

Project 3

Classification and inference with machine learning

This notebook is arranged in cells. Texts are usually written in the markdown cells, and here you can use html tags (make it bold, italic, colored, etc). You can double click on this cell to see the formatting.

The ellipsis (...) are provided where you are expected to write your solution but feel free to change the template (not over much) in case this style is not to your taste.

Hit "Shift-Enter" on a code cell to evaluate it. Double click a Markdown cell to edit.

Link Okpy

```
In [ ]: 1 from client.api.notebook import Notebook
        2 ok = Notebook('Project3_U.ok')
        3 _ = ok.auth(inline = True)
```

Imports

```
In [ ]: 1 import numpy as np
        2 from scipy.integrate import quad
        3 #For plotting
        4 import matplotlib.pyplot as plt
        5 %matplotlib inline
        6 import warnings
        7 warnings.filterwarnings('ignore')
```

Problem 1 - Using Keras - MNIST

The goal of this notebook is to introduce deep neural networks (DNNs) and convolutional neural networks (CNNs) using the high-level Keras package and to become familiar with how to choose its architecture, cost function, and optimizer in Keras. We will also learn how to train neural networks.

We will once again work with the MNIST dataset of hand written digits introduced in HW8. The goal is to find a statistical model which recognizes and distinguishes between the ten handwritten digits (0-9).

The MNIST dataset comprises handwritten digits, each of which comes in a square image, divided into a 28×28 pixel grid. Every pixel can take on 256 nuances of the gray color, interpolating between white and black, and hence each data point assumes any value in the set $\{0, 1, \dots, 255\}$. Since there are 10 categories in the problem, corresponding to the ten digits, this problem represents a generic classification task.

In this Notebook, we show how to use the Keras python package to tackle the MNIST problem with the help of deep neural networks.

Creating DNNs with Keras

Constructing a Deep Neural Network to solve ML problems is a multiple-stage process. Quite generally, one can identify the key steps as follows:

- step 1: Load and process the data
- step 2: Define the model and its architecture
- step 3: Choose the optimizer and the cost function
- step 4: Train the model
- step 5: Evaluate the model performance on the unseen test data
- step 6: Modify the hyperparameters to optimize performance for the specific data set

We would like to emphasize that, while it is always possible to view steps 1-5 as independent of the particular task we are trying to solve, it is only when they are put together in step 6 that the real gain of using Deep Learning is revealed, compared to less sophisticated methods such as the regression models. With this remark in mind, we shall focus predominantly on steps 1-5 below. We show how one can use grid search methods to find optimal hyperparameters in step 6.

Step 1: Load and Process the Data

Keras knows to download automatically the MNIST data from the web. All we need to do is import the `mnist` module and use the `load_data()` class, and it will create the training and test data sets or us.

The MNIST set has pre-defined test and training sets, in order to facilitate the comparison of the performance of different models on the data.

Once we have loaded the data, we need to format it in the correct shape $((N_{\text{samples}}, N_{\text{features}}))$.

The size of each sample, i.e. the number of bare features used is `N_features` (whis is 784 because we have a 28×28 pixel grid), while the number of potential classification categories is "num_classes" (which is 10, number of digits).

Each pixel contains a greyscale value quantified by an integer between 0 and 255. To standardize the dataset, we normalize the input data in the interval $[0, 1]$.

```
In [ ]: 1 from __future__ import print_function
2 import keras,sklearn
3 # suppress tensorflow compilation warnings
4 import os
5 import tensorflow as tf
6 os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'
7 seed=0
8 np.random.seed(seed) # fix random seed
9 tf.set_random_seed(seed)
10
11 from keras.datasets import mnist
12
13 # input image dimensions
14 num_classes = 10 # 10 digits
15
16 img_rows, img_cols = 28, 28 # number of pixels
17
18 # the data, shuffled and split between train and test sets
19 (X_train, Y_train), (X_test, Y_test) = mnist.load_data()
20 X_train = X_train[:40000]
21 Y_train = Y_train[:40000]
22
23 # reshape data, depending on Keras backend
24 X_train = X_train.reshape(X_train.shape[0], img_rows*img_cols)
25 X_test = X_test.reshape(X_test.shape[0], img_rows*img_cols)
26
27 # cast floats to single precision
28 X_train = X_train.astype('float32')
29 X_test = X_test.astype('float32')
30
31 # rescale data in interval [0,1]
32 X_train /= 255
33 X_test /= 255
34
```

1. Make a plot of one MNIST digit (2D plot using X data - make sure to reshape it into a 28×28 matrix) and label it (which digit does it correspond to?).

```
In [ ]: 1 ...
```

Last, we cast the label vectors y to binary class matrices (a.k.a. one-hot format).

```
In [ ]: 1 # convert class vectors to binary class matrices
2
3 print("before conversion - ")
4 print("y vector : ", Y_train[0:10])
5
6 Y_train = keras.utils.to_categorical(Y_train, num_classes)
7 Y_test = keras.utils.to_categorical(Y_test, num_classes)
8
9 print("after conversion - ")
10 print("y vector : ", Y_train[0:10])
```

Here in this template, we use 40000 training samples and 10000 test samples. Remember that we preprocessed data into the shape ($N_{\text{samples}}, N_{\text{features}}$).

```
In [ ]: 1 print('X_train shape:', X_train.shape)
2 print('Y_train shape:', Y_train.shape)
3 print()
4 print(X_train.shape[0], 'train samples')
5 print(X_test.shape[0], 'test samples')
```

Step 2: Define the Neural Net and its Architecture

We can now move on to construct our deep neural net. We shall use Keras's `Sequential()` class to instantiate a model, and will add different deep layers one by one.

Let us create an instance of Keras' `Sequential()` class, called `model`. As the name suggests, this class allows us to build DNNs layer by layer. (<https://keras.io/getting-started/sequential-model-guide/> (<https://keras.io/getting-started/sequential-model-guide/>))

```
In [ ]: 1 from keras.models import Sequential
2 from keras.layers import Dense, Dropout, Flatten
3 from keras.layers import Conv2D, MaxPooling2D
4
5 # instantiate model
6 model = Sequential()
```

We use the `add()` method to attach layers to our model. For the purposes of our introductory example, it suffices to focus on Dense layers for simplicity. (<https://keras.io/layers/core/> (<https://keras.io/layers/core/>)) Every `Dense()` layer accepts as its first required argument an integer which specifies the number of neurons. The type of activation function for the layer is defined using the `activation` optional argument, the input of which is the name of the activation function in string format. Examples include `relu`, `tanh`, `elu`, `sigmoid`, `softmax`.

In order for our DNN to work properly, we have to make sure that the numbers of input and output neurons for each layer match. Therefore, we specify the shape of the input in the first layer of the model explicitly using the optional argument `input_shape=(N_features,)`. The sequential construction of the model then allows Keras to infer the correct input/output dimensions of all hidden layers automatically. Hence, we only need to specify the size of the softmax output layer to match the number of categories.

