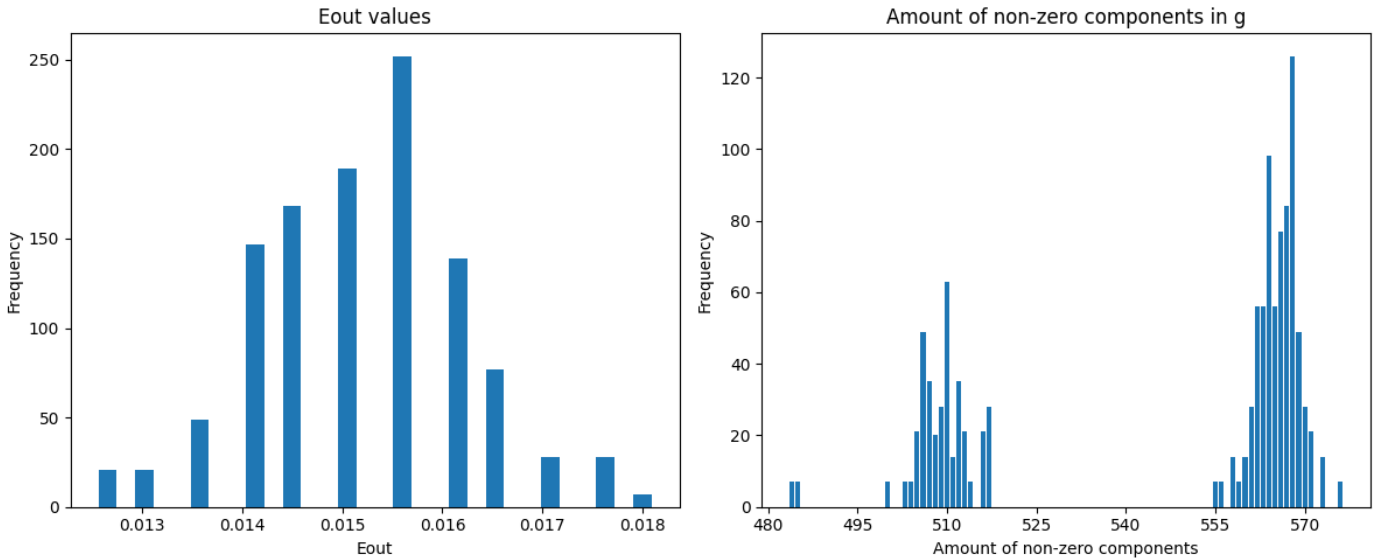# HTML homework 5: q10 report

The resulting plot is as follows:



Some additional information that I printed in order to do the comparison between this question and question 11:

```
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
```

## Code[1]:

### Preprocessing the dataset:

```
1   from liblinear.liblinearutil import *
2   from itertools import combinations_with_replacement
3   import numpy as np
4   import scipy.sparse
5   import matplotlib.pyplot as plt
6   from tqdm import tqdm
7   from joblib import Parallel, delayed
8
```
[28]  ✓  0.4s

## Dataset

- `mnist.scale` for training
- `mnist.scale.t` for testing

```
1   train_data = 'mnist.scale'
2   test_data = 'mnist.scale.t'
```
[2]

*Empty markdown cell, double-click or press enter to edit.*

## Read and save data

The training data is saved in the `X_train` list, with each element being a dictionary of features. example format: {153: 0.0117647, 154: 0.0705882, 155: 0.0705882, ...}

The corresponding label list is saved in `y_train`.

Similarily, the testing data is saved in `X_test` and the corresponding label list is saved in `y_test`.

```
1   def read_linear_format(file_path):
2       X, y = [], []
3       with open(file_path, 'r') as f:
4           for line in f:
5               parts = line.strip().split()
6               if parts[0] == '2' or parts[0] == '6':
7                   y.append(int(parts[0]))
8                   features = {}
9                   for item in parts[1:]:
10                      index, value = item.split(":")
11                      features[int(index)] = float(value)
12                  X.append(features)
13      return X, y
14
15  X_train, y_train = read_linear_format(train_data)
16  X_test, y_test = read_linear_format(test_data)
17
```

### Error function:

## 0/1 error

```
1   def ZeroOneError(predictions, y):
2       return sum(p_i != y_i for p_i, y_i in zip(predictions, y)) / len(predictions)
```
[6]

---

[1] Most of the screenshots here are from  my .ipynb file, since I wrote the explanations here, however, I ran the code in a .py file (which has equivalent contents), and part of the screenshot is from this file.

