# lab2_neural_partial

May 2, 2021

# 1 Lab: Model Order Selection for Neural Data

Machine learning is a key tool for neuroscientists to understand how sensory and motor signals are encoded in the brain. In addition to improving our scientific understanding of neural phenomena, understanding neural encoding is critical for brain machine interfaces. In this lab, you will use model selection for performing some simple analysis on real neural signals.

In doing this lab, you will learn to: * Represent neural time-series data in arrays * Load data from a pickle file * Describe and fit memoryless linear models * Describe and fit linear time-series models with delays * Fit linear models with multiple target outputs * Select the optimal delay via cross-validation

## 1.1 Loading the data

The data in this lab comes from neural recordings described in:

Stevenson, Ian H., et al. "Statistical assessment of the stability of neural movement representations." Journal of neurophysiology 106.2 (2011): 764-774

Neurons are the basic information processing units in the brain. Neurons communicate with one another via *spikes* or *action potentials* which are brief events where voltage in the neuron rapidly rises then falls. These spikes trigger the electro-chemical signals between one neuron and another. In this experiment, the spikes were recorded from 196 neurons in the primary motor cortex (M1) of a monkey using an electrode array implanted onto the surface of a monkey's brain. During the recording, the monkey performed several reaching tasks and the position and velocity of the hand was recorded as well.

The goal of the experiment is to try to *read the monkey's brain*: That is, predict the hand motion from the neural signals from the motor cortex.

We first load the key packages.

```
[1]: import numpy as np
     import matplotlib.pyplot as plt
     import pickle

     from sklearn.linear_model import LinearRegression
     from sklearn.metrics import r2_score
```

The full data is available on the CRCNS website http://crcns.org/data-sets/movements/dream. This website has a large number of great datasets and can be used for projects as well. However,

the raw data files can be quite large. To make the lab easier, the Kording lab at UPenn has put together an excellent repository where they have created simple pre-processed versions of the data. For this lab, I have created a simplified version of the dataset that you can download with the following code. This may take a little while to download since the file is 26 MB.

```python
fn_src = 'https://www.dropbox.com/sh/n4924ipcfjqc0t6/AADOv9JYMUBK1tlg9P71gSSra/
    →example_data_s1.pickle?dl=1'
fn_dst = 'example_data_s1.pickle'

import os
from six.moves import urllib

if os.path.isfile(fn_dst):
    print('File %s is already downloaded' % fn_dst)
else:
    urllib.request.urlretrieve(fn_src, fn_dst)
```

```
File example_data_s1.pickle is already downloaded
```

The file is a *pickle* data structure, which is a package to serialize python objects into data files. Once you have downloaded the file, you can run the following command to retrieve the data from the pickle file.

```python
with open('example_data_s1.pickle', 'rb') as fp:
    X,y = pickle.load(fp)
```

The matrix `X` is matrix of spike counts where `X[i,j]` is the number of spikes from neuron `j` in time bin `i`. The matrix `y` has two columns: * `y[i,0]` = velocity of the monkey's hand in the x-direction * `y[i,1]` = velocity of the monkey's hand in the y-direction Our goal will be to predict `y` from `X`.

Each time bin represent `tsamp=0.05` seconds of time. Using `X.shape` and `y.shape` compute and print: * `nt` = the total number of time bins * `nneuron` = the total number of neurons * `nout` = the total number of output variables to track = number of columns in `y` * `ttotal` = total time of the experiment is seconds.

```python
tsamp = 0.05   # sampling time in seconds

nt = X.shape[0]
print('The total number of time bins is %d.' % nt)

nneuron = X.shape[1]
print('The total number of neurons is %d.' % nneuron)

nout = y.shape[1]
print('The total number of output variables to track is %d.' % nout)

ttotal = tsamp * nt
print('Total time of the experiment is %f seconds.' % ttotal)
```

```
The total number of time bins is 61339.
The total number of neurons is 52.
The total number of output variables to track is 2.
Total time of the experiment is 3066.950000 seconds.
```

## 1.2  Fitting a Memoryless Linear Model

Let's first try a simple linear regression model to fit the data.