First, add a Dense layer with 400 output neurons and `relu` activation function.

```
In [ ]: 1 model.add(Dense(400, input_shape=(img_rows*img_cols,), activation='relu'))
```

Add another layer with 100 output neurons. Then, we will apply "dropout," a regularization scheme that has been widely adopted in the neural networks literature: during the training procedure neurons are randomly "dropped out" of the neural network with some probability p giving rise to a thinned network. It prevents overfitting by reducing spurious correlations between neurons within the network by introducing a randomization procedure.

```
In [ ]: 1 model.add(Dense(100, activation='relu'))
2 # apply dropout with rate 0.5
3 model.add(Dropout(0.5))
```

Lastly, we need to add a soft-max layer since we have a multi-class output.

```
In [ ]: 1 model.add(Dense(num_classes, activation='softmax'))
```

Step 3: Choose the Optimizer and the Cost Function

Next, we choose the loss function according to which to train the DNN. For classification problems, this is the cross entropy, and since the output data was cast in categorical form, we choose the `categorical_crossentropy` defined in Keras' `losses` module. Depending on the problem of interest one can pick any other suitable loss function. To optimize the weights of the net, we choose SGD. This algorithm is already available to use under Keras' `optimizers` module (<https://keras.io/optimizers/> (<https://keras.io/optimizers/>)), but we could use `Adam()` or any other built-in one as

well. The parameters for the optimizer, such as `lr` (learning rate) or `momentum` are passed using the corresponding optional arguments of the `SGD()` function.

While the loss function and the optimizer are essential for the training procedure, to test the performance of the model one may want to look at a particular metric of performance. For instance, in categorical tasks one typically looks at their accuracy, which is defined as the percentage of correctly classified data points.

To complete the definition of our model, we use the `compile()` method, with optional arguments for the optimizer, loss, and the validation metric as follows:

```
In [ ]: 1 # compile the model
        2 model.compile(loss=keras.losses.categorical_crossentropy, optimizer='SGD', metrics=['accuracy'])
        3
```

Step 4: Train the model

We train our DNN in minibatches. Shuffling the training data during training improves stability of the model. Thus, we train over a number of training epochs.

(The number of epochs is the number of complete passes through the training dataset, and the batch size is a number of samples propagated through the network before the model is updated.)

Training the DNN is a one-liner using the `fit()` method of the `Sequential` class. The first two required arguments are the training input and output data. As optional arguments, we specify the `mini_batch_size`, the number of training epochs, and the test or validation data. To monitor the training procedure for every epoch, we set `verbose=True`.

Let us set `batch_size = 64` and `epochs = 10`.

```
In [ ]: 1 # training parameters
        2 batch_size = 64
        3 epochs = 10
        4
        5 # train DNN and store training info in history
        6 history=model.fit(X_train, Y_train, batch_size=batch_size, epochs=epochs,
        7                 verbose=1, validation_data=(X_test, Y_test))
```

Step 5: Evaluate the Model Performance on the Unseen Test Data

Next, we evaluate the model and read of the loss on the test data, and its accuracy using the `evaluate()` method.

```
In [ ]: 1 # evaluate model
        2 score = model.evaluate(X_test, Y_test, verbose=1)
        3
        4 # print performance
        5 print('Test loss:', score[0])
        6 print('Test accuracy:', score[1])
        7
        8 # look into training history
        9
        10 # summarize history for accuracy
        11 plt.plot(history.history['acc'])
        12 plt.plot(history.history['val_acc'])
        13 plt.ylabel('model accuracy')
        14 plt.xlabel('epoch')
        15 plt.legend(['train', 'test'], loc='best')
        16 plt.show()
        17
        18 # summarize history for loss
        19 plt.plot(history.history['loss'])
        20 plt.plot(history.history['val_loss'])
        21 plt.ylabel('model loss')
        22 plt.xlabel('epoch')
        23 plt.legend(['train', 'test'], loc='best')
        24 plt.show()
```

Step 6: Modify the Hyperparameters to Optimize Performance of the Model

Last, we show how to use the grid search option of scikit-learn (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html) to optimize the hyperparameters of our model.

First, define a function for crating a DNN:

```
In [ ]: 1 def create_DNN(optimizer=keras.optimizers.Adam()):
        2     model = Sequential()
        3     model.add(Dense(400, input_shape=(img_rows*img_cols,), activation='relu'))
        4     model.add(Dense(100, activation='relu'))
        5     model.add(Dropout(0.5))
        6     model.add(Dense(num_classes, activation='softmax'))
        7     model.compile(loss=keras.losses.categorical_crossentropy,
        8               optimizer=optimizer,
        9               metrics=['accuracy'])
        10     return model
```

With `epochs = 1` and `batch_size = 64`, do grid search over the following optimization schemes: ['SGD', 'RMSprop', 'Adagrad', 'Adadelata', 'Adam', 'Adamax', 'Nadam'].

```
In [ ]: 1 from sklearn.model_selection import GridSearchCV
        2 from keras.wrappers.scikit_learn import KerasClassifier
        3
        4 batch_size = 64
        5 epochs = 1
        6 model_gridsearch = KerasClassifier(build_fn=create_DNN,
        7                                 epochs=epochs, batch_size=batch_size, verbose=1)
        8
        9 # list of allowed optional arguments for the optimizer, see `compile_model()`
        10 optimizer = ['SGD', 'RMSprop', 'Adagrad', 'Adadelata', 'Adam', 'Adamax', 'Nadam']
        11 # define parameter dictionary
        12 param_grid = dict(optimizer=optimizer)
        13
        14 # call scikit grid search module
        15 grid = GridSearchCV(estimator=model_gridsearch, param_grid=param_grid, n_jobs=1, cv=4)
        16 grid_result = grid.fit(X_train, Y_train)
```

Show the mean test score of all optimization schemes and determine which scheme gives the best accuracy.

```
In [ ]: 1 # summarize results
        2 print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
        3 means = grid_result.cv_results_['mean_test_score']
        4 stds = grid_result.cv_results_['std_test_score']
        5 params = grid_result.cv_results_['params']
        6 for mean, stdev, param in zip(means, stds, params):
        7     print("%f (%f) with: %r" % (mean, stdev, param))
```

2. Create a DNN with one Dense layer having 200 output neurons. Do the grid search over any 5 different activation functions from <https://keras.io/activations/> (<https://keras.io/activations/>). Let epochs = 1, batches = 64, p_dropout=0.5, and optimizer=keras.optimizers.Adam(). Make sure to print the mean test score of each case and determine which activation functions gives the best accuracy.

Doing the grid search requires quite a bit of memory. Please restart the kernel ("Kernel"- "Restart") and re-load the data before doing a new grid search.

```
In [ ]: 1 ...
```

3. Now, do the grid search over different combination of batch sizes (10, 30, 50, 100) and number of epochs (1, 2, 5). Make sure to print the mean test score of each case and determine which activation functions gives the best accuracy. Here, you have a freedom to create your own DNN (assume an arbitrary number of Dense layers, optimization scheme, etc).

Doing the grid search requires quite a bit of memory. Please restart the kernel ("Kernel"- "Restart") and re-load the data before doing a new grid search.

Hint: To do the grid search over both batch_size and epochs, you can do:

```
param_grid = dict(batch_size=batch_size, epochs=epochs)
```

```
In [ ]: 1 ...
```

4. Do the grid search over the number of neurons in the Dense layer and make a plot of mean test score as a function of num_neurons. Again, you have a freedom to create your own DNN.

