 % Finds the mean vector and variances for the testing values when fitting
3 % to a Gaussian Process with parameters theta formed from training data
4 % X and y.
5 % OUTPUTS: ytest, predicted mean vector of the testing data
6 %           bounds, 95% confidence values
7 %           K, The kernel for theta and the training data
8 %           ytestvar, the variance computed for the testing data.
9
10 if nargin < 5, sigma_n = 0; end
11
12 % Compute Kernel
13 ym = mean(y); y = y - ym; % Make data to have zero mean.
14 N = size(X,1); Nxtest = length(xstar);
15 K = zeros(N); kstar = zeros(N,1);
16 for i = 1:N
17     for j = 1:N
18         K(i,j) = k(X(i),X(j),theta);
19     end
20 end
21
22 % Factorize the Kernel for easier computations.
23 L = chol (K,'lower');
24 alpha = L\'(L\y);
25
26 % Compute bounds and estimate for xstar
27 ytestmean = zeros(Nxtest,1); ytestvar = zeros(Nxtest,1);
28 for q = 1:Nxtest
29     for i = 1:N % Compute Kernel components
30         kstar(i) = k(X(i),xstar(q),theta);
31     end
32     ytestmean(q) = kstar' * alpha; % Mean values
33     v = L\kstar;
34     ystar_noise = sigma_n^2;
35     ytestvar(q) = k (xstar(q),xstar(q),theta) - v'*v + ystar_noise;
36 end
37 % Compute bounds to 2 SDs (i.e. 95%)
38 bounds = [ytestmean + 1.96 * sqrt(ytestvar) ytestmean - 1.96 * sqrt(ytestvar)] + ym;
39 ytest = ytestmean + ym;
40
41 end

```

## computeHP.m



Computes the hyper-parameters for a GP given training data.

```

1 function [ theta ] = computeHP( X, y, Nrepeats)
2 % Computes the hyperparameters for inputted data X and y from Nrepeats
3 % OUTPUTS: theta: vector of hyperparameters [l1, sigma_n, sigma_f, f, l2]
4
5 if nargin < 3, Nrepeats = 1; end
6
7 n = length(X);           % Set n
8 sigma_n = sqrt(var(y));   % Set noise variance
9 toBeComputed = [1 1 1 1]'; % 0 = known
10 % Log to prevents negative values
11 constants = [log([1 sigma_n 1 1 1])'; toBeComputed; n; X; y];
12
13 options = optimset('GradObj','on');
14 % Several varied starting points for minimization function
15 l1_samp = hypSample ([30 40], Nrepeats); % for l
16 sf_samp = hypSample ([0.1 2], Nrepeats); % for sigma_f
17 sigmaN_samp = hypSample ([0.1*sigma_n 0.2*sigma_n], Nrepeats); % for sigma_n
18 f_samp = hypSample ([2 4], Nrepeats); % for frequency
19 l2_samp = hypSample ([0.9 1], Nrepeats); % for l
20 startVals = log ([l1_samp sigmaN_samp sf_samp f_samp l2_samp]);
21
22 % Loop over the random starting values to find minimum negative log
23 % likelihood.
24 thetaVec = [];
25 for s = 1:Nrepeats
26     [theta, fval] = fminunc(@(hyperparameters) optimHP(hyperparameters, constants), startVals(s,:), options);
27     thetaVec = [thetaVec; fval startVals(s,:) theta];
28 end
29 thetaVec(:,5:end) = exp(thetaVec(:,5:end)); % Return data to non log format
30 thetaVec(:,1) = thetaVec(:,1)/max(abs(thetaVec(:,1))); % Normalize
31 thetaVec = sortrows(thetaVec); % Sort so smallest is first
32 theta = thetaVec(1,size(startVals,2)+2:end); % Choose smallest
33
34 end

```

plotGP.m

creates plots of the computed Gaussian Process

