# Analysis
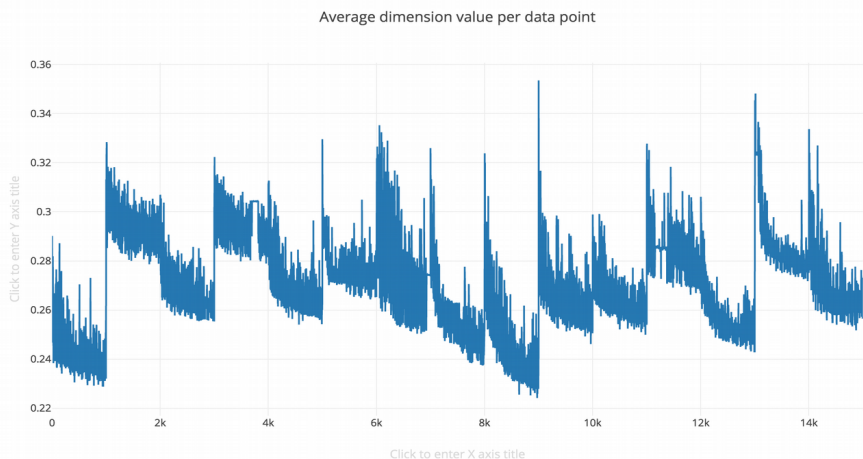
motivation(1): make CSV files that can be read easily into python or excel spread sheet
method: write a string manipulation program called reformat_data.py, get rid of the redundant
characters, such as 'qid:WT04:' , "1:" etc...



| QID | n | no. 1s | no. 0s |
|---|---|---|---|
| 170 | 1000 | 3 | 997 |
| 176 | 1000 | 13 | 987 |
| 177 | 1000 | 23 | 977 |
| 178 | 1000 | 11 | 989 |
| 185 | 1000 | 10 | 990 |
| 189 | 1000 | 14 | 986 |
| 191 | 1000 | 13 | 987 |
| 192 | 1000 | 1 | 999 |
| 194 | 1000 | 32 | 968 |
| 197 | 1000 | 3 | 997 |
| 200 | 1000 | 26 | 974 |
| 204 | 1000 | 13 | 987 |
| 207 | 1000 | 7 | 993 |
| 209 | 1000 | 12 | 988 |
| 221 | 1000 | 9 | 991 |
| TOTAL | 15000 | 190 | 14810 |

Observation: Preliminary observations show that the Query's are organized into 1000 documents
each. All together 15000 documents output, on 15 different query's. Relatively few relevant
documents across the query's. Qid:192. Only has one relevant documets.

Motivation: plot to see obvious overall structure.
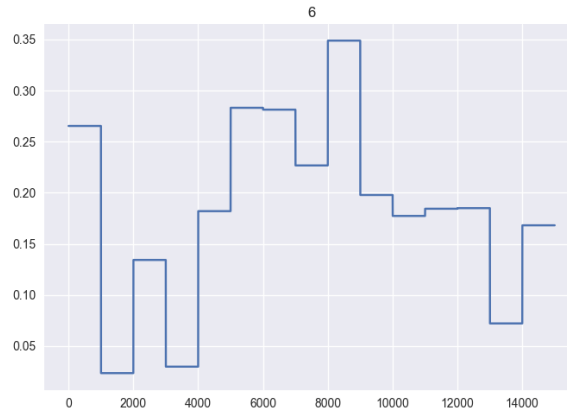Method: plots the horizontal mean vs document indexes



Observation: Quite obvious mean structural change between querys, similar level of signal strength
across the querys, with the 7th query (qid:191) show noticeably high strength. A curious decay in

