

# 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.

Before doing this lab, you should review the ideas in the [polynomial model selection demo](#). In addition to the concepts in that demo, 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

## 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.

In [187...

```
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. You can download the file `example_data_s1.pickle` from the [Dropbox link](#).

Alternatively, you can directly run the following code. This may take a little while to download since the file is 26 MB.

In [188...

```
fn_src = 'https://www.dropbox.com/sh/n4924ipcfjqc0t6/AAD0v9JYMUBK1t1g9P71gSSra/exampl
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.

In [189...

```
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 in seconds.

In [190...

```
tsamp = 0.05 # sampling time in seconds

# TODO
nt = X.shape[0]
nneuron = X.shape[1]
nout = y.shape[1]
ttotal = nt * tsamp
print("The total number of time bins: ", nt)
print("The total numbner of neurons: ", nneuron)
print("The total number of output variables to track", nout)
print("The total time of the experiment in seconds", ttotal)
```

```
The total number of time bins: 61339
The total numbner of neurons: 52
The total number of output variables to track 2
The total time of the experiment in seconds 3066.9500000000003
```

## 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.

In [191...

```
from sklearn.model_selection import train_test_split

# TODO
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`.

In [192...

```
# TODO
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
regr = LinearRegression()
regr.fit(Xtr, ytr)
yhat = regr.predict(Xts)

rsq = r2_score(yts, yhat)
print("The Coefficient of Determination is {}".format(rsq))
```

The Coefficient of Determination is 0.4629643766736126

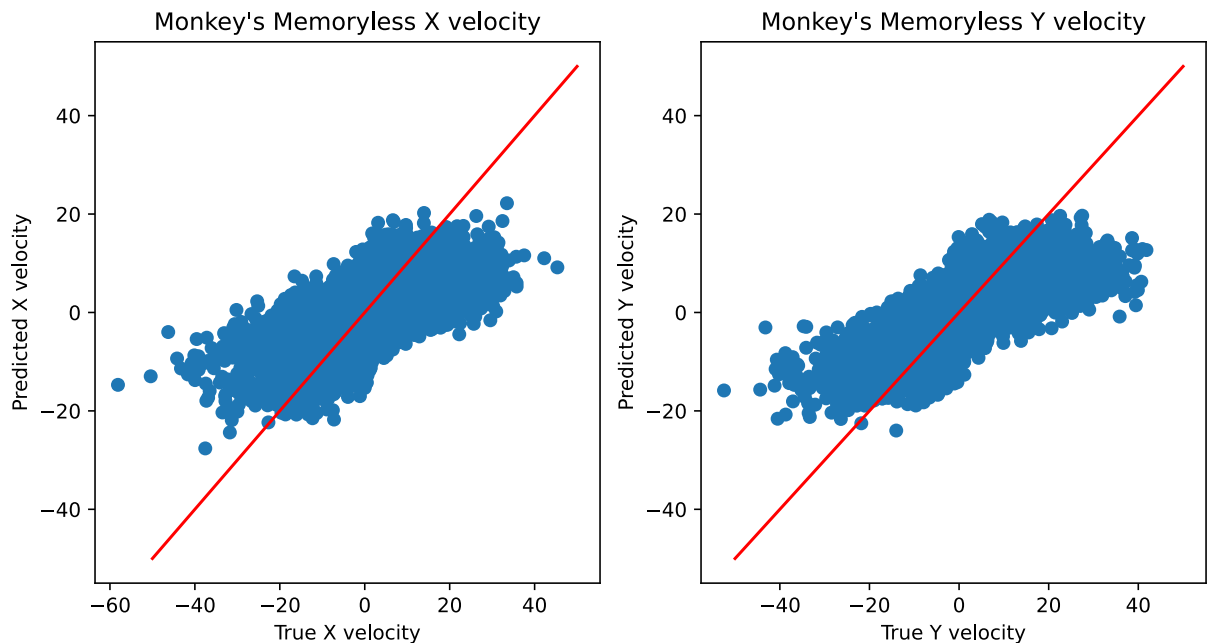
It is useful to plot the predicted vs. true values. Since we have two outputs, create two subplots using the `plt.subplot()` command. In plot `i=0,1`, plot `yhat[:,i]` vs. `yts[:,i]` with a scatter plot. Label the axes of the plots. You may also use the command:

```
plt.figure(figsize=(10,5))
```

to make the figures a little larger.

In [193...

```
# TODO
plt.figure(figsize=(10,5))
plt.subplot(1, 2, 1)
plt.scatter(yts[:,0], yhat[:,0])
plt.plot(np.linspace(-50, 50, 100), np.linspace(-50, 50, 100), '-r')
plt.xlabel("True X velocity")
plt.ylabel("Predicted X velocity")
plt.title("Monkey's Memoryless X velocity")
plt.subplot(1, 2, 2)
plt.scatter(yts[:,1], yhat[:,1])
plt.plot(np.linspace(-50, 50, 100), np.linspace(-50, 50, 100), '-r')
plt.xlabel("True Y velocity")
plt.ylabel("Predicted Y velocity")
plt.title("Monkey's Memoryless Y velocity")
plt.show()
```



## Fitting Models with Delay

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

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

where  $p$  is the number of features and  $w[j,k]$  is a matrix of coefficients. In this model,  $\hat{y}[i,:]$  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,

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

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  $i, i+1, \dots, i+d$ . 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 features for the target `yhat`. Note that if `X` is  $n \times p$  then `Xdly` will be  $n-dly \times (dly+1)*p$ .