First, use the `train_test_split` function to split the data into training and test. Let `Xtr,ytr` be the training data set and `Xts,yts` be the test data set. Use `test_size=0.33` so 1/3 of the data is used for test.

```
[5]: from sklearn.model_selection import train_test_split

Xtr, Xts, ytr, yts = train_test_split(X, y, test_size = 0.33)
```

Now, fit a linear model using `Xtr,ytr`. Make a prediction `yhat` using `Xts`. Compare `yhat` to `yts` to measure `rsq`, the R^2. You can use the `r2_score` method. Print the `rsq` value. You should get `rsq` of around 0.45.

```
[6]: lm = LinearRegression()
lm.fit(Xtr, ytr)

yhat = lm.predict(Xts)
rsq = r2_score(yts, yhat)
print(np.round(rsq, 2))
```

```
0.46
```

It is useful to plot the predicted vs. true values. Since we have two outputs, create two `subplots` using the `plt.subplot()` command. In plot `i=0,1`, plot `yhat[:,i]` vs. `yts[:,i]` with a scatter plot. Label the axes of the plots. You may also use the command:

```
plt.figure(figsize=(10,5))
```
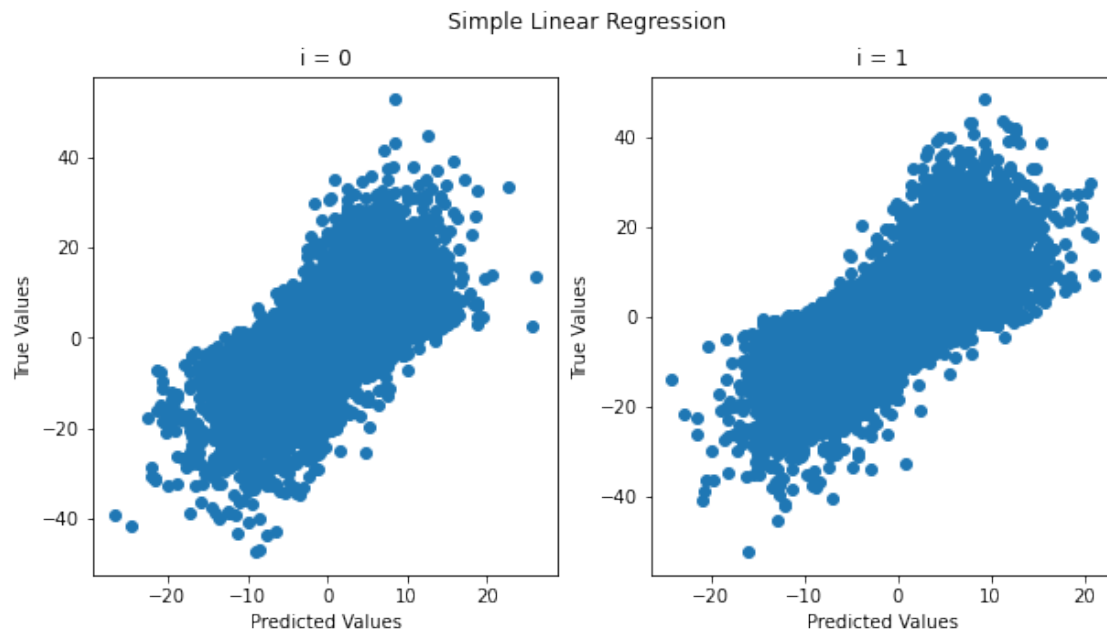
to make the figures a little larger.

```
[7]: plt.figure(figsize = (10, 5))

plt.subplot(1, 2, 1)
plt.suptitle('Simple Linear Regression')

plt.scatter(yhat[:, 0], yts[:, 0])
plt.title('i = 0')
plt.xlabel('Predicted Values')
plt.ylabel('True Values')

plt.subplot(1, 2, 2)
plt.scatter(yhat[:, 1], yts[:, 1])
plt.title('i = 1')
```

```
plt.xlabel('Predicted Values')
plt.ylabel('True Values')

plt.show()
```

Simple Linear Regression

## 1.3 Fitting Models with Delay

One way we can improve the model accuracy is to used delayed version of the features. Specifically, the model we used above mapped the features

```
yhat[i,k] = \sum_{j=0}^{p-1} X[i,j]*w[j,k] + b[k]
```

where `p` is the number of features and `w[j,k]` is a matrix of coefficients. In this model, `yhat[i,:]` at time `i` was only dependent on the inputs `X[i,:]` at time `i`. In signal processing, this is called a *memoryless* model. However, in many physical systems, such as those that arise in neuroscience, there is a delay between the inputs `X[i,:]` and the outputs `y[i]`. For such cases, we can use a model of the form,

```
yhat[i+d,k] = \sum_{j=0}^{p-1} \sum_{m=0}^d X[i+m,j]*W[j,m,k] + b[k]
```

where `W` is a 3-dim array of coefficients where:

```
W[j,m,k] is the influence of the input X[i+m,j] onto output y[i+d,k]
```

In signal processing, this model is called an *FIR* filter and `W[j,:,k]` is the *impulse response* from the `j`-th input to the `k`-th output. The point is that the output at time `i+d` depends on the inputs at times `i,i+1,...,i+d`. Hence, it depends on the last `d+1` time steps, not just the most recent time.