```

1 function [ ] = plotGP( X, y, xtest, ytest, bounds, linestyle )
2 % Plots the data for a Gaussian Process.
3 % INPUTS: X, y: Training Data
4           xtest, ytest: Fit of the Testing Data
5           bounds: Upper and lower 2SD bounds on the testing data
6           linestyle: Can decide the linestyle of the training data.
7 % OUTPUTS: Plot.
8
9 if nargin < 6
10     linestyle = 'b+' ;
11 end
12
13 hold on
14
15 stdRegion(xtest,bounds);
16 plot (xtest,ytest,'c','DisplayName','Best Estimate');
17 plot (X,y,linestyle,'DisplayName','Data Points');
18 legend('show');
19 end

```

stdRegion.m

From [1]. Creates shaded bounds around the mean

```

1 % Fill a graph with standard deviations; t is nx1 and Range is nx2
2 % sC is the colour to use, in RGB vector format
3 % reaxisVar = 0 to leave the axes to be overdrawn, or 1 to redraw them
4 % Mark Ebden, 2008
5
6 function theRange = stdRegion (t, theRange, sC, reaxisVar)
7 if exist('sC') == 0
8     sC = [.6 .6 .6];
9 end
10 if exist('reaxisVar') == 0
11     reaxisVar = 0;
12 end
13 t = t(:);
14 % Create enclosed shape using bounds
15 fill ([t; flipud(t)], [theRange(:,1); flipud(theRange(:,2))], sC, 'EdgeColor', sC,'DisplayName','Bounds on GP');
16
17 if reaxisVar == 1
18     v = axis;
19     line (v([1 1]), v([3 4]), 'Color', 'k');
20     line (v([1 2]), v([3 3]), 'Color', 'k');
21 end

```

HypSample.m

From [10]. Samples a hyperbola N times for a parameter with given bounds

```

1 % Sample a hyperbola N times for a parameter bounded by the two components of 'bounds'.
2 % Example usage: if p(x) = k/x from 0.1 to 4, and p(x) = 0 otherwise:
3 %               x = hypSample ([.1 4], 1e5); hist(x,100)
4 % Bounds can be calculated as per Section 5.5.1 in 'Data Analysis: A Bayesian Tutorial',
5 % 2nd edition, D.S. Sivia and J. Skilling, 2006.
6
7 function x = hypSample (bounds, N)
8
9 xmin = bounds(1); xmax = bounds(2);
10 F = rand(N,1);
11 x = xmin.*(1-F) .* xmax.*F;
12
13 end

```

### 4.3.3 Q4

testingGeneral.m

for generating the data visualisations at the start of Q4.