mean across the documents indexing within each query. This suggest that the documents indexing has a informative ordering.
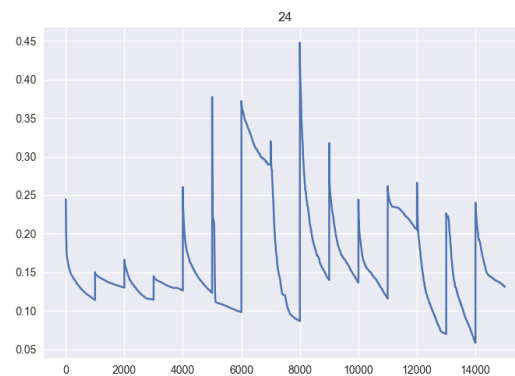
Motivation: find out the source of this decay, and mean structure change, and documents signal
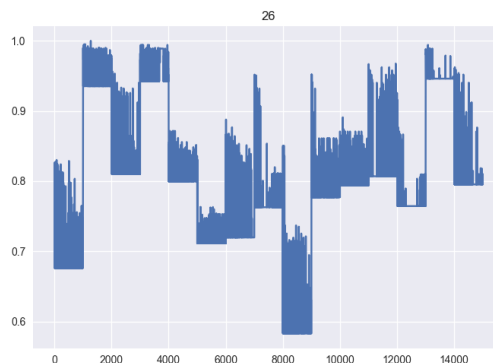Method: plot the individual features vs documents
Observations: By looking at the 64 plots, it looks like we can seperate into 6 class
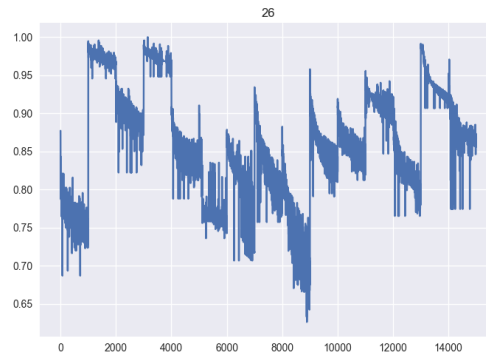


Class1: [6,7,8,9,10] these features show zero document signal, changes only when query changes. These are certainly query based only features.
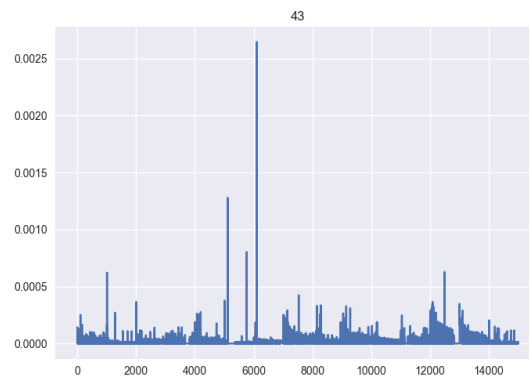


Class2: [25] This singular feature, capture no documents signal, respond to query change only and has a decay across the document indexing. We can use this feature to learn the decay rate of signal strength for each Querys or the average decay rate of all Querys. This learnt decay rate can be used to offset the signal decay of other features.
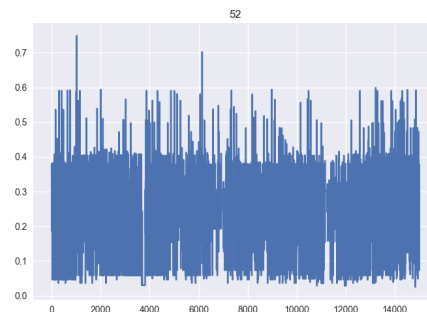


Class3: [27,28,29,32,33,34,37,38,39,45,61,62,63,64] Captures everything and no decay. These are likely to be the most informative assessment features for the interactions between query and documents in-order to learn relevance.

Class4: [21,26,30,31,35,36,40,41,42,44,46,47,48] captures everything and included the decay. We can offset this decay by the decay rate learnt from class 2.
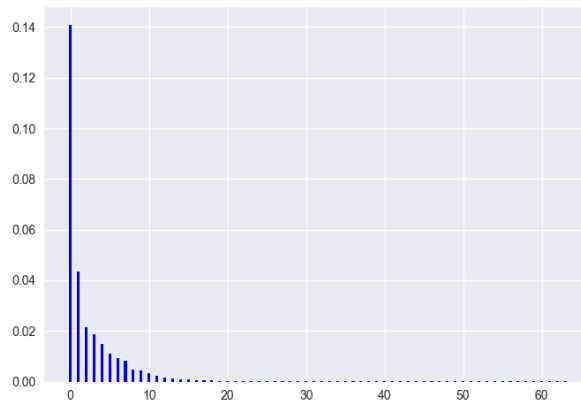


