# HTML homework 5: q11 report

The resulting picture is as follows:

**Distribution of Eout**



In order to compare with the previous question, some statistics are:

```
Statistics for Eout:
Mean: 0.0150
Median: 0.0151
Standard Deviation: 0.0010
Min: 0.0126
Max: 0.0191
```

```
Statistics for Eout:
Mean: 0.0152
Median: 0.0151
Standard Deviation: 0.0010
Min: 0.0126
Max: 0.0181

Statistics for Non-zero Components:
Mean: 546.5
Median: 563.0
Standard Deviation: 27.2
Min: 484
Max: 576
```

[1]

---

[1] Left: statistics of q11 result; Right: statistics of q10 result

In this problem, the difference from the previous is that we further divide the training set into a subtrain set (with 8000 exampls), and a validation set (3876 examples).

From the statistics, we can see that most of them are quite similar, I think part of the reason is that, even though 3876 examples are taken out from the original training set, so one may expect that as we train on a smaller subtrain set, the out of sample error should be bigger than in q10 while we have the whole training set used, <u>we got an expected value slightly smaller because we retrain the whole training set after the optimal $\lambda$ is chosed</u>. We did not use g⁻ to predict the test set, but g instead.

Apart from this reason, I think this result may also due to the precision, if more digits are taken into consideration, maybe we can observe the difference.

Another cause may be the limited number of possible $\lambda$, this constraints the variability of our model, even though there are infinitely many possible models, most of them are alike, while only slightly different by the $\lambda$ value chosen.

## Code:

```
1   Eouts = []
2
3   for experiment in tqdm(range(1126)):
4
5       np.random.seed(experiment)
6
7       total_size = len(X_train)
8       indices = np.random.permutation(total_size)
9
10      subtrain_indices = indices[:8000]
11      validation_indices = indices[8000:]
12
13      X_subtrain = [X_train[i] for i in subtrain_indices]
14      X_validation = [X_train[i] for i in validation_indices]
15      y_subtrain = [y_train[i] for i in subtrain_indices]
16      y_validation = [y_train[i] for i in validation_indices]
17
18      subtrain_prob = problem(y_subtrain, X_subtrain)
19
20      min_validation_err = np.inf
21      opt_log10_lambda = 0
22      for log10_lambda in (-2, -1, 0, 1, 2, 3):
23          subtrain_pred_res = []
24          c = 1 / (10 ** log10_lambda)
25          param = parameter('-s 6 -c ' + str(c))
26          model_by_subtrain = train(subtrain_prob, param)
27
28          validation_label, _, _ = predict(y_validation, X_validation, model_by_subtrain)
29          validation_err = ZeroOneError(validation_label, y_validation)
30          if validation_err == min_validation_err:
31              opt_log10_lambda = max(opt_log10_lambda, log10_lambda)     # break tie by choosing the larger lambda
32          elif validation_err < min_validation_err:
33              min_validation_err = validation_err
34              opt_log10_lambda = log10_lambda
35
36      # rerun the model with the best lambda on the whole training set
37      whole_train_prob = problem(y_train, X_train)
38      param_with_opt_lambda = parameter('-s 6 -c ' + str(1 / (10 ** opt_log10_lambda)))
39      whole_train_model = train(whole_train_prob, param_with_opt_lambda)
40
41      test_label, _, _ = predict(y_test, X_test, whole_train_model)
42      Eout = ZeroOneError(test_label, y_test)
43
44      Eouts.append(Eout)
✓   61m 42.5s
```

2

---

[2] Other part of codes that are similar to problem 10 (like reading in the data, split the data to get proper form of X and y, the ZeroOneError function…) are omitted in this report, can check for problem 10 for details.