In [194...

```
def create_dly_data(X,y,dly):
    """
    Create delayed data
    """
    n = X.shape[0]
    p = X.shape[1]
    print("Shape of Matrix X before transformation: ", X.shape)
    print("Shape of Matrix y before transformation: ", y.shape)
    # TODO
    Xdly = np.zeros((n - dly, (dly + 1) * p))
    ydly = np.zeros((n - dly, y.shape[1]))
    for i in range(Xdly.shape[0]):
        Xdly[i,:] = np.hstack([X[i+j,:] for j in range(dly+1)])
        ydly[i,:] = y[i+dly,:]
    print("Shape of Matrix X after transformation: ", Xdly.shape)
    print("Supppose shape of Xdly", (n - dly, (dly + 1) * p))
    print("Shape of Matrix y after transformation: ", ydly.shape)
    print("Supppose shape of ydly", (n - dly, y.shape[1]))
    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.

In [195...

```
# TODO
Xdly, ydly = create_dly_data(X, y, dly=6)
Xdly_tr, Xdly_ts, ydly_tr, ydly_ts = train_test_split(Xdly, ydly, test_size=0.33)
regr.fit(Xdly_tr, ydly_tr)
yhat_dly = regr.predict(Xdly_ts)
rsq_dly = r2_score(ydly_ts, yhat_dly)
print("The Coefficient of Determination with a Delay of 6 seconds is: ", rsq_dly)
```

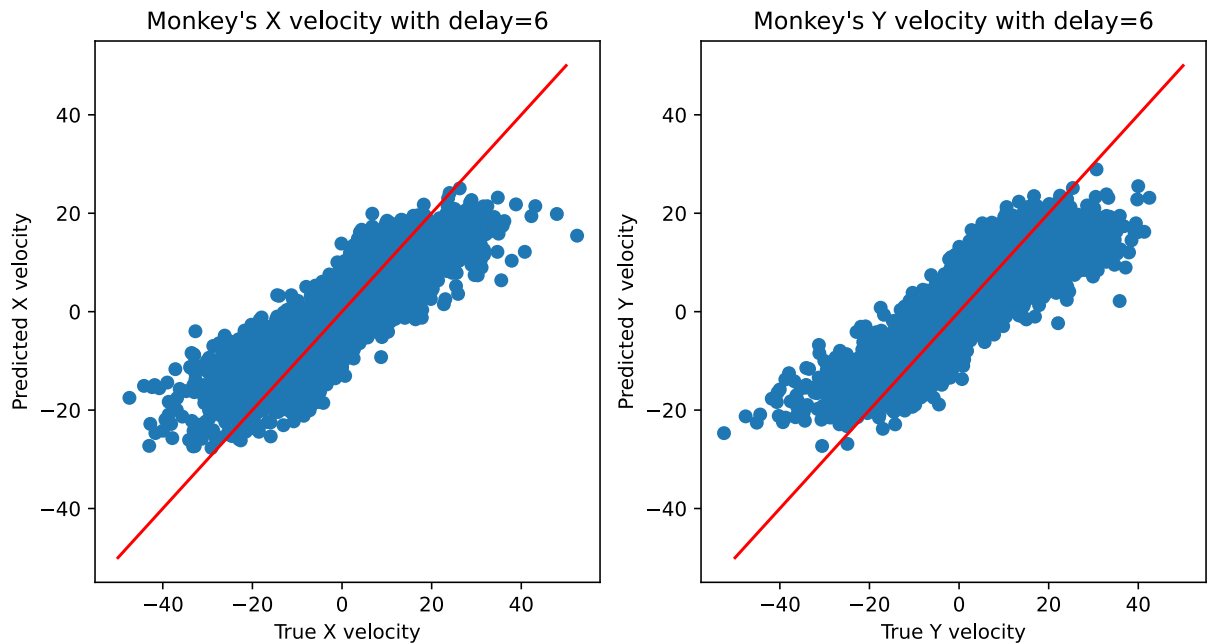
```
Shape of Matrix X before transformation: (61339, 52)
Shape of Matrix y before transformation: (61339, 2)
Shape of Matrix X after transformation: (61333, 364)
Supppose shape of Xdly (61333, 364)
Shape of Matrix y after transformation: (61333, 2)
Supppose shape of ydly (61333, 2)
The Coefficient of Determination with a Delay of 6 seconds is: 0.6907844460153875
```

Plot the predicted vs. true values as before. You should visually see a better fit.

In [196...