Doing the grid search requires quite a bit of memory. Please restart the kernel ("Kernel"- "Restart") and re-load the data before doing a new grid search.

```
In [ ]: 1 ...
```

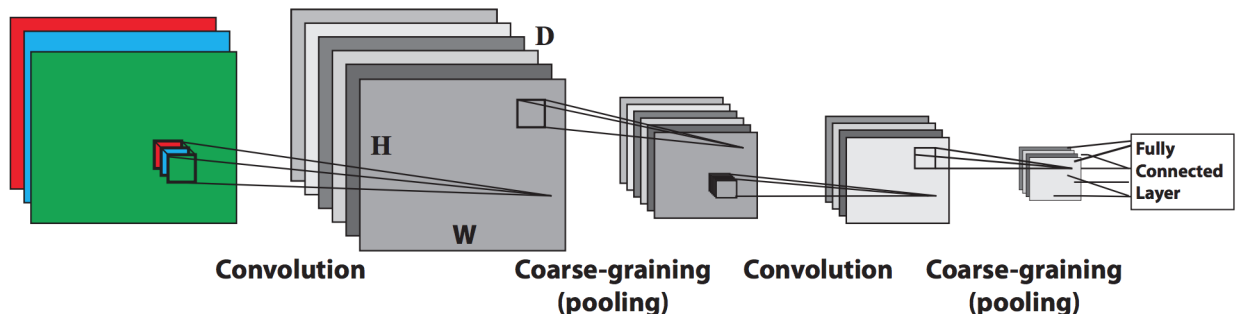
Creating CNNs with Keras

Please restart the kernel ("Kernel"- "Restart") and re-load the data.

We have so far considered each MNIST data sample as a (28×28) -long 1d vector. This approach neglects any spatial structure in the image. On the other hand, we do know that in every one of the hand-written digits there are local spatial correlations between the pixels, which we would like to take advantage of to improve the accuracy of our classification model. To this end, we first need to reshape the training and test input data as follows

```
In [ ]: 1 # reshape data, depending on Keras backend
2 if keras.backend.image_data_format() == 'channels_first':
3     X_train = X_train.reshape(X_train.shape[0], 1, img_rows, img_cols)
4     X_test = X_test.reshape(X_test.shape[0], 1, img_rows, img_cols)
5     input_shape = (1, img_rows, img_cols)
6 else:
7     X_train = X_train.reshape(X_train.shape[0], img_rows, img_cols, 1)
8     X_test = X_test.reshape(X_test.shape[0], img_rows, img_cols, 1)
9     input_shape = (img_rows, img_cols, 1)
10
11 print('X_train shape:', X_train.shape)
12 print('Y_train shape:', Y_train.shape)
13 print()
14 print(X_train.shape[0], 'train samples')
15 print(X_test.shape[0], 'test samples')
```

One can ask the question of whether a neural net can learn to recognize such local patterns. This can be achieved by using convolutional layers. Luckily, all we need to do is change the architecture of our DNN.



After we instantiate the model, add the first convolutional layer with 10 filters, which is the dimensionality of output space. (<https://keras.io/layers/convolutional/>) Here, we will be concerned with local spatial filters that take as inputs a small spatial patch of the previous layer at all depths. We consider a three-dimensional kernel of size $5 \times 5 \times 1$. Check out this visualization of the convolution procedure for a square input of unit depth: https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md (https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md) The convolution consists of running this filter over all locations in the spatial plane. After computing the filter, the output is passed through a non-linearity, a ReLU.

```
In [ ]: 1 from keras.models import Sequential
2 from keras.layers import Dense, Dropout, Flatten
3 from keras.layers import Conv2D, MaxPooling2D
4
5 model = Sequential()
6 model.add(Conv2D(10, kernel_size=(5, 5),
7                 activation='relu',
8                 input_shape=input_shape))
```

Subsequently, add a 2D pooling layer. (<https://keras.io/layers/pooling/>) This pooling layer coarse-grain spatial information by performing a subsampling at each depth. Here, we use the the max pool operation. In a max pool, the spatial dimensions are coarse-grained by replacing a small region (say 2×2 neurons) by a single neuron whose output is the maximum value of the output in the region.

```
In [ ]: 1 model.add(MaxPooling2D(pool_size=(2, 2)))
```

Add another convolutional layers with 20 filters and apply dropout. Then, add another pooling layer and flatten the data. You can do DNNs afterwards and compile the model.

```
In [ ]: 1 # add second convolutional layer with 20 filters
2 model.add(Conv2D(20, (5, 5), activation='relu'))
3 # apply dropout with rate 0.5
4 model.add(Dropout(0.5))
5 # add 2D pooling layer
6 model.add(MaxPooling2D(pool_size=(2, 2)))
7 # flatten data
8 model.add(Flatten())
9 # add a dense all-to-all relu layer
10 model.add(Dense(20*4*4, activation='relu'))
11 # apply dropout with rate 0.5
12 model.add(Dropout(0.5))
13 # soft-max layer
14 model.add(Dense(num_classes, activation='softmax'))
15
16 # compile the model
17 model.compile(loss=keras.losses.categorical_crossentropy,
18               optimizer='Adam',
19               metrics=['accuracy'])
20
```

Lastly, train your CNN and evaluate the model.

```
In [ ]: 1 # training parameters
2 batch_size = 64
3 epochs = 10
4
5
6 # train CNN
7 model.fit(X_train, Y_train,
8           batch_size=batch_size,
9           epochs=epochs,
10          verbose=1,
11          validation_data=(X_test, Y_test))
12
13 # evaluate model
14 score = model_CNN.evaluate(X_test, Y_test, verbose=1)
15
16 # print performance
17 print()
18 print('Test loss:', score[0])
19 print('Test accuracy:', score[1])
```

5. Do the grid search over any 3 different optimization schemes and 2 activation functions. Suppose that we have a 2 convolutional layers with 10 neurons. Let $p_{\text{dropout}} = 0.5$, $\text{epochs} = 1$, and $\text{batch_size} = 64$. Determine which combination of optimization scheme, activation function, and number of neurons gives the best accuracy.

Doing the grid search requires quite a bit of memory. Please restart the kernel ("Kernel"- "Restart") and re-load the data before doing a new grid search.

```
In [ ]: 1 ...
```

6. Create an arbitrary DNN (you are free to choose any activation function, optimization scheme, etc) and evaluate its performance. Then, add two convolutional layers and pooling layers and evaluate its performance again. How do they compare?

```
In [ ]: 1 ...
```

Problem 2 - Using Tensorflow - Ising Model

You should restart the kernel for Problem 2.

Next, we show how one can use deep neural nets to classify the states of the 2D Ising model according to their phase. This should be compared with the use of logistic-regression in HW8.

The Hamiltonian for the classical Ising model is given by

$$H = -J \sum_{\langle ij \rangle} S_i S_j, \quad S_j \in \{\pm 1\}$$

where the lattice site indices i, j run over all nearest neighbors of a 2D square lattice, and J is some arbitrary interaction energy scale. We adopt periodic boundary conditions. Onsager proved that this model undergoes a phase transition in the thermodynamic limit from an ordered ferromagnet with all spins aligned to a disordered phase at the critical temperature $T_c/J = 2/\log(1 + \sqrt{2}) \approx 2.26$. For any finite system size, this critical point is expanded to a critical region around T_c .

Step 1: Load and Process the Data

We begin by writing a `DataSet` class and two functions `read_data_sets` and `load_data` to process the 2D Ising data.