```

1 % Load and format data
2 load('bitcoinData.mat');
3 bid = string(bid);
4 bid = (bid == 'TRUE');
5 symbol = string(symbol);
6 exchange = string(exchange);
7
8 N = length(price); % get length of data
9
10 % Plot the prices the transactions at bid and ask price occurred
11 fig.prices = figure(); hold on;
12 plot(date1(bid == 1), price(bid == 1), 'rx', 'DisplayName', 'Bid Prices');
13 plot(date1(bid == 0), price(bid == 0), 'g.', 'DisplayName', 'Ask Prices');
14 xlabel('Date Value'); ylabel('Price'); legend('show');
15 title('Bid and Ask Prices for Bitcoin Transactions')
16
17 U = unique(date1); % Unique Days
18 NU = length(U); % # Unique Days
19 NpDay = zeros(NU,1); % N per Day
20 for day = 1:NU
21     NpDay(day) = length(date1(date1 == U(day)));
22 end
23
24 Bdate = date1(bid == 1); % Dates of Bid transactions
25 Bprice = price(bid == 1); % Price of Bid transactions
26 Bamount = amount(bid == 1); % Amount of Ask transactions
27 Adate = date1(bid == 0); % Dates of Ask transactions
28 Aprice = price(bid == 0); % Price of Ask transactions
29 Aamount = amount(bid == 0); % Amount of Ask transactions
30
31 % fig.next = figure(); hold on;
32 % subplot(1,2,1); plot(Bprice);
33 % subplot(1,2,2); plot(Aprice);
34
35 VpDay = zeros(NU,2);
36 AvgpDay = zeros(NU,2);
37
38 for day = 1:NU
39     dayAmount = Bamount(Bdate == U(day));
40     dayPrice = Bprice(Bdate == U(day));
41     VpDay(day,1) = sum(dayAmount);
42     AvgpDay(day,1) = sum(dayAmount.*dayPrice) / VpDay(day,1);
43
44     dayAmount = Aamount(Adate == U(day));
45     dayPrice = Aprice(Adate == U(day));
46     VpDay(day,2) = sum(dayAmount);
47     AvgpDay(day,2) = sum(dayAmount.*dayPrice) / VpDay(day,2);
48 end
49
50 % fig.next = figure(); hold on;
51 plot(U, AvgpDay(:,1), 'x-r');
52 plot(U, AvgpDay(:,2), 'x-g');
53
54 figure(); hold on;
55 plot(id(bid == 1), price(bid == 1), 'rx');
56 plot(id(bid == 0), price(bid == 0), 'gx');
57 xlabel('Price'); ylabel('ID'); title('Price (coloured in bid/ask) vs ID');

```

testingRegression.m

For generating the data regression plot visualisations

```

1 % Load Data

```

```

2 load('bitcoinData.mat');
3 bid = string(bid);
4 bid = (bid == 'TRUE');
5 symbol = string(symbol);
6 exchange = string(exchange);
7
8 X = date1;
9 y = price;
10
11 % Extract length of data
12 nx = length(X);
13 % Confidence interval size: 2SD = 95% for normal distribution
14 alpha = 0.05;
15
16 % Fit data
17 [ y_OLS, bounds, ~ ] = OLS(X, y, alpha);
18
19 % Plot data
20 figure(); hold on;
21 % plotRegression(X,y,X,y_OLS,bounds,'bx');
22 xlabel('Date'); ylabel('Price');
23 title('Simple Linear Regression with 2 standard deviation bounds');
24
25 U = unique(date1); % Unique Days
26 NU = length(U); % # Unique Days
27 NpDay = zeros(NU,1); % N per Day
28 for day = 1:NU
29     NpDay(day) = length(date1(date1 == U(day)));
30 end
31
32 Bdate = date1(bid == 1); % Dates of Bid transactions
33 Bprice = price(bid == 1); % Price of Bid transactions
34 Bamount = amount(bid == 1); % Amount of Ask transactions
35 Adate = date1(bid == 0); % Dates of Ask transactions
36 Aprice = price(bid == 0); % Price of Ask transactions
37 Aamount = amount(bid == 0); % Amount of Ask transactions
38
39 % fig.next = figure(); hold on;
40 % subplot(1,2,1); plot(Bprice);
41 % subplot(1,2,2); plot(Aprice);
42
43 VpDay = zeros(NU,2);
44 AvgpDay = zeros(NU,2);
45
46 for day = 1:NU
47     dayAmount = Bamount(Bdate == U(day));
48     dayPrice = Bprice(Bdate == U(day));
49     VpDay(day,1) = sum(dayAmount);
50     AvgpDay(day,1) = sum(dayAmount.*dayPrice) / VpDay(day,1);
51
52     dayAmount = Aamount(Adate == U(day));
53     dayPrice = Aprice(Adate == U(day));
54     VpDay(day,2) = sum(dayAmount);
55     AvgpDay(day,2) = sum(dayAmount.*dayPrice) / VpDay(day,2);
56 end
57
58 AvgpDayT = sum(VpDay.*AvgpDay,2)./sum(VpDay)';
59
60 [ y_OLS.dayAVG, bounds, ~ ] = OLS(U, AvgpDayT, alpha);
61
62 % Plot data
63 plotRegression( U, AvgpDayT, U, y_OLS.dayAVG, bounds, 'bx');
64 plot(X,y,'k');

```

## OLS

### For generating the data regression plot visualisations