Class5: [1,2,5,11,12,13,14,15,16,20,43,49,50,51,53,54,55,56,60] These features captures the documents, but only very minor level of decay and query (if at all)
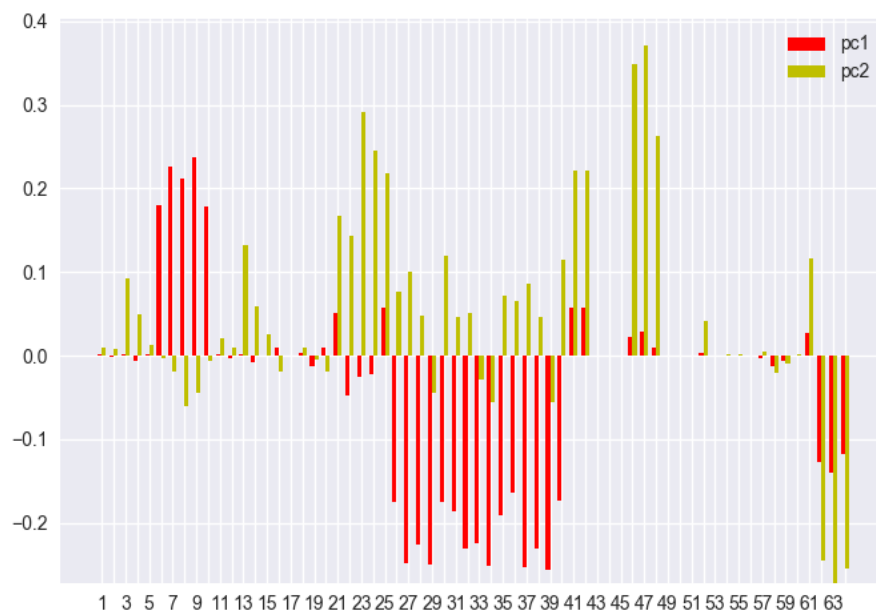


Class6: [3,4,17,18,19,22,23,24,52,57,58,59]. These features doesn't seem to capture any query information or decay at all. It is very likely that these are documents based only features. However, the precise set member of this class are harder to tell than class 1, since every documents are unique, we can't hold the documents constant and check if they really respond to query for sure.

Observation: It is very possible that class member of Class 5 and Class 6 get mixed up. In addition, we can deduce that, of the features that is based on both Documents and Querys, with the exception of class 3; the more they react to query, the more they are affected by decay.
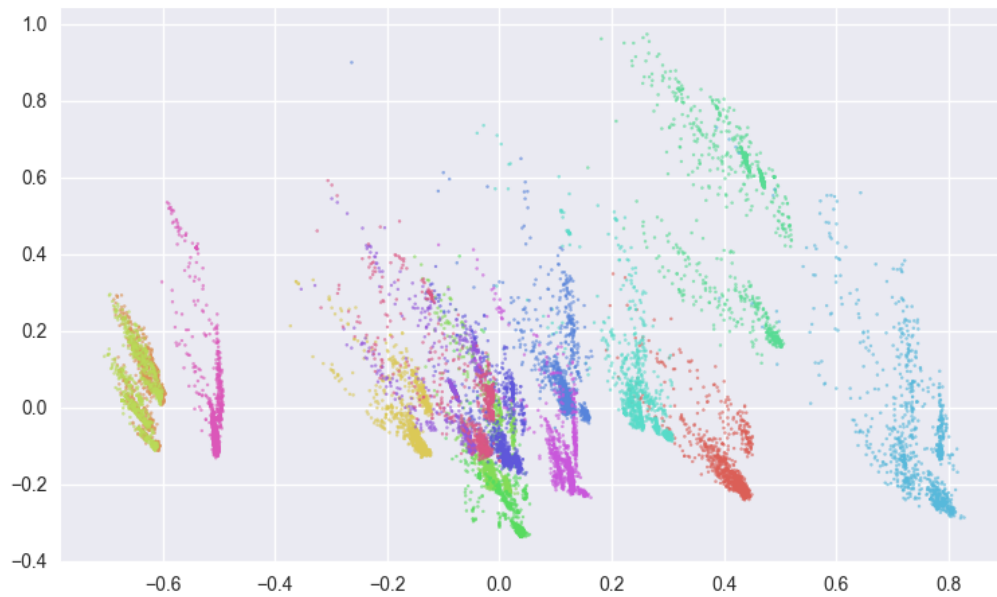
Motivation: Use Dimensionality reduction to visualize data
Method: PCA