The `DataSet` class performs checks on the data shape and casts the data into the correct data type for the calculation. It contains a function method called `next_batch` which shuffles the data and returns a mini-batch of a pre-defined size. This structure is particularly useful for the training procedure in TensorFlow.

```
In [ ]: 1 # -*- coding: utf-8 -*-
2 from __future__ import absolute_import, division, print_function
3 import numpy as np
4 seed=12
5 np.random.seed(seed)
6 import sys, os, argparse
7 import tensorflow as tf
8 from tensorflow.python.framework import dtypes
9 # suppress tf compilation warnings
10 os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'
11
12 tf.set_random_seed(seed)
```

```

In [ ]: 1 class DataSet(object):
2
3     def __init__(self, data_X, data_Y, dtype=dtypes.float32):
4         """Checks data and casts it into correct data type. """
5
6         dtype = dtypes.as_dtype(dtype).base_dtype
7         if dtype not in (dtypes.uint8, dtypes.float32):
8             raise TypeError('Invalid dtype %r, expected uint8 or float32' % dtype)
9
10        assert data_X.shape[0] == data_Y.shape[0], ('data_X.shape: %s data_Y.shape: %s' % (data_X.shape, data_Y.shape))
11        self.num_examples = data_X.shape[0]
12
13        if dtype == dtypes.float32:
14            data_X = data_X.astype(np.float32)
15            self.data_X = data_X
16            self.data_Y = data_Y
17
18            self.epochs_completed = 0
19            self.index_in_epoch = 0
20
21        def next_batch(self, batch_size, seed=None):
22            """Return the next `batch_size` examples from this data set."""
23
24            if seed:
25                np.random.seed(seed)
26
27            start = self.index_in_epoch
28            self.index_in_epoch += batch_size
29            if self.index_in_epoch > self.num_examples:
30                # Finished epoch
31                self.epochs_completed += 1
32                # Shuffle the data
33                perm = np.arange(self.num_examples)
34                np.random.shuffle(perm)
35                self.data_X = self.data_X[perm]
36                self.data_Y = self.data_Y[perm]
37                # Start next epoch
38                start = 0
39                self.index_in_epoch = batch_size
40                assert batch_size <= self.num_examples
41            end = self.index_in_epoch
42
43            return self.data_X[start:end], self.data_Y[start:end]

```

Now, load the Ising dataset, and splits it into three subsets: ordered, critical and disordered, depending on the temperature which sets the distribution they are drawn from. Once again, we use the ordered and disordered data to create a training and a test data set for the problem. Classifying the states in the critical region is expected to be harder and we only use this data to test the performance of our model in the end.

```

In [ ]: 1
2 import pickle
3 from sklearn.model_selection import train_test_split
4 from keras.utils import to_categorical
5 import collections
6
7 L=40 # linear system size
8
9 # load data
10 fac = 25
11 file_name = "Ising2DFM_reSample_L40_T=All.pkl" # this file contains 16*10000 samples taken in T=np.arange(0.25,4.0001,0.25)
12 data = pickle.load(open(file_name,'rb')) # pickle reads the file and returns the Python object (1D array, compressed bits)
13 data = data[:,fac]
14 data = np.unpackbits(data).reshape(-1, 1600) # Decompress array and reshape for convenience
15 data=data.astype('int')
16 data[np.where(data==0)]=-1 # map 0 state to -1 (Ising variable can take values +/-1)
17
18 file_name = "Ising2DFM_reSample_L40_T=All_labels.pkl" # this file contains 16*10000 samples taken in T=np.arange(0.25,4.0001,0.25)
19 labels = pickle.load(open(file_name,'rb')) # pickle reads the file and returns the Python object (here just a 1D array with the binary labels)
20
21 # divide data into ordered, critical and disordered
22 X_ordered=data[:int(70000/fac),:]
23 Y_ordered=labels[:70000][::fac]
24
25 X_critical=data[int(70000/fac):int(100000/fac),:]
26 Y_critical=labels[70000:100000][::fac]
27
28 X_disordered=data[int(100000/fac):,: ]
29 Y_disordered=labels[100000:][::fac]
30
31 del data, labels
32
33 # define training and test data sets
34 X=np.concatenate((X_ordered,X_disordered)) #np.concatenate((X_ordered,X_critical,X_disordered))
35 Y=np.concatenate((Y_ordered,Y_disordered)) #np.concatenate((Y_ordered,Y_critical,Y_disordered))
36
37 del X_ordered, X_disordered, Y_ordered, Y_disordered
38

```

```

In [ ]: 1 # pick random data points from ordered and disordered states to create the training and test sets
2 X_train,X_test,Y_train,Y_test=train_test_split(X,Y,train_size=0.6)
3
4
5 # make data categorical
6 Y_train=to_categorical(Y_train)
7 Y_test=to_categorical(Y_test)
8 Y_critical=to_categorical(Y_critical)
9
10
11 # create data sets
12 train = DataSet(X_train, Y_train, dtype=dtypes.float32)
13 test = DataSet(X_test, Y_test, dtype=dtypes.float32)
14 critical = DataSet(X_critical, Y_critical, dtype=dtypes.float32)
15
16 Datasets = collections.namedtuple('Datasets', ['train', 'test', 'critical'])
17 Dataset = Datasets(train=train, test=test, critical=critical)

```

You can load the training data in the following way: (Dataset.train.data_X, Dataset.train.data_Y).

Steps 2+3: Define the Neural Net and its Architecture, Choose the Optimizer and the Cost Function

We can now move on to construct our deep neural net using TensorFlow.

Unique for TensorFlow is creating placeholders for the variables of the model, such as the feed-in data X and Y or the dropout probability `dropout_keepprob` (which has to be set to unity explicitly during testing). Another peculiarity is using the `with scope` to give names to the most important operators. While we do not discuss this here, TensorFlow also allows one to visualise the computational graph for the model (see package documentation on <https://www.tensorflow.org/> (<https://www.tensorflow.org/>)).

The shape of X is only partially defined. We know that it will be a matrix, with instances along the first dimension and features along the second dimension, and we know that the number of features is going to be 28×28 , but we don't know yet how many instances each training batch will contain. So the shape of X is `(None, n_inputs)`. Similarly, we know that Y will be a vector with one entry per instance, but again we don't know the size of the training batch, so the shape is `(None)`.

```
In [ ]: 1 L=40 # system linear size
        2 n_feats=L**2 # 40x40 square lattice
        3 n_categories=2 # 2 Ising phases: ordered and disordered
        4
        5 n_hidden1 = 300
        6 n_hidden2 = 100
        7 n_outputs = 2
        8
        9 with tf.name_scope('data'):
        10 X=tf.placeholder(tf.float32, shape=(None,n_feats))
        11 Y=tf.placeholder(tf.float32, shape=(None,n_categories))
        12 dropout_keepprob=tf.placeholder(tf.float32)
        13
```

To classify whether a given spin configuration is in the ordered or disordered phase, we construct a minimalistic model for a DNN with a single hidden layer containing N_{neurons} (which is kept variable so we can try out the performance of different sizes for the hidden layer).

Let us use a `neuron_layer()` function to create layers in the neural nets.

