# COMP5212 Machine Learning Project 1
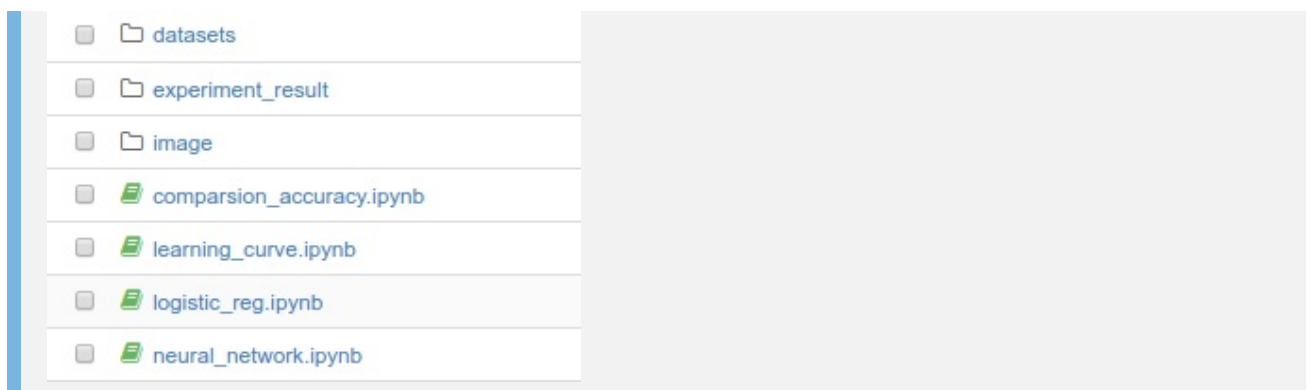
## Requirement

1. Anaconda
2. jupyter notebook

> http://jupyter.readthedocs.io/en/latest/install.html

## How to use

1. File structure



2. There are three main modules for this projects: logistic_reg.ipynb, neural network.ipynb, svm.ipynb, while comparsion_accuracy.ipynb is for drawing comparsion curves between the four models. All operations like building classificer, training and testing are wrote in the three main files with same code structure.
   - neural_network.ipynb
     1. Module 1: Setting the parameters

```
In [43]:   1  # Module 1: setting the parameters
           2  m_i = 1000
           3  c = 0.001
           4  p = "l2"
           5  eta = 0.1
           6  learning_rate = 'adaptive'
           7  repeat_times = 5
           8  file_count = 5
           9  class_names = ['0', '1']
          10  max_H = 11
```