```

1 function [ fitted_data, bounds, beta_coeffs, sigma ] = OLS(x, y, alpha)
2 %OLS – Ordinary Least Squares
3 % Detailed explanation goes here
4
5 if nargin < 3, alpha = 0.05; end
6
7 % Extract length of the data of interest
8 N = length(x);
9
10 xm = mean(x);
11 ym = mean(y);
12
13 betal = sum((x - xm).*(y-ym));
14 betal = betal / sum( (x - xm).^2 );
15 beta0 = ym - betal * xm;
16 beta_coeffs = [beta0, betal];
17 % Find the estimated values
18 fitted_data = beta_coeffs(2) * x + beta_coeffs(1);
19
20 r = fitted_data - y;
21 SSe = sum(r.^2); % Squared sum of errors
22 var_e = SSe/(N-2); % Variance of the error
23 sigma = sqrt(var_e);
24 tval = tinv(1 - alpha/2, N-2); % Student t value with given degrees of freedom and (1-alpha) Confidence interval
25 c = sqrt( 1 + 1/N + ((x - xm).^2) / sum( (x - xm).^2 ) );

```

```

26 bound = tval * sigma * c;
27
28 bounds = [fitted_data + bound, fitted_data - bound ];
29
30 end

```

## NeuralNetwork.py

This file contains some of the initial python plots and some basic analysis of the system using iterative validation

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 #####
4
5 def buildModel(NIL,NHL,NOL):
6
7     # Initialize matrices and arrays to random values
8     np.random.seed(0) # For testing purposes
9     W1 = np.random.randn(NHL, NIL) / np.sqrt(NIL)
10    c1 = np.zeros((NHL,1))
11    W2 = np.random.randn(NOL, NHL) / np.sqrt(NHL)
12    c2 = np.zeros((NOL,1))
13
14    return { 'W1': W1, 'c1': c1, 'W2': W2, 'c2': c2}
15 #####
16
17 def trainModel(model, X, y, repeats = 1):
18     W1, c1, W2, c2 = model['W1'], model['c1'], model['W2'], model['c2']
19     # Compute solution and then propagate error back
20     for i in range(repeats):
21         # Forward propagation
22         u1 = W1.dot(X) + c1
23         v1 = np.tanh(u1)
24         u2 = W2.dot(v1) + c2
25         exp_scores = np.exp(u2)
26         probs = exp_scores / np.sum(exp_scores)
27         # Backpropagation
28         delta3 = probs - y
29         dW2 = delta3.dot(v1.T)
30         dc2 = np.sum(delta3)
31         delta2 = (W2.T).dot(delta3) * (1 - np.power(v1, 2))
32         dW1 = np.dot(delta2,X.T)
33         dc1 = np.sum(delta2)
34         # Compute level of error to add
35         dW2 += eps2 * W2
36         dW1 += eps2 * W1
37         # Add Error to the model matrices/arrays
38         W1 += -eps1 * dW1
39         c1 += -eps1 * dc1
40         W2 += -eps1 * dW2
41         c2 += -eps1 * dc2
42     return { 'W1': W1, 'c1': c1, 'W2': W2, 'c2': c2}
43 #####
44
45 # Helper function to predict an output (0 or 1)
46 def predict(model, X):
47     W1, c1, W2, c2 = model['W1'], model['c1'], model['W2'], model['c2']
48     # Forward propagation
49     u1 = W1.dot(X) + c1
50     v1 = np.tanh(u1)
51     u2 = W2.dot(v1) + c2
52     exp_scores = np.exp(u2)
53     probs = exp_scores / np.sum(exp_scores)
54     return probs
55 #####
56 def moving_average(x, n) :
57     ret = np.cumsum(x, dtype=float)
58     ret[n:] = ret[n:] - ret[:-n]
59     for i in range(n-1):
60         ret[i] /= (i+1)
61     ret[n-1:] /= n
62     return ret
63 # Cite: http://stackoverflow.com/users/110026/jaime
64 #####
65
66 def testResult(model, Xtest, ytest):
67     # Specefic to this situation. Can be changed for different models
68     probs = predict(model, Xtest)
69     if np.argmax(probs) == np.argmax(ytest):
70         return 1
71     else:
72         return 0
73 #####
74 def trainVals(model, start, num):
75     for i in range(start, start+num):
76         vals = np.arange(i-Ntau+1,i+1) # indices of prices of interest
77         Xt = np.array([price[vals].T, bid[vals].T])
78         Xt = Xt.reshape(NIL,1) # Make into a vector
79         yt = np.array([priceUp[i], 1 - priceUp[i]]) # Vector of [up, down] booleans
80         yt = yt.reshape((2,1))
81         # Train the model of this value