1. First, create a name scope using the name of the layer.
2. Get the number of inputs by looking up the input matrix's shape and getting the size of the second dimension.
3. Create a W variable which holds the weight matrix (i.e. kernel). Initialize it randomly, using a truncated normal distribution.
4. Create a b variable for biases, initialized to 0.
5. Create a subgraph to compute $Z = XW + b$
6. Use activation function if provided.

```
In [ ]: 1 def neuron_layer(X, n_neuron, name, activation = None):
        2     with tf.name_scope(name):
        3         n_inputs = int(X.get_shape()[1])
        4         stddev = 2 / np.sqrt(n_inputs + n_neuron)
        5         init = tf.truncated_normal([n_inputs, n_neuron], stddev = stddev)
        6         W = tf.Variable(init, name = "kernel")
        7         b = tf.Variable(tf.zeros([n_neuron]), name = "bias")
        8         Z = tf.matmul(X, W) + b
        9         if activation is not None:
        10             return activation(Z)
        11         else:
        12             return Z
```

Using a `neuron_layer()` function, create two hidden layers and an output layer. The first hidden layer takes X as its input, and the second takes the output of the first hidden layer as its input. Finally, the output layer takes the output of the second hidden layer as its input.

```
In [ ]: 1 with tf.name_scope("dnn"):
        2     hidden1 = tf.layers.dense(X, n_hidden1, activation = tf.nn.relu)
        3     hidden2 = tf.layers.dense(hidden1, n_hidden2, activation = tf.nn.relu)
        4     logits = tf.layers.dense(hidden2, n_outputs)
```

Then, define the cost function that we will use to train the neural net model. Here, use the cross entropy to penalize models that estimate a low probability for the target class.

```
In [ ]: 1 with tf.name_scope('loss'):
        2     xentropy = tf.nn.softmax_cross_entropy_with_logits(labels = Y, logits = logits)
        3     loss = tf.reduce_mean(xentropy)
```

Then, define a `GradientDescentOptimizer` that will tweak the model parameters to minimize the cost function. Now, set `learning_rate = 1e-6`.

```
In [ ]: 1 learning_rate = 1e-6
        2 with tf.name_scope('optimiser'):
        3     optimizer = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
        4
```

Lastly, specify how to evaluate the model. Let us simply use accuracy as our performance measure.

```
In [ ]: 1
        2 with tf.name_scope('accuracy'):
        3     correct_prediction = tf.equal(tf.argmax(Y, 1), tf.argmax(logits, 1))
        4     correct_prediction = tf.cast(correct_prediction, tf.float64) # change data type
        5     # correct_prediction = tf.nn.in_top_k(logits, Y, 1)
        6     accuracy = tf.reduce_mean(correct_prediction)
        7
```

Steps 4-5: Train the Model and Evaluate its Performance

We train our DNN using mini-batches of size 100 over a total of 100 epochs, which we define first. We then set up the optimizer parameter dictionary `opt_params`, and use it to create a DNN model.

Running TensorFlow requires opening up a `Session` which we abbreviate as `sess` for short. All operations are performed in this session by calling the `run` method. First, we initialize the global variables in TensorFlow's computational graph by running the `global_variables_initializer`. To train the DNN, we loop over the number of epochs. In each fix epoch, we use the `next_batch` function of the `DataSet` class we defined above to create a mini-batch. The forward and backward passes through the weights are performed by running the `loss` and `optimizer` methods. To pass the mini-batch as well as any other external parameters, we use the `feed_dict` dictionary. Similarly, we evaluate the model performance, by getting `accuracy` on the same minibatch data. Note that the dropout probability for testing is set to unity.

Once we have exhausted all training epochs, we test the final performance on the entire training, test and critical data sets. This is done in the same way as above.

Last, we return the loss and accuracy for each of the training, test and critical data sets.

```

In [ ]: 1 training_epochs=100
        2 batch_size=100
        3
        4 with tf.Session() as sess:
        5
        6     # initialize the necessary variables, in this case, w and b
        7     sess.run(tf.global_variables_initializer())
        8
        9     # train the DNN
        10    for epoch in range(training_epochs):
        11
        12        batch_X, batch_Y = Dataset.train.next_batch(batch_size)
        13
        14        sess.run(optimizer, feed_dict={X: batch_X, Y: batch_Y, dropout_keepprob: 0.5})
        15
        16
        17    # test DNN performance on entire train test and critical data sets
        18    train_loss, train_accuracy = sess.run([loss, accuracy],
        19                                         feed_dict={X: Dataset.train.data_X,
        20                                                  Y: Dataset.train.data_Y,
        21                                                  dropout_keepprob: 0.5}
        22                                         )
        23    print("train loss/accuracy:", train_loss, train_accuracy)
        24
        25    test_loss, test_accuracy = sess.run([loss, accuracy],
        26                                         feed_dict={X: Dataset.test.data_X,
        27                                                  Y: Dataset.test.data_Y,
        28                                                  dropout_keepprob: 1.0}
        29                                         )
        30
        31    print("test loss/accuracy:", test_loss, test_accuracy)
        32
        33    critical_loss, critical_accuracy = sess.run([loss, accuracy],
        34                                                 feed_dict={X: Dataset.critical.data_X,
        35                                                  Y: Dataset.critical.data_Y,
        36                                                  dropout_keepprob: 1.0}
        37                                                 )
        38    print("critical loss/accuracy:", critical_loss, critical_accuracy)
        39
        40

```

Step 6: Modify the Hyperparameters to Optimize Performance of the Model

To study the dependence of our DNN on some of the hyperparameters, we do a grid search over the number of neurons (initially set as 100) in the hidden layer, and different SGD learning rates (initially set as 1e-6). These searches are best done over logarithmically-spaced points.

To do this, define a function for creating a DNN model: `create_DNN` and for evaluating the performance: `evaluate_model`.

The function `grid_search` will output 2D heat map to show how accuracy changes with learning rate and number of neurons.

```

In [ ]: 1 def create_DNN(n_hidden1=100, n_hidden2=100, learning_rate=1e-6):
        2     with tf.name_scope('data'):
        3         X=tf.placeholder(tf.float32, shape=(None,n_feats))
        4         Y=tf.placeholder(tf.float32, shape=(None,n_categories))
        5         dropout_keepprob=tf.placeholder(tf.float32)
        6
        7     with tf.name_scope("dnn"):
        8         hidden1 = tf.layers.dense(X, n_hidden1, activation = tf.nn.relu)
        9         hidden2 = tf.layers.dense(hidden1, n_hidden2, activation = tf.nn.relu)
        10        logits = tf.layers.dense(hidden2, n_outputs)
        11
        12    with tf.name_scope('loss'):
        13        xentropy = tf.nn.softmax_cross_entropy_with_logits(labels = Y, logits = logits)
        14        loss = tf.reduce_mean(xentropy)
        15
        16    with tf.name_scope('optimiser'):
        17        optimizer = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
        18
        19    with tf.name_scope('accuracy'):
        20        correct_prediction = tf.equal(tf.argmax(Y, 1), tf.argmax(logits, 1))
        21        correct_prediction = tf.cast(correct_prediction, tf.float64) # change data type
        22        # correct_prediction = tf.nn.in_top_k(logits, Y, 1)
        23        accuracy = tf.reduce_mean(correct_prediction)
        24
        25    return X, Y, dropout_keepprob, loss, optimizer, accuracy
        26