2. Module 2: When you get a optimal H, this module could help you to **train and test a classifier**.

```
In [44]:   1  # Module 2: when you get a optimal H, this module could help you to train and test a clas
           2  import time
           3
           4  import numpy as np
           5  import matplotlib.pyplot as plot
           6  from sklearn.metrics import confusion_matrix
           7
           8  from sklearn.neural_network import MLPClassifier
           9  from sklearn.preprocessing import MinMaxScaler
          10  from sklearn.metrics import accuracy_score
          11  from sklearn.metrics import log_loss
          12
          13  # initialize some system value
          14  input_data_filename = ['breast-cancer', 'diabetes',
          15                         'digit', 'iris', 'wine']
          16  result = np.array([['Dataset', '$m\_i$', '$c$', '$\eta$', '$p$', 'H',
          17                      '$a_{train}(\%)$', '$a_{test}(\%)$', '$l_{train}$',
          18                      '$l_{test}$', '$time(ms)$']])
          19
          20  clf = MLPClassifier(solver='sgd', activation='logistic', alpha=1e-5,
          21                      random_state=0, max_iter=m_i, tol=c,
          22                      learning_rate_init=eta, learning_rate=learning_rate, verbose=False)
          23
          24  H_star = [3,9,10,8,4]
          25
          26  # Create a plot
          27  fig, axes=plt.subplots(2, 5, figsize=(15, 10))
          28
          29  for i in range(0, file_count):
          30      clf.set_params(hidden_layer_sizes=H_star[i])
          31      print('Reading data from: ', input_data_filename[i])
          32      data = np.load('datasets/' + input_data_filename[i] + '.npz')
          33
          34      accu_train, accu_test, loss_train, loss_test, t_last, cnf_matrix_train, cnf_matrix_te
          35
          36      ''' report parameters '''
          37      print('m_i', m_i, 'c', c, 'learning rate', eta, 'p', p, 'H', H_star[i],
          38            'a_train', accu_train*100, 'a_test', accu_test*100,
          39            'loss_train', loss_train, 'loss_test', loss_test, 'time=train+test', t_last, 'm
          40
          41      plt.subplot(2,5,i+1)
          42      plot_confusion_matrix(cnf_matrix_train, classes=class_names, normalize=False, title=i
          43      plt.subplot(2,5,5+i+1)
          44      plot_confusion_matrix(cnf_matrix_test, classes=class_names, normalize=False, title=ir
          45
          46      #     print('auc', compute_auc(train_Y, posterior_train_Y[:,1]))
          47
          48      newrow = np.array([[input_data_filename[i], m_i, c, eta, p, H_star[i], accu_train*100
          49                          accu_test*100, loss_train, loss_test, t_last]])
          50      result = np.append(result, newrow, axis=0)
          51
          52  output_csv(result)
          53
          54  plt.savefig('neural_confusion_matrix.eps', dpi=300)
          55  plt.show()
```

3. Module 3: **Cross-validation for choosing a optimal H** by implementing a cross-validation

```
In [45]:   1  # Module 3: Cross-validation for choosing a optimal H by implementing a cross-validation
           2  import time
           3
           4  import numpy as np
           5  import matplotlib.pyplot as plot
           6
           7  from sklearn.neural_network import MLPClassifier
           8  from sklearn.preprocessing import MinMaxScaler
           9  from sklearn.metrics import accuracy_score
          10  from sklearn.metrics import log_loss
          11
          12  # initialize some system value
          13  input_data_filename = ['breast-cancer', 'diabetes',
          14                         'digit', 'iris', 'wine']
          15  result = np.array([['Dataset', '$m\_i$', '$c$', '$\eta$', '$p$', 'H',
          16                      '$a_{test}(\%)$', '$l_{train}$', '$l_{test}$', '$time(ms)$']])
          17
          18  # Create a mlp classifier
          19  # learning_rate='invscaling', 'adaptive'
          20  clf = MLPClassifier(solver='sgd', activation='logistic', alpha=1e-5,
          21                      random_state=0, max_iter=m_i, tol=c, learning_rate_init=eta,
          22                      learning_rate='adaptive', verbose=False)
          23  print('Parameters: ', clf.get_params(True))
          24
          25  # Create a plot
          26  fig, axes = plt.subplots(2, 3, figsize=(15, 10))
          27
          28  for i in range(0, file_count):
          29      print('Reading data from: ', input_data_filename[i])
          30      data = np.load('datasets/' + input_data_filename[i] + '.npz')
          31
          32      # cross validataion(80% train, 20% test) for H_star
          33      H_star = cross_validation_(data, max_H, i)
          34      print('H_star:', H_star)
          35
          36  plt.savefig('H_selection.eps', dpi=600)
          37  plt.show()
```