Explanation of deriving the relationship between C, $\lambda$:

# Regularized logistic regression
## Original Formulation in README.md

As the explanation in the liblinear README file, for L1-regularized logistic regression `(-s 6)`, we solve

$$\min_{\vec{w}} \sum |w_j| + C \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i))$$

## Components of this equation

In this equation, the L1 regularizer is defined as:

$$\Omega(\vec{w}) = \sum |w_j| = ||\vec{w}||_1$$

> info in Lec14 slide 18

- Note: using L1 regularization encourages sparsity, which means fewer parameters.

And the in-sample error for logistic regression is:

$$E_{in}(\vec{w}) = \frac{1}{N} \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i))$$

> in textbook p.91

## Derivation of $C$

Since we knew that:

$$E_{aug}(\vec{w}) = E_{in}(\vec{w}) + \frac{\lambda}{N} \Omega(\vec{w})$$

We can derive that:

$$E_{aug}(\vec{w}) = E_{in}(\vec{w}) + \frac{\lambda}{N} \sum |w_j|$$

For minimization:

$$\arg\min_{\vec{w}} E_{aug}(\vec{w}) = \arg\min_{\vec{w}}(E_{in}(\vec{w}) + \frac{\lambda}{N} \sum |w_j|)$$
$$= \arg\min_{\vec{w}} \frac{N}{\lambda}(\frac{1}{N} \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i)) + \frac{\lambda}{N} \sum |w_j|)$$
$$= \arg\min_{\vec{w}} \frac{1}{\lambda} \sum \log(1 + \exp(-y_i \vec{w}^T \vec{x}_i)) + \sum |w_j|$$

Compare to the original equation marked with ($\star$), we can see that:

$$C = \frac{1}{\lambda}$$

Some additional information:

# Select the best $\lambda^*$

Selecting optimal $\lambda^*$ using the following equation:

$$\lambda^* = \arg\min_{\log_{10}(\lambda) \in \{-2,-1,0,1,2,3\}} E_{in}(\vec{w}_\lambda)$$

## get amount of non-zeros

To get the amount of non-zero components in the model, we can use the function `.get_decfun()`:

```
[W, b] = model_.get_decfun()
```

Main code:

```python
1   def run_single_experiment(experiment):
2       np.random.seed(experiment)
3       min_Ein = np.inf
4       opt_log10_lambda = 0
5
6       for log10_lambda in (-2, -1, 0, 1, 2, 3):
7           c = 1 / (10 ** log10_lambda)
8           prob = problem(y_train, X_train)
9           param = parameter('-s 6 -c ' + str(c))
10          model = train(prob, param)
11
12          train_label, _, _ = predict(y_train, X_train, model)
13          Ein = ZeroOneError(train_label, y_train)
14          if Ein == min_Ein:
15              opt_log10_lambda = max(opt_log10_lambda, log10_lambda)
16              if opt_log10_lambda == log10_lambda:
17                  opt_model = model
18          elif Ein < min_Ein:
19              min_Ein = Ein
20              opt_log10_lambda = log10_lambda
21              opt_model = model
22
23      test_label, _, _ = predict(y_test, X_test, opt_model)
24      Eout = ZeroOneError(test_label, y_test)
25
26      W = np.array(opt_model.get_decfun()[0])
27      non_zero_count = np.count_nonzero(W)
28
29      return Eout, non_zero_count
30
31  experiment_amount = 10
32  results = Parallel(n_jobs=-1)(                                    # use all cores by setting n_jobs to -1
33      delayed(run_single_experiment)(i) for i in tqdm(range(experiment_amount))
34  )
35
36  Eouts, non_zero_count_list = zip(*results)
37  Eouts = list(Eouts)
38  non_zero_count_list = list(non_zero_count_list)
```

Plotting:

```python
# aim: plotting
# subaim: First subplot (Eout values)
ax1.hist(Eouts, bins=30)
ax1.set_title('Eout values')
ax1.set_xlabel('Eout')
ax1.set_ylabel('Frequency')

# subaim: Second subplot (non-zero components)
min_val = int(min(non_zero_count_list))
max_val = int(max(non_zero_count_list))
integer_bins = np.arange(min_val, max_val + 2) - 0.5  # +2 to include max_val, -0.5 for bin edges

ax2.hist(non_zero_count_list,
         bins=integer_bins,
         align='mid',
         rwidth=0.8)

ax2.xaxis.set_major_locator(plt.MaxNLocator(integer=True))
ax2.set_title('Amount of non-zero components in g')
ax2.set_xlabel('Amount of non-zero components')
ax2.set_ylabel('Frequency')

plt.tight_layout()
plt.show()
```