```



```

In [ ]: 1 def evaluate_model(neurons,lr):
2
3     training_epochs=100
4     batch_size=100
5
6
7     X, Y, dropout_keepprob, loss, optimizer, accuracy = create_DNN(n_hidden1=neurons, n_hidden2=neurons, learning_rate=lr)
8     with tf.Session() as sess:
9
10        # initialize the necessary variables, in this case, w and b
11        sess.run(tf.global_variables_initializer())
12
13        # train the DNN
14        for epoch in range(training_epochs):
15
16            batch_X, batch_Y = Dataset.train.next_batch(batch_size)
17
18            sess.run(optimizer, feed_dict={X: batch_X,Y: batch_Y,dropout_keepprob: 0.5})
19
20
21        # test DNN performance on entire train test and critical data sets
22        train_loss, train_accuracy = sess.run([loss, accuracy],
23                                             feed_dict={X: Dataset.train.data_X,
24                                                         Y: Dataset.train.data_Y,
25                                                         dropout_keepprob: 0.5}
26                                             )
27        print("train loss/accuracy:", train_loss, train_accuracy)
28
29        test_loss, test_accuracy = sess.run([loss, accuracy],
30                                             feed_dict={X: Dataset.test.data_X,
31                                                         Y: Dataset.test.data_Y,
32                                                         dropout_keepprob: 1.0}
33                                             )
34
35        print("test loss/accuracy:", test_loss, test_accuracy)
36
37        critical_loss, critical_accuracy = sess.run([loss, accuracy],
38                                                    feed_dict={X: Dataset.critical.data_X,
39                                                            Y: Dataset.critical.data_Y,
40                                                            dropout_keepprob: 1.0}
41                                                    )
42        print("critical loss/accuracy:", critical_loss, critical_accuracy)
43
44    return train_loss,train_accuracy,test_loss,test_accuracy,critical_loss,critical_accuracy

```

```

In [ ]: 1 def grid_search():
2     """This function performs a grid search over a set of different learning rates
3     and a number of hidden layer neurons."""
4
5     # perform grid search over learnign rate and number of hidden neurons
6     N_neurons=[100, 200, 300, 400, 500]
7     learning_rates=np.logspace(-6,-1,6)
8
9     # pre-allocate variables to store accuracy and loss data
10    train_loss=np.zeros((len(N_neurons),len(learning_rates)),dtype=np.float64)
11    train_accuracy=np.zeros_like(train_loss)
12    test_loss=np.zeros_like(train_loss)
13    test_accuracy=np.zeros_like(train_loss)
14    critical_loss=np.zeros_like(train_loss)
15    critical_accuracy=np.zeros_like(train_loss)
16
17    # do grid search
18    for i, neurons in enumerate(N_neurons):
19        for j, lr in enumerate(learning_rates):
20
21            print("training DNN with %4d neurons and SGD lr=%0.6f." %(neurons,lr) )
22
23            train_loss[i,j],train_accuracy[i,j],\
24            test_loss[i,j],test_accuracy[i,j],\
25            critical_loss[i,j],critical_accuracy[i,j] = evaluate_model(neurons,lr)
26
27
28    plot_data(learning_rates,N_neurons,train_accuracy, "training data")
29    plot_data(learning_rates,N_neurons,test_accuracy, "test data")
30    plot_data(learning_rates,N_neurons,critical_accuracy, "critical data")

```

```

In [ ]: 1 %matplotlib notebook
2 import matplotlib.pyplot as plt
3
4 def plot_data(x,y,data, title):
5
6     # plot results
7     fontsize=16
8
9
10    fig = plt.figure()
11    ax = fig.add_subplot(111)
12    cax = ax.matshow(data, interpolation='nearest', vmin=0, vmax=1)
13    fig.colorbar(cax)
14
15    # put text on matrix elements
16    for i, x_val in enumerate(np.arange(len(x))):
17        for j, y_val in enumerate(np.arange(len(y))):
18            c = "{0:.1f}%".format( 100*data[j,i])
19            ax.text(x_val, y_val, c, va='center', ha='center')
20
21    # convert axis vaues to to string labels
22    x=[str(i) for i in x]
23    y=[str(i) for i in y]
24
25
26    ax.set_xticklabels(['']+x)
27    ax.set_yticklabels(['']+y)
28
29    ax.set_xlabel('$\\mathrm{learning\\ rate}$',fontsize=fontsize)
30    ax.set_ylabel('$\\mathrm{hidden\\ neurons}$',fontsize=fontsize)
31
32    ax.set_title(title,fontsize=fontsize)
33
34    plt.tight_layout()
35
36    plt.show()

```

```
In [ ]: 1 grid_search()
```

1. Do the grid search over 5 different types of activation functions (https://www.tensorflow.org/api_guides/python/nn#Activation_Functions (https://www.tensorflow.org/api_guides/python/nn#Activation_Functions)). Evaluate the performance for each case and determine which gives the best accuracy. You can assume an arbitrary DNN. Show results for training, test, and critical data.

```
In [ ]: 1 ...
```

2. Do the grid search over 5 different numbers of epochs and batch sizes. Make a 2D heat map as shown in the example. You can assume an arbitrary DNN. Show results for training, test, and critical data.

```
In [ ]: 1 ...
```

Problem 3 - SDSS galaxies

You should restart the kernel for Problem 2.

The data is provided in the file "specz_data.txt". The columns of the file (length of 13) correspond to - spectroscopic redshift ('zspec'), RA, DEC, magnitudes in 5 bands - u, g, r, i, z (denoted as 'mu,' 'mg,' 'mr,' 'mi,' 'mz' respectively); Exponential and de Vaucouleurs model magnitude fits ('logExp' and 'logDev' <http://www.sdss.org/dr12/algorithms/magnitudes/> (<http://www.sdss.org/dr12/algorithms/magnitudes/>); zebra fit ('pz_zebra'); Neural Network fit ('pz_NN') and its error estimate ('pz_NN_Err')

We will undertake 2 exercises -

- Regression
 - We will use the magnitude of object in different bands ('mu, mg, mr, mi, mz') and do a regression exercise to estimate the redshift of the object. Hence our feature space is 5.
 - The correct redshift is given by 'zspec', which is the spectroscopic redshift of the object. We will use this for training and testing purpose.

Sidenote: Photometry vs. Spectroscopy

The amount of energy we receive from celestial objects – in the form of radiation – is called the flux, and an astro- nomical technique of measuring the flux is photometry. Flux is usually measured over broad wavelength bands, and with the estimate of the distance to an object, it can infer the object's luminosity, temperature, size, etc. Usually light is passed through colored filters, and we measure the intensity of the filtered light.

On the other hand, spectroscopy deals with the spectrum of the emitted light. This tells us what the object is made of, how it is moving, the pressure of the material in it, etc. Note that for faint objects making photometric observation is much easier.

Photometric redshift (photoz) is an estimate of the distance to the object using photometry. Spectroscopic redshift observes the object's spectral lines and measures their shifts due to the Doppler effect to infer the distance.

- Classification
 - We will use the same magnitudes and now also the redshift of the object ('zspec') to classify the object as either Elliptical or Spiral. Hence our feature space is now 6.
 - The correct class is given by comparing 'logExp' and 'logDev' which are the fits for Exponential and Devocular profiles. If logExp > logDev, its a spiral and vice-versa. We will use this for training and testing purpose. Since the classes are not explicitly given, generate a column for those (Classes can be ± 1. If it is 0, it does not belong to either of the class.)

Cleaning

Read in the files to create the data (X and Y) for both regression and classification.