4. Other modules: Defining some functions

```
In [29]:   1  def train_test(clf, data):
           2      accu_train = []
           3      accu_test = []
           4      loss_train = []
           5      loss_test = []
           6      t_last = []
           7
           8      train_X, train_Y, test_X, test_Y =
           9          data['train_X'], data['train_Y'], data['test_X'], data['test_Y']
          10
          11      for j in range(0, repeat_times):
          12          t_begin = time.time()
          13
          14          clf.fit(train_X, train_Y)
          15
          16          predict_train_Y = clf.predict(train_X)
          17          predict_test_Y = clf.predict(test_X)
          18          t_last.append(time.time() - t_begin)
          19
          20          posterior_train_Y = clf.predict_proba(train_X)
          21          posterior_test_Y = clf.predict_proba(test_X)
          22
          23          accu_train.append(accuracy_score(predict_train_Y, train_Y))
          24          accu_test.append(accuracy_score(predict_test_Y, test_Y))
          25
          26          loss_train.append(log_loss(train_Y, posterior_train_Y, normalize=True))
          27          loss_test.append(log_loss(test_Y, posterior_test_Y, normalize=True))
          28
          29      cnf_matrix_train = confusion_matrix(predict_train_Y, data['train_Y'])
          30      cnf_matrix_test = confusion_matrix(predict_test_Y, data['test_Y'])
          31
          32      accu_train = round(np.mean(accu_train), 4)
          33      accu_test = round(np.mean(accu_test), 4)
          34      loss_train = round(np.mean(loss_train), 4)
          35      loss_test = round(np.mean(loss_test), 4)
          36      t_last = round(np.mean(t_last)*1000, 4)
          37
          38      ''' report train and test error '''
          39      print('Average training data accuracy:', accu_train)
          40      print('Average testing data accuracy:', accu_test)
          41
          42      ''' report train and test log loss'''
          43      print('Average training data log loss:', loss_train)
          44      print('Average testing data log loss:', loss_test)
          45      print('Average Time ms', t_last)
          46
          47      return accu_train, accu_test, loss_train, loss_test, t_last,
          48          cnf_matrix_train, cnf_matrix_test
```

- svm.ipynb
  1. Module 1: Setting the parameters
  2. Module 2: When you get a optimal gamma, this module could help you to **train and test a classifier**.
  3. Module 3: **Cross-validation for choosing a**

> **optimal gamma** by implementing a cross-validation

4. Other modules: Defining some functions.
- logistic_reg.ipynb
  1. Module 1: Setting the parameters
  2. Module 2: This module could help you to **train and test a classifier**
  3. Module 3: measuring model's performance overtimes

```python
iteraitons = 30
batch_num = 20
for X, Y in zip(np.array_split(train_X, batch_num),
                np.array_split(train_Y, batch_num)):

    clf.partial_fit(X, Y, classes=classes)
    # loss
    loss.append(log_loss(train_Y, clf.predict(train_X),
                         normalize=True))
    # accuracy on train & test data
    accu_train.append(clf.score(train_X, train_Y))
    accu_test.append(clf.score(test_X, test_Y))
    # time required for learning and testing
    t_last.append(time.time() - t_begin)
print(clf.score(train_X, train_Y)*100)
print(log_loss(train_Y, clf.predict(train_X), normalize=True))
print(clf.score(test_X, test_Y)*100)
print(log_loss(test_Y, clf.predict(test_X), normalize=True))
print(np.sum(t_last)*1000, 'ms')
```
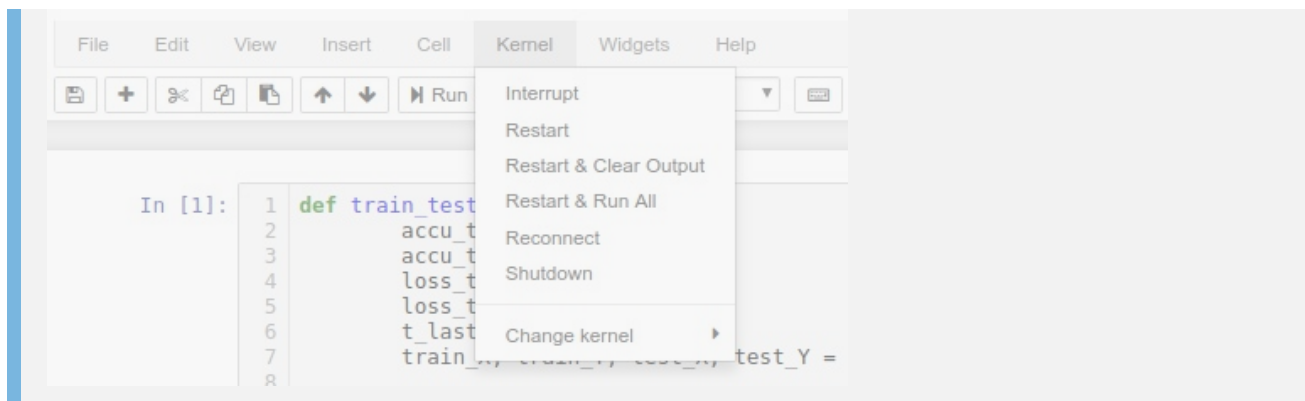
4. Other modules: Defining some functions
- Notes
  1. Different from neural_network.ipynb and svm.ipynb, logistic_reg.ipynb does not have the **cross-validation module**, but have a **performance measurement module**.