To translate this into a linear regression problem, complete the following function that creates a new feature and target matrix where:

```
Xdly[i,:] has the rows X[i,:], X[i++1,:], ..., X[i+dly,:]
ydly[i,:] = y[i+dly,:]
```

Thus, `Xdly[i,:]` contains all the delayed fetaures for the target `yhat`. Note that if X is `n x p` then `Xdly` will be `n-dly x (dly+1)*p`.

```
[8]: def create_dly_data(X, y, dly):

         """
         Create delayed data
         """

         Xdly = np.zeros((X.shape[0] - dly, (dly+1) * X.shape[1]))
         ydly = np.zeros((X.shape[0] - dly, y.shape[1]))

         for i in range(Xdly.shape[0]):
             X_new = X[i,:]
             for j in range(1, dly+1):
                 X_new = np.hstack([X_new, X[i+j, :]])
             Xdly[i,:] = X_new
             ydly[i,:] = y[i+dly, :]

         return Xdly, ydly
```

Now fit an linear delayed model with `dly=6` additional delay lags. That is, * Create delayed data `Xdly,ydly=create_dly_data(X,y,dly=6)` * Split the data into training and test as before * Fit

the model on the training data * Measure the `R^2` score on the test data

If you did this correctly, you should get a new `R^2` score around 0.69. This is significantly better than the memoryless models.

```
[9]: Xdly,ydly = create_dly_data(X, y, dly = 6)
```

```
[10]: Xtr, Xts, ytr, yts = train_test_split(Xdly, ydly, test_size = 0.33)
```

```
[11]: lm = LinearRegression()
      lm.fit(Xtr, ytr)

      yhat = lm.predict(Xts)
      rsq = r2_score(yts, yhat)
      print(np.round(rsq, 2))
```

```
0.69
```

Plot the predicted vs. true values as before. You should visually see a better fit.
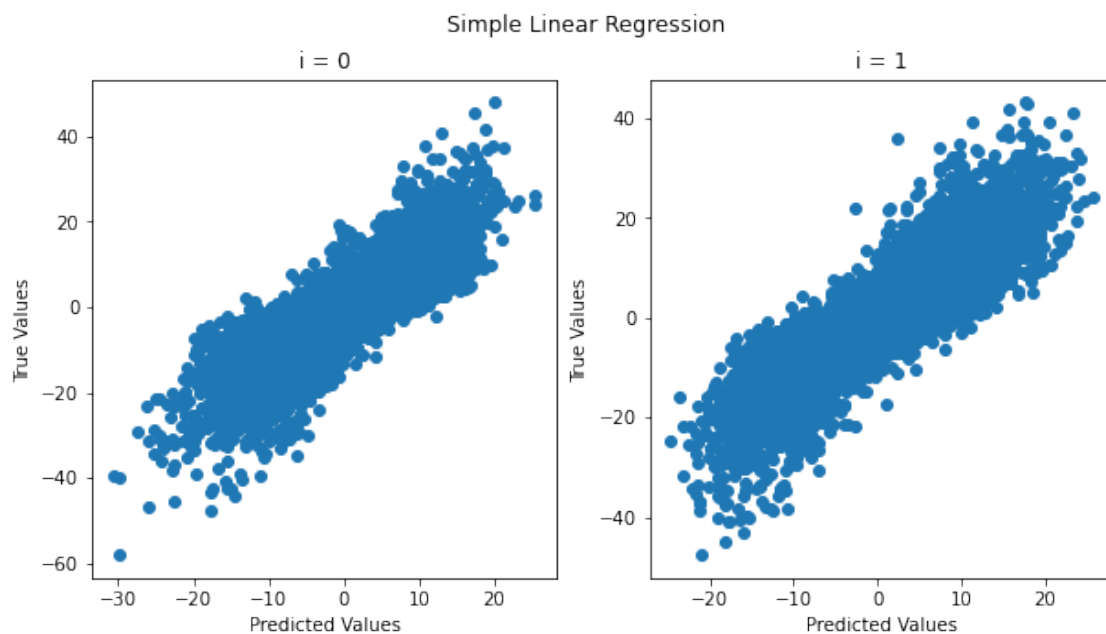
```
[12]: plt.figure(figsize = (10, 5))

      plt.subplot(1, 2, 1)
      plt.suptitle('Simple Linear Regression')

      plt.scatter(yhat[:, 0], yts[:, 0])
      plt.title('i = 0')
      plt.xlabel('Predicted Values')
      plt.ylabel('True Values')

      plt.subplot(1, 2, 2)
      plt.scatter(yhat[:, 1], yts[:, 1])
      plt.title('i = 1')
      plt.xlabel('Predicted Values')
      plt.ylabel('True Values')

      plt.show()
```