***These are the eigenvalues of the 64 PCs***: The first 10 principle component almost explains almost all the variance. The First and $2^{nd}$ principle components combined is sufficient to explain 63% of all the variance.
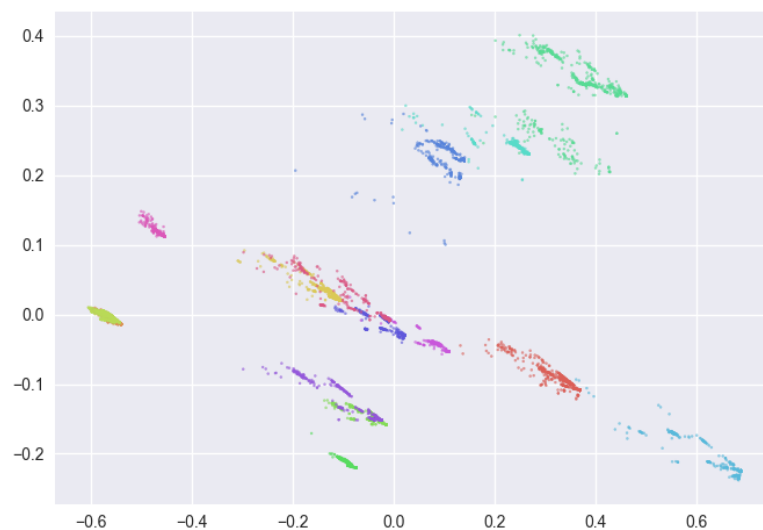


***The direction of projection of the $1^{st}$ and $2^{nd}$ principle component.*** As expected from previous plots, class 1,2,3,4 has dominant presence. With Class 1,2,3 mostly in the PC1, and class 4 mostly in PC4.
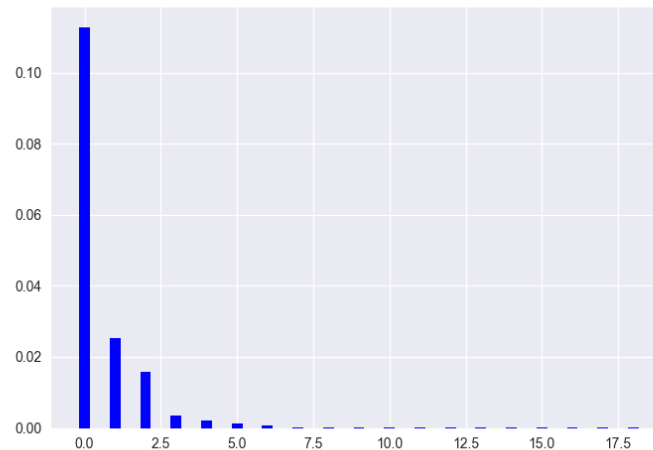
***Scatterplot of all data points projected into a plain span by PC1 (x-axis) and PC2 (y-axis).*** Notice that qid176 and qid178 is almost identical in shape and position. And the Furtherest away.

Motivation: Lets focus on the features most closely related to Qids, ie: class1 (query based only without decay), class3 (strong query based relative to DocBased without decay). To to better differentiate between the Qids.
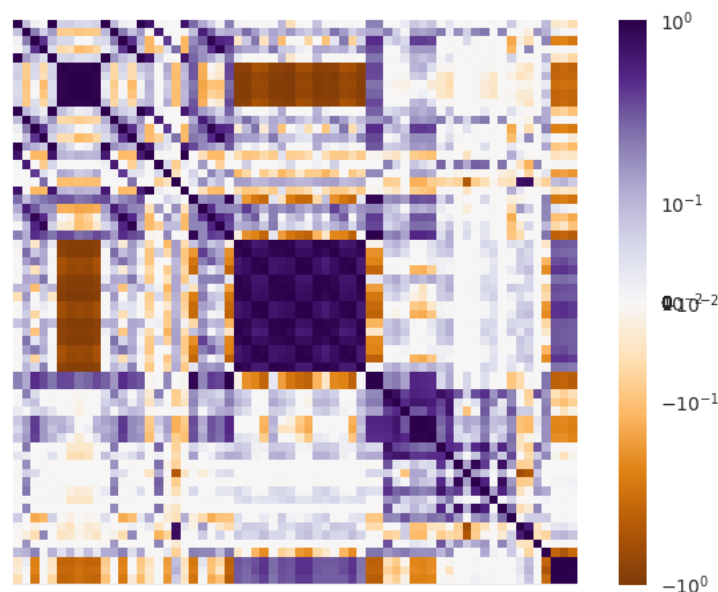
Incase, if I want to do query specific supervised learning and there isn't sufficient relevancy count, I can group them based on the above clustering.
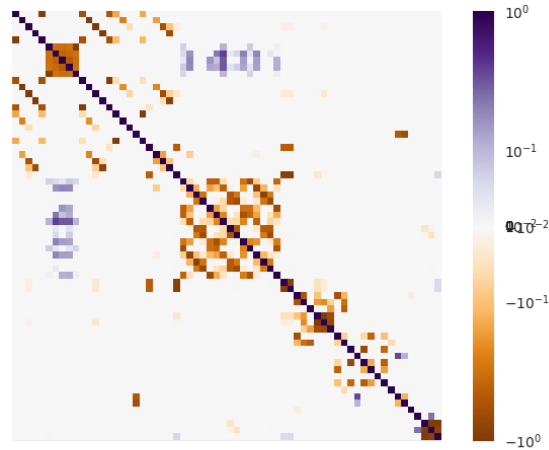