You will have to clean the data -

- Drop the entries that are nan or infinite
- Drop the unrealistic numbers such as 999, -999; and magnitudes that are unrealistic. Since these are absolute magnitudes, they should be positive and high. Lets choose a magnitude limit of 15 as safe bet.
- For classification, drop the entries that do not belong to either of the class

```
In [ ]: 1 #Read in and create data
2
3 fname = 'specz_data.txt'
4 spec_dat=np.genfromtxt(fname,names=True)
5 print(spec_dat.dtype.fields.keys())
6 #convenience variable
7 zspec = spec_dat['zspec']
8 pzNN = spec_dat['pz_NN']
9 #some N redshifts are not defined
10 pzNN[pzNN < 0] = np.nan
11
12 #For Regression
13 bands = ['u', 'g', 'r', 'i', 'z' ]
14 mlim = 15
15
16 xdata = np.concatenate([spec_dat['m%s%i' for i in bands]].T
17 bad = (xdata[:, 0] < mlim) | (xdata[:, 1] < mlim) | (xdata[:, 2] < mlim) & (xdata[:, 3] < mlim) | (xdata[:, 4] < mlim)
18 xdata = xdata[~bad]
19 xdata[xdata<0] = 0
20 ydata = zspec[~bad]
21
22 #For classification
23 classes = np.sign(spec_dat['logExp'] - spec_dat['logDev'])
24 tmp = np.concatenate([spec_dat['m%s%i' for i in bands]].T
25 xxdata = np.concatenate([tmp, zspec.reshape(-1, 1)], axis=1)
26 bad = (classes==0) | (xxdata[:, 0] < mlim) | (xxdata[:, 1] < mlim) | (xxdata[:, 2] < mlim) & (xxdata[:, 3] < mlim) | (xxdata[:, 4] < mlim)
27 xxdata = xxdata[~bad]
28 classes = classes[~bad]
```

For regression, X and Y data (called "xdata" and "ydata," respectively) is cleaned magnitudes (5 feature space) and spectroscopic redshifts respectively. For classification, X and Y data (called "xxdata" and "classes" respectively) is cleaned magnitudes+spectroscopic redshifts respectively (6 feature space) and classees respectively.

```
In [ ]: 1 print('For Regression:')
2 print('Before: Size of datasets is ', zspec.shape[0])
3 print('After: Size of datasets is ', xdata.shape[0])
4 print('')
5 print('For Classification:')
6 print('Before: Size of datasets is ', zspec.shape[0])
7 print('After: Size of datasets is ', xxdata.shape[0])
```

Visualization

The next step should be to visualize the data.

For regression

- Make a histogram for the distribution of the data (spectroscopic redshift).
- Make 5 2D histograms of the distribution of the magnitude as function of redshift (Hint: https://matplotlib.org/devdocs/api/_as_gen/matplotlib.axes.Axes.hist2d.html (https://matplotlib.org/devdocs/api/_as_gen/matplotlib.axes.Axes.hist2d.html))

For classification

- Make 6 1-d histogram for the distribution of the data (6 features - zspec and 5 magnitudes) for both class 1 and -1 separately

1. Make histograms for both regression and classification.

```
In [ ]: 1 ...
```

2. Do the following preprocessing:

Preprocessing:

- Next, split the sample into training data and the testing data. We will be using the training data to train different algorithms and then compare the performance over the testing data. In this project, keep 80% data as training data and uses the remaining 20% data for testing.
- Often, the data can be ordered in a specific manner, hence shuffle the data prior to splitting it into training and testing samples.
- Many algorithms are also not scale invariant, and hence scale the data (different features to a uniform scale). All this comes under preprocessing the data. <http://scikit-learn.org/stable/modules/preprocessing.html#preprocessing> Use StandardScaler from sklearn (or write your own routine) to center the data to 0 mean and 1 variance. Note that you only center the training data and then use its mean and variance to scale the testing data before using it.

Hint: How to get a scaled training data:

1. Let the training data be: train = ("training X data", "training Y data")
2. You can first define a StandardScaler:
scale_xdata, scale_ydata = preprocessing.StandardScaler(), preprocessing.StandardScaler()
3. Then, do the fit:
for regression: scale_xdata.fit(train_regression[0]), scale_ydata.fit(train_regression[1].reshape(-1, 1))
for classification: scale_xdata.fit(train_classification[0])
Here, no need to fit for y data for classification (it's either +1 or -1. Already scaled)
4. Next, transform:
for regression: scaled_train_data = (scale_xdata.fit_transform(train_regression[0]), scale_ydata.fit_transform(train_regression[1].reshape(-1, 1)))
for classification: scaled_train_data = (scale_xdata.fit_transform(train_classification[0]), train_classification[1])
Again, y data is already scaled for classification.

Do this for test data as well.

```
In [ ]: 1 from sklearn import preprocessing
```

```
In [ ]: 1 ...
```

Metrics

The last remaining preparatory step is to write metric for gauging the performance of the algorithm. Write a function to calculate the 'RMS' error given (y_predict, y_truth) to gauge regression and another function to evaluate accuracy of classification.

In addition, for classification, we will also use confusion matrix.

Below is an example you can use. Feel free to write you own.

```
In [ ]: 1 from sklearn.metrics import confusion_matrix
2
3 def rms(x, y, scale1=None, scale2=None):
4     '''Calculate the rms error given the truth and the prediction
5     '''
6     mask = np.isfinite(x[:]) & np.isfinite(y[:])
7     if scale1 is not None:
8         x = scale1.inverse_transform(x)
9     if scale2 is not None:
10        y = scale2.inverse_transform(y)
11    return np.sqrt(np.mean((x[mask] - y[mask])** 2))
12
13 def acc(x, y):
14     '''Calculate the accuracy given the truth and the prediction
15     '''
16    mask = np.isfinite(x[:]) & np.isfinite(y[:])
17    return (x == y).sum()/x.size
18
```

Hyperparameter method

Now, we will be varying hyperparameters to get the best model and build some intuition. There are various ways to do this and we will use Grid Search methodology (as you did in Problem 1 and 2) which simply tries all the combinations along with some cross-validation scheme. For most part, we will use 4-fold cross validation.

Sklearn provides GridSearchCV functionality for this purpose.

Its recommended to spend some time to go through output format of GridSearchCV and write some utility functions to make the recurring plots for every parameter.

Grid Search returns a dictionary with self explanatory keys for the most part. Mostly, the keys correspond to (masked) numpy arrays of size = #(all possible combination of parameters). The value of individual parameter in every combination is given in arrays with keys starting from 'param_-' and this should help you to match the combination with the corresponding scores.

For masked arrays, you can access the data values by using ".data"

Do not overwrite these grid search-ed variables (and not only their result) since we will compare all the models together in the end

```
In [ ]: 1 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
2 # http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
3 # http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.RandomizedSearchCV.html
```

Method 1. k Nearest Neighbors

For regression, let us play with grid search using knn to tune hyperparameters. (<https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsRegressor.html>) Consider the following 3 hyperparameters -

- Number of neighbors ([2, 3, 5, 10, 15, 20, 25, 50, 100])
- Weights of leaves (Uniform or Inverse Distance weighing)
- Distance metric (Euclidian or Manhattan distance - parameter 'p')

1. Do a grid search on these parameters. List the combination of hyperparameters you tried and evaluate the accuracy (mean test score) and its standard deviation. Which gives the highest accuracy value?

```
In [ ]: 1 from sklearn.neighbors import KNeighborsRegressor
```

Hint: (Read the documentations carefully for more detail.)