## Resulting

1. If you need to verify my result, **please choose any ipynb file, and select Kernel -> Restart & Run All, all the modules in this file will run automatically in order.** You can check my output of the module 2 and module 3.

2. When you run the code successfully, you could see the result directly like this:

```
Parameters:  {'C': 1.0, 'cache_size': 200, 'class_weight': None, 'coef0': 0.0, 'decision_fun
ction_shape': 'ovr', 'degree': 3, 'gamma': 'auto', 'kernel': 'rbf', 'max_iter': 2000, 'proba
bility': False, 'random_state': 0, 'shrinking': True, 'tol': 0.001, 'verbose': False}
Reading data from:  breast-cancer
gamma: [1, 0.1, 0.01, 0.001]
g scores: [96.90909091 96.90909091 96.18181818 94.90909091]
loss_train: [0.04146357990578099, 0.08044230522059559, 0.08779153988272645, 0.09302266582939
434]
loss_test: [0.09149580916363458, 0.08151004722497966, 0.09416992130071798, 0.090628809044320
43]
gamma_star: 1
Reading data from:  diabetes
gamma: [1, 0.1, 0.01, 0.001]
g scores: [75.44715447 75.28455285 66.99186992 68.61788618]
loss_train: [0.41023770039454116, 0.45587392011590994, 0.4882479154117044, 0.493551122092130
26]
loss_test: [0.4837010736506908, 0.5058437423441038, 0.5155455409104033, 0.4992173633034017]
gamma_star: 1
Reading data from:  digit
gamma: [1, 0.1, 0.01, 0.001]
g scores: [52.25 52.75 94.5  99.75]
loss_train: [4.074427049150545, 4.064574892274672, 5.171877225667428e-08, 0.0041303289507558
05]
loss_test: [0.6929971672851531, 0.6930771600879907, 0.07868577993953338, 0.01905411812063057
7]
gamma_star: 4
Reading data from:  iris
gamma: [1, 0.1, 0.01, 0.001]
g scores: [100. 100. 100.  70.]
loss_train: [0.016753443802707425, 0.019292949029276817, 0.024712069496450908, 0.02897778998
85727]
loss_test: [0.01830075442806535, 0.018003746683217844, 0.030492577630699808, 0.0446722603927
2951]
gamma_star: 1
Reading data from:  wine
gamma: [1, 0.1, 0.01, 0.001]
g scores: [63.44827586 64.82758621 76.55172414 81.37931034]
loss_train: [4.493428126865702, 0.031194129093737026, 0.08264005223288591, 0.356101913402283
25]
loss_test: [0.6653871802788931, 0.5920967915997062, 0.5086273491819073, 0.46437840030289823]
gamma_star: 4

<matplotlib.figure.Figure at 0x7ff12832ff60>
```

# Contact

1. If you do not about jupyter, or meet some problems about running it, please feel free to contact me: jjiao@ust.hk.