```

```

82     model = trainModel(model, Xt, yt)
83
84     return model
85
86 #####
87 def testVals(model, start, num):
88     correct = 0
89     for i in range(start, start+num):
90         vals = np.arange(i-Ntau+1,i+1) # indices of prices of interest
91         Xt = np.array([price[vals].T, bid[vals].T])
92         Xt = Xt.reshape((NIL,1)) # Make into a vector
93         yt = np.array([priceUp[i], 1 - priceUp[i]]) # Vector of [up, down] booleans
94         yt = yt.reshape((2,1))
95         # Compute if the result is correct
96         correct += testResult(model, Xt, yt)
97     return correct / float(num) # Return % of correct values
98
99 #####
100 if __name__ == '__main__':
101     # Load and format the data
102     my_data = np.genfromtxt('data.csv', delimiter = ',', dtype = None)
103
104     ID = my_data[:,0]
105     ID = ID.astype(np.float)
106     ID = ID[np.argsort(ID)]
107
108     date = my_data[:,3]
109     date = date.astype(np.float)
110     date = date[np.argsort(ID)]
111
112     price = my_data[:,4]
113     price = price.astype(np.float)
114     price = price[np.argsort(ID)]
115
116     amount = my_data[:,5]
117     amount = amount.astype(np.float)
118     amount = amount[np.argsort(ID)]
119
120     bid = my_data[:,6]
121     bid[bid == 'TRUE'] = 1
122     bid[bid == 'FALSE'] = 0
123     bid = bid.astype(np.float)
124     bid = bid[np.argsort(ID)]
125
126     # Comment out the statements with np.argsort to see the impact
127     # plt.plot(ID,price, label = 'Bitcoin Price Data')
128     # plt.legend(loc='upper right')
129     # plt.xlabel('Transaction ID')
130     # plt.ylabel('Price')
131     # plt.show()
132
133     # Different definition of if the price went up.
134     # 1. Firstly if the next transaction value is higher or lower
135     # Not very useful in terms of understanding anything because of financial fluctuations
136     N = np.size(price)
137     D1 = price[1:N] - price[:N-1] # Level 1 different in price
138     priceUp1 = np.ones(N-1)
139     priceUp1[D1 < 0] = 0 # I.e. if the price went up or down
140     # 2. See if the average price in the next avgSize transactions is higher or lower
141     avgSize = 50 # Chosen arbitrarily
142     movingAvg = moving.average(price, avgSize)
143     ND = avgSize
144     D = movingAvg[ND:] - movingAvg[:ND]
145     priceUp = np.ones(N-ND)
146     priceUp[D < 0] = 0
147     # Plot to visualize
148     # plt.plot(ID,price, label = 'Bitcoin Price Data')
149     # plt.plot(ID,movingAvg, '--r', label='Moving Average (50 Transaction)')
150     # plt.legend(loc='upper right')
151     # plt.xlabel('Transaction ID')
152     # plt.ylabel('Price')
153     # plt.show()
154
155     Ntau = 100
156     dim_X = 2
157     NX = Ntau * dim_X
158
159     # Parameters for Neural Network
160     eps1 = 1e-6 # learning rate (chosen)
161     eps2 = 1e-6 # regularization strength (chosen)
162     NIL = NX # input layer size
163     NHL = NIL
164     NOL = 2 # output layer size
165
166     # Initialize matrices and vectors for the model
167     model = buildModel(NIL,NHL,NOL)
168
169     # Basic testing if trained on roughly half the values is prediction on the
170     # rest of the data set
171     model = trainVals(model, Ntau,10000)
172     print testVals(model, Ntau+10000, N-2*Ntau-10000)
173
174     # Iterative cross validation but without randomness dealt with.
175     N_samples = 4000
176     N_runs = N / N_samples