***The eigenvalue decay from 1st PC on the left to 19th PC***. With features from class1 and class 3 only. Observation: Only 3 orthogonal component is enough to describe the qids type. With 1st PCs contributing to 96% of all the variances.

Motivation: discover more relationship betweens features
Methods: Normalized the entire feature spaces, and make covariance matrix and then use GraphicalLasso to discover sparsed percision matrix.
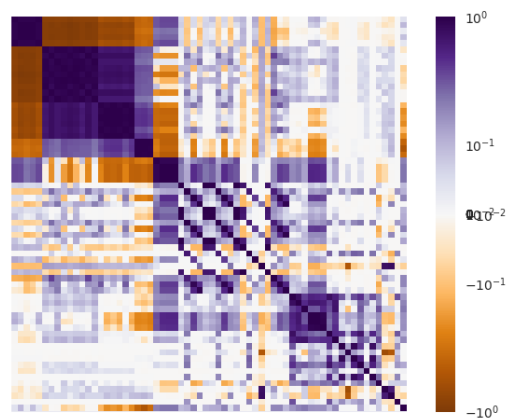


***Covariance matrix:*** Observation: Since we know that 6,7,8,9,10. Belong to class 1 and has the cleanest query based features we can deduce from that, block [25 to 39] has the closest interaction with the query (which made sense because they are majority class 3) Follow with block [61 to 64]. Notice that there is a checker board pattern in black 25 to 39. This mean that they are ordered in an interesting phase oscillation pattern, and they can be further divided into 2 subgroub. Subgroub 1: [25,26,29,30,31,34,35,36,39], subgroup 2: [27,28,32,33,37,38]. Subgroup 1 has slightly stronger tide with the query than subgroup 2. They each have stronger relationship with features within the subgroup, less but still strong with the other subgroup.
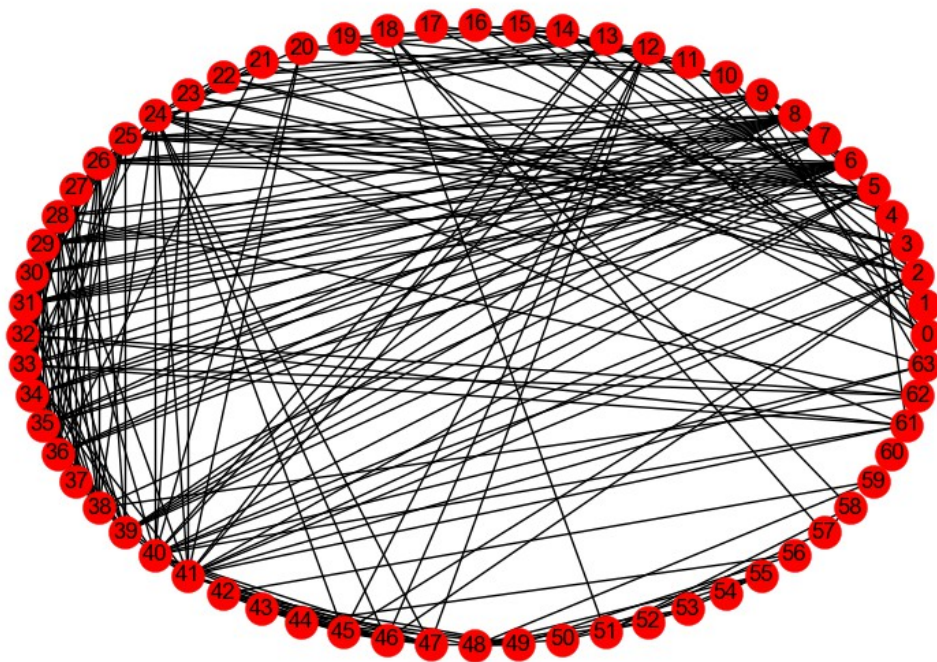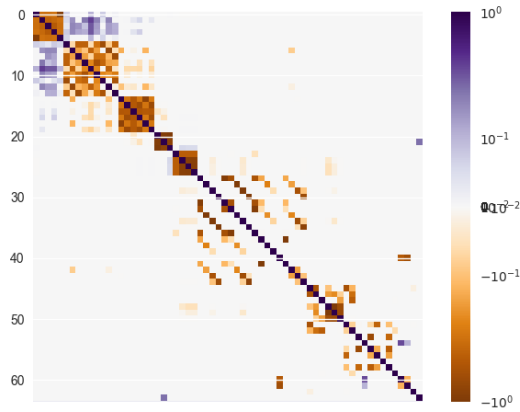
***Sparsed precision matrix from Graphical Lasso:*** Observation: Notice that the there a repetitive pattern between in a period of 5 indexes between [10,11,12,13,14] →[15,16,17,18,19]-→ [20,21,22,23,24]. This could be another oscillatory relationship therefore related subgroups. Also notice that block 61-64 is now conditionally independent to block 6 to10 and the big middle block. Therefore 61-64 could be capturing the redundant part of the query. Also, block[49-59] do is conditionally independent on almost everything else. This could be a strong indicator that they are capturing the document based only information.