Simple Linear Regression

*Note*: Fitting an FIR model with the above method is very inefficient when the number of delays, `dly` is large. In the above method, the number of columns of `X` grows from `p` to `(dly+1)*p` and the computations become expensive with `dly` is large. We will describe a much faster way to fit such models using gradient descent when we talk about convolutional neural networks.

## 1.4 Selecting the Optimal Delay via Model Order Selection

In the previous example, we fixed `dly=6`. We can now select the optimal delay using model order selection. Since we have a large number of data samples, it turns out that the optimal model order uses a very high delay. Using the above fitting method, the computations take too long. So, to simplify the lab, we will first just pretend that we have a very limited data set.

Compute `Xred` and `yred` by taking the first `nred=6000` samples of the data `X` and `y`. This is about 10% of the overall data.

```
[13]: nred = 6000

Xred = X[:nred, :]
yred = y[:nred, :]
```

Now complete the following code to implement K-fold cross validation with `nfold=5` and values of delays `dtest = [0,1,...,dmax]`.

The code also includes a progress bar using the `tqdm` package. This is very useful when you have a long computation.

Note: Some students appeared to use the `mse` metric (i.e. RSS per sample) instead of `R^2`. That is fine. For the solution, I have computed both.

```
[14]: import  sklearn.model_selection
import tqdm.notebook

nfold = 5  # Number of folds
dmax = 15  # maximum number of delays

# Create a k-fold object
kf = sklearn.model_selection.KFold(n_splits = nfold)

# Model orders to be tested
dtest = np.arange(dmax + 1)
nd = len(dtest)

# Initialize a matrix Rsq to hold values of the R^2 across the model orders and
 ↪folds.
# Alternatively, you can also create an RSS matrix

Rsq = np.zeros((nd, nfold))

# Create a progress bar.  Note there are nd*nfold total fits.
pbar = tqdm.notebook.tqdm(
    total=nfold*nd, initial=0,
    unit='fits', unit_divisor=nd, desc='Model order test')


for it, d in enumerate(dtest):
```

```
    # Create the delayed data using the create_dly_function from the reduced
    →data Xred, yred
    Xdly, ydly = create_dly_data(Xred, yred, d)

    # Loop over the folds
    for isplit, Ind in enumerate(kf.split(Xdly)):

        # Get the training data in the split
        Itr, Its = Ind

        # Split the data (Xdly,ydly) into training and test
        Xtr = Xdly[Itr]
        ytr = ydly[Itr]
        Xts = Xdly[Its]
        yts = ydly[Its]

        # Fit data on training data
        lm = LinearRegression()
        lm.fit(Xtr, ytr)
        yhat = lm.predict(Xts)

        # Measure the R^2 vale on test data and store in the matrix Rsq
        Rsq[it, isplit] = r2_score(yts, yhat)

        pbar.update(1)

pbar.close()
```

Model order test:    0%|              | 0/80 [00:00<?, ?fits/s]

Compute the mean and standard error of the R^2 values as a function of the model order d. Use a plt.errorbar plot. Label your axes.

```
[15]: mean_rsq = np.mean(Rsq, axis = 1)
      std_rsq = np.std(Rsq, axis = 1) / np.sqrt(nfold - 1)

      plt.errorbar(dtest, mean_rsq, yerr = std_rsq)
      plt.xlabel('Number of Delay')
      plt.ylabel('R Square')
      plt.title('Optimal Delay via Model Order Selection')
      plt.grid()
```

Optimal Delay via Model Order Selection

Find the optimal order `d` with the normal rule (i.e. highest test `R^2`)

```
[16]: d_max = np.argmax(mean_rsq)
      print(d_max)
```

8

Now find the optimal model order via the one SE rule (i.e. highest test `R^2` within on SE)

```
[17]: tgt_rsq = mean_rsq[d_max] - std_rsq[d_max]
      sorted(mean_rsq >= tgt_rsq)
      d_oneSE = np.where(mean_rsq >= tgt_rsq)
      print(d_oneSE[0][0])
```

6