```

```

177 # Train model on initial data set
178 model = trainVals(model, Ntau, N.samples)
179 # Store model performance in each validation set
180 performance = np.array([])
181 # Test model on the next N.sample values then train and repeat
182 for run in range(1,N.runs):
183     if run == N.runs: # Make sure all values are tested in the last batch
184         N.samples += N - N.runs * N.samples
185     # Compute performance
186     performance = np.append(performance, testVals(model, Ntau + run*N.samples, N.samples))
187     if run != N.runs: # If not the last run then train the model further
188         model = trainVals(model, Ntau + run*N.samples, N.samples)
189 print performance
190 print np.average(performance)
191
192 plt.figure()
193 plt.plot(ID, price)
194 for run in range(1,N.runs):
195     midval = run*N.samples+ N.samples/10
196     plt.annotate( str( performance[run-1]) , ( ID[midval], 680))
197     plt.plot((ID[run*N.samples], ID[run*N.samples]), (610, 690), 'k-')
198 plt.xlabel('Transaction ID')
199 plt.ylabel('Price')
200 plt.show()

```

## NeuralNetwork.py

### Python file for K-folds cross validation

```

1 perfs = np.array([])
2 K.fold = 10
3 # Set number of values in each fold
4 N.samples = (N - Ntau - avgSize) / K.fold
5 for K in range(K.fold):
6     # Initialize matrices and vectors for the model
7     model = buildModel(NIL,NHL,NOL)
8
9     # Train model on 9 of the 10 data sets
10    for run in np.random.permutation(K.fold):
11        if run != K:
12            model = trainVals(model, Ntau + run * N.samples, N.samples)
13        # Test on the Kth data set
14        performance = testVals(model, Ntau + K * N.samples, N.samples)
15        perfs = np.append(perfs, performance)
16 print perfs
17 print np.average(perfs)
18
19 plt.figure()
20 plt.plot(ID, price)
21 for K in range(K.fold):
22     midval = K*N.samples+ N.samples / 10
23     plt.annotate( str.format('{0:3f}', perfs[K]) , ( ID[midval], 680))
24     plt.plot((ID[K*N.samples], ID[K*N.samples]), (610, 690), 'k-')
25 plt.plot((ID[N-1], ID[N-1]), (610, 690), 'k-')
26 plt.xlabel('Transaction ID')
27 plt.ylabel('Price')
28 plt.show()

```

## NeuralNetwork.py

### Python file for Iterative Batch Cross Validation with multiple runs

```

1 perfs = np.array([])
2 N.random = 10
3 for rand in range(N.random):
4     # Initialize matrices and vectors for the model
5     model = buildModel(NIL,NHL,NOL)
6
7     # Iterative cross validation:
8     N.samples = 2000
9     N.runs = N / N.samples
10    # Train model on initial data set
11    model = trainVals(model, Ntau, N.samples)
12    # Store model performance in each validation set
13    performance = np.array([])
14    # Test model on the next N.sample values then train and repeat
15    for run in range(1,N.runs):
16        if run == N.runs: # Make sure all values are tested in the last batch
17            N.samples += N - N.runs * N.samples
18        # Compute performance
19        performance = np.append(performance, testVals(model, Ntau + run*N.samples, N.samples))
20        if run != N.runs: # If not the last run then train the model further
21            model = trainVals(model, Ntau + run*N.samples, N.samples)
22    perfs = np.append(perfs, performance)
23
24 perfs = np.reshape(perfs, (N.runs-1,N.random))
25 avgPerfs = np.average(perfs, axis = 1)
26 TotalPerf = np.average(avgPerfs)
27 print TotalPerf
28