First, define the hyperparameters: parameters = {'n_neighbors':[2, 3, 5, 10, 15, 20, 25, 50, 100], 'weights':['uniform', 'distance'], 'p':[1, 2]}

Specify the algorithm you want to use: e.g. knnr = KNeighborsRegressor()

Then, Do a grid search on these parameters using 4 fold cross validation: gcknn = GridSearchCV(knnr, parameters, cv=4)

Do the fit: gcknn.fit("scaled_training_data")

(Let "scaled_training_data" be the training data where "scaled_training_data = ("train X data", "train Y data")"

Get results: results = gcknn.cv_results_

cv_results_ has the following dictionaries: "rank_test_score," "mean_test_score," "std_test_score," and "params" (See http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html (http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html))

Then, you can evaluate the models based on "rank_test_score" and print out their "params," along with their "mean_test_score" and "std_test_score".

```
In [ ]: 1 ...
```

2. Also print out fitting and scoring times for all hyperparameter combinations.

Plot timings for fitting and scoring

Hint: Assume that you got results from: results = gcknn.cv_results_

Then, get the scoring time: results['mean_score_time']

and the fitting time: results['mean_fit_time']

```
In [ ]: 1 ...
```

3. Based on the results you obtained in Part 1 and 2, answer the following questions

- Is it always better to use more neighbors?
- Is it better to weigh the leaves, if yes, which distance metric performs better?
- GridCV returns fitting and scoring time for every combination. You will find that scoring time is higher than training time. Why do you think is that the case?

Answer:

4. Which parameters seem to affect the performance most? To better answer this question, make plots of the mean test score for each hyperparameter.

Hint: Suppose you have two types of hyperparameters: A and B. Let A = [1, 2] and B = [1, 2, 4, 7, 10].

Then, you have 20 different combination of hyperparameters.

Let A = 1. Then, you can try (A,B) = (1,1), (1,2), (1,4), (1,7), (1,10) Suppose that the mean score you got for the above combination is [0.7, 0.72, 0.75, 0.77, 0.8]. Similarly, for A = 2, you tried (A,B) = (2,1), (2,2), (2,4), (2,7), (2,10) and obtained the mean score of [0.8, 0.82, 0.85, 0.87, 0.9].

To better see how changing the value of parameter A affects the performance, you can make the following plot:

```
In [ ]: 1 A_1 = [0.7, 0.72, 0.75, 0.77, 0.8]
2 A_2 = [0.8, 0.82, 0.85, 0.87, 0.9]
3
4 plt.plot(A_1, label = "A=1")
5 plt.plot(A_2, label = "A=2")
6 plt.ylabel("mean test score")
7 plt.legend()
8 plt.show()
```

This is the plot of the mean test score for A marginalizing over B.

Similarly, make a plot of the mean test score for each kNN hyperparameter.

```
In [ ]: 1 ...
```

5. You have determined the best combination of hyperparameters and CV schemes. Predict the test y data using the GridSearchCV method. Use the "rms" metric function we defined earlier and calculate the rms error on the test data.

Hint: To determine the rms error, you need:

Truth: given from data (test_data[1])

Prediction: gridsearch.predict(test_data[0]) (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html))

```
In [ ]: 1 ...
```

Classification

```
In [ ]: 1 from sklearn.neighbors import KNeighborsClassifier
2 # http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html
```

Here we will look at 4 different type of cross-validation schemes -

- Kfold
- Stratified Kfold
- Shuffle Split
- Stratified Shuffle Split

6. Assuming the list of hyperparameters from Part 1, do 4 different grid searches. From Part 1, take top 5 combination of hyperparameters which gives you the highest accuracy value. Rank the performance of CV schemes for each combination.

```
In [ ]: 1 from sklearn.model_selection import KFold, StratifiedKFold, ShuffleSplit, StratifiedShuffleSplit
```

```
In [ ]: 1 parameters = {'n_neighbors':[2, 3, 5, 10, 15, 20, 25, 50, 100], 'weights':['uniform', 'distance'], 'p':[1, 2]}
        2 knnc = KNeighborsClassifier()
        3
        4 #Grid Search
        5 gc = GridSearchCV(knnc, parameters, cv=KFold(4, random_state=100))
        6 #Do the fit
        7 ...
        8
        9 gc2 = GridSearchCV(knnc, parameters, cv=StratifiedKFold(4, random_state = 100))
        10 #Do the fit
        11 ...
        12
        13 gc3 = GridSearchCV(knnc, parameters, cv=ShuffleSplit(4, 0.1, random_state = 100))
        14 #Do the fit
        15 ...
        16
        17 gc4 = GridSearchCV(knnc, parameters, cv=StratifiedShuffleSplit(4, 0.1, random_state = 100))
        18 #Do the fit
        19 ...
```

7. Answer the following questions:

- Are the conclusions different for any parameter from the regression case?
- Does the mean accuracy change for different CV scheme?
- Does the standard deviation in mean accuracy change?

```
In [ ]: 1 ...
```

Answer:

8. Using the best combination of hyperparameters and CV schemes you have found, compute the confusion matrix (https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html) and evaluate the accuracy.

Hint: To get a confusion matrix, you need both truth (available from data) and prediction (can be computed using .predict function from GridSearchCV (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)).

```
In [ ]: 1 ...
```

Method 2. Random Forests

The most important feature of the random forest is the number of trees in the ensemble. We will also play with the maximum depth of the trees.

Try:

```
n_estimators = [10, 50, 150, 200, 300]
max_depth = [10, 50, 100]
```

```
In [ ]: 1 from sklearn.ensemble import RandomForestRegressor
        2 # http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html
```

1. Do the grid search over n_estimators and max_depth. List the combination of hyperparameters you tried and evaluate the accuracy (mean test score) and its standard deviation. Which gives the highest accuracy value?

```
In [ ]: 1 rf = RandomForestRegressor()
        2 parameters = ...
        3
        4 gcrf = GridSearchCV(rf, parameters, cv=5)
        5
        6 ...
```

2. Which parameters seem to affect the performance most? To better answer this question, make plots of the mean test score for each hyperparameter. (plot the mean test score of n_estimators marginalizing over max_depth, etc)

```
In [ ]: 1 ...
```

3. Based on the results you obtained in Part 1, answer the following questions:

- Are the scores of these models statistically different? Based on this, which architecture will you choose for your model?
- For every parameter, make the plot for fitting time. Based on this and the previous question, how many trees do you recommend keeping in the ensemble?

```
In [ ]: 1 ...
```

Answer:

4. You have determined the best combination of hyperparameters. Predict the test y data using the GridSearchCV method. Use the "rms" metric function we defined earlier and calculate the rms error on the test data.

```
In [ ]: 1 ...
```

Classification

```
In [ ]: 1 from sklearn.ensemble import RandomForestClassifier
        2 # http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html
```

```
In [ ]: 1 #Grid search (This will take few minutes)
        2
        3 rfc = RandomForestClassifier()
        4 parameters = ...
        5
        6 gcrfc = GridSearchCV(rfc, parameters, cv=StratifiedShuffleSplit(4, 0.1, random_state = 100))
        7
        8 ...
```

5. Assuming the list of hyperparameters from Part 1, do the grid search using StratifiedShuffleSplit CV scheme. List the combination of hyperparameters you tried and evaluate the accuracy (mean test score) and its standard deviation. Which gives the highest accuracy value?

```
In [ ]: 1 ...
```

6. Using the best combination of hyperparameters, compute the confusion matrix (https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html) and evaluate the accuracy.

```
In [ ]: 1 ...
```

To Submit

Execute the following cell to submit. If you make changes, execute the cell again to resubmit the final copy of the notebook, they do not get updated automatically.

We recommend that all the above cells should be executed (their output visible) in the notebook at the time of submission.

Only the final submission before the deadline will be graded.

```
In [ ]: 1 _ = ok.submit()
```