Motivation: Move the coulombs around to better visualize relationship
Methods: Manually move block [6-10] to left most. Follow with subgroup from the middle block, [25,26,29,30,31,34,35,36,39], then [27,28,32,33,37,38], and then [61-64].
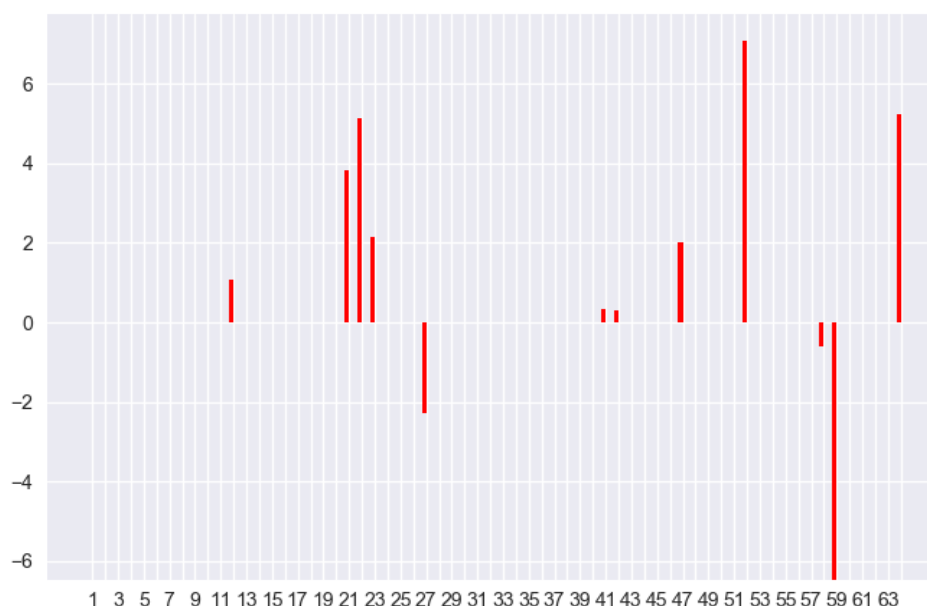
Above is a network built by the Precision Matrix as a weighted adjacency matrix. If given more time. I could use Hierarchical clustering with linkage to better group the features.

Motivation: Try to predict the relevancy by the features. See which features contributes the most to a documents relevancy.

Method: Logistic regression using 'L₁' penalty to force sparsity in the prediction co-efficient. Use 5-fold cross validation to tune the regularization parameter so as to not overfit or under-fit.



Above shows the co-efficient of the features in the decision function. Learnt by the logistic regression solver. It has the average score of 0.98 according to my 5-fold cross validation algorithm. However, this may not be a high score at all since the there are far too many zero's and not enough ones in the relevancy column. Interestingly, the majority of the dominant predictor are what I considered class 5 or class 6; the classes of features which I believe are almost all documents based only. This could mean that, the Queries given in this data set are not sufficiently different to make a difference to the probability of relevancy.

If I have more time, I will try to use a different function approximator to predict relevancy. I would approximate the decay using feature 25, use it to mediate the decay found in class 2,4,5 before using any supervised learning algorithm. I would try centering and normalizing the data set before applying the learning algorithm. I would use Hierarchical clustering to better group the features.