```

```

29 N_samples = 2000
30 plt.figure()
31 plt.plot(ID, price)
32 for run in range(N_runs-1):
33     midval = (run+1)*N_samples+ N_samples/10
34     plt.annotate( str.format( '{0:.3f}', avgPerfs[run] ) , ( ID[midval], 680))
35     plt.plot((ID[(run+1)*N_samples], ID[(run+1)*N_samples]), (610, 690), 'k-')
36 plt.xlabel('Transaction ID')
37 plt.ylabel('Price')
38 plt.show()

```

## References

- [1] M. Ebden. Gaussian Processes: A Quick Introduction. *ArXiv e-prints*, May 2015.
- [2] C.E. Rasmussen and C.K.I. Williams. *Gaussian Processes for Machine Learning*. Adaptive computation and machine learning series. University Press Group Limited, 2006. ISBN 9780262182539.
- [3] Tesla self-driving demonstration. URL [https://www.tesla.com/en\\_GB/videos/autopilot-self-driving-hardware-neighborhood-short](https://www.tesla.com/en_GB/videos/autopilot-self-driving-hardware-neighborhood-short). Accessed: March 2017.
- [4] Google waymo. URL <https://waymo.com/press/>. Accessed: March 2017.
- [5] Sebastian Thrun, Mike Montemerlo, Hendrik Dahlkamp, David Stavens, Andrei Aron, James Diebel, Philip Fong, John Gale, Morgan Halpenny, Gabriel Hoffmann, Kenny Lau, Celia Oakley, Mark Palatucci, Vaughan Pratt, Pascal Stang, Sven Strohband, Cedric Dupont, Lars-Erik Jendrossek, Christian Koelen, Charles Markey, Carlo Rummel, Joe van Niekerk, Eric Jensen, Philippe Alessandrini, Gary Bradski, Bob Davies, Scott Ettinger, Adrian Kaehler, Ara Nefian, and Pamela Mahoney. Stanley: The robot that won the darpa grand challenge. *Journal of Field Robotics*, 23(9):661–692, 2006. ISSN 1556-4967. doi: 10.1002/rob.20147. URL <http://dx.doi.org/10.1002/rob.20147>.
- [6] Alberto Broggi, Claudio Caraffi, Pier Paolo Porta, and Paolo Zani. The single frame stereo vision system for reliable obstacle detection used during the 2005 darpa grand challenge on terramax. In *Intelligent Transportation Systems Conference, 2006. ITSC'06. IEEE*, pages 745–752. IEEE, 2006.
- [7] G Pasaoglu, D Fiorello, A Martino, G Scarcella, A Alemanno, A Zubaryeva, and C Thiel. Driving and parking patterns of european car drivers-a mobility survey. 2012.
- [8] Aurangzeb Khan, Baharum Baharudin, Lam Hong Lee, and Khairullah Khan. A review of machine learning algorithms for text-documents classification. *Journal of advances in information technology*, 1(1):4–20, 2010.
- [9] Mckinsey: Capturing the advanced driver-assistance systems opportunity. URL <http://www.mckinsey.com/industries/automotive-and-assembly/our-insights/capturing-the-advanced>. Accessed: March 2017.
- [10] D. Sivia and J. Skilling. *Data Analysis: A Bayesian Tutorial*. Oxford science publications. OUP Oxford, 2006. ISBN 9780198568315.