```
# TODO
plt.figure(figsize=(10,5))
plt.subplot(1, 2, 1)
plt.scatter(ydly_ts[:,0], yhat_dly[:,0])
plt.plot(np.linspace(-50, 50, 100),np.linspace(-50, 50, 100), '-r')
plt.xlabel("True X velocity")
plt.ylabel("Predicted X velocity")
plt.title("Monkey's X velocity with delay=6")
plt.subplot(1, 2, 2)
```

```
plt.scatter(ydly_ts[:,1], yhat_dly[:,1])
plt.plot(np.linspace(-50, 50, 100), np.linspace(-50, 50, 100), '-r')
plt.xlabel("True Y velocity")
plt.ylabel("Predicted Y velocity")
plt.title("Monkey's Y velocity with delay=6")
plt.show()
```



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.

## 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  $X_{red}$  and  $y_{red}$  by taking the first  $n_{red}=6000$  samples of the data  $X$  and  $y$ . This is about 10% of the overall data.

In [197...

```
nred = 6000

# TODO
Xred = X[:nred]
yred = y[:nred]
```

Now complete the following code to implement K-fold cross validation with  $n_{fold}=5$  and values of delays  $d_{test} = [0, 1, \dots, d_{max}]$ .

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.

In [198...

```
import sklearn.model_selection
import tqdm.notebook

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

# TODO: Create a k-fold object
kf = sklearn.model_selection.KFold(n_splits=nfold, shuffle=True)

# TODO: Model orders to be tested
dtest = range(dmax+1)
nd = len(dtest)

# TODO.
# Initialize a matrix Rsq to hold values of the R^2 across the model orders and fold
# Alternatively, you can also create an RSS matrix
RSQ_kfold_test = 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):
    # TODO:
    # Create the delayed data using the create_dly_function from the reduced
    # data Xred, yred
    Xdly, ydly = create_dly_data(Xred, yred, dly=d)

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

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

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

        # TODO: Fit data on training data
        regr.fit(Xtr, ytr)
        yhat_kfold = regr.predict(Xts)
        # TODO: Measure the R^2 vale on test data and store in the matrix Rsq
        rsq_kfold = r2_score(yts, yhat_kfold)
        RSQ_kfold_test[it, isplit] = rsq_kfold

    pbar.update(1)
pbar.close()
```

```
Shape of Matrix X before transformation: (6000, 52)
Shape of Matrix y before transformation: (6000, 2)
Shape of Matrix X after transformation: (6000, 52)
Supppose shape of Xdly (6000, 52)
Shape of Matrix y after transformation: (6000, 2)
Supppose shape of ydly (6000, 2)
Shape of Matrix X before transformation: (6000, 52)
Shape of Matrix y before transformation: (6000, 2)
```

Shape of Matrix X after transformation: (5999, 104)  
 Suppose shape of Xdly (5999, 104)  
 Shape of Matrix y after transformation: (5999, 2)  
 Suppose shape of ydly (5999, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5998, 156)  
 Suppose shape of Xdly (5998, 156)  
 Shape of Matrix y after transformation: (5998, 2)  
 Suppose shape of ydly (5998, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5997, 208)  
 Suppose shape of Xdly (5997, 208)  
 Shape of Matrix y after transformation: (5997, 2)  
 Suppose shape of ydly (5997, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5996, 260)  
 Suppose shape of Xdly (5996, 260)  
 Shape of Matrix y after transformation: (5996, 2)  
 Suppose shape of ydly (5996, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5995, 312)  
 Suppose shape of Xdly (5995, 312)  
 Shape of Matrix y after transformation: (5995, 2)  
 Suppose shape of ydly (5995, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5994, 364)  
 Suppose shape of Xdly (5994, 364)  
 Shape of Matrix y after transformation: (5994, 2)  
 Suppose shape of ydly (5994, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5993, 416)  
 Suppose shape of Xdly (5993, 416)  
 Shape of Matrix y after transformation: (5993, 2)  
 Suppose shape of ydly (5993, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5992, 468)  
 Suppose shape of Xdly (5992, 468)  
 Shape of Matrix y after transformation: (5992, 2)  
 Suppose shape of ydly (5992, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5991, 520)  
 Suppose shape of Xdly (5991, 520)  
 Shape of Matrix y after transformation: (5991, 2)  
 Suppose shape of ydly (5991, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5990, 572)  
 Suppose shape of Xdly (5990, 572)  
 Shape of Matrix y after transformation: (5990, 2)  
 Suppose shape of ydly (5990, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5989, 624)  
 Suppose shape of Xdly (5989, 624)  
 Shape of Matrix y after transformation: (5989, 2)  
 Suppose shape of ydly (5989, 2)  
 Shape of Matrix X before transformation: (6000, 52)  
 Shape of Matrix y before transformation: (6000, 2)  
 Shape of Matrix X after transformation: (5988, 676)  
 Suppose shape of Xdly (5988, 676)  
 Shape of Matrix y after transformation: (5988, 2)



```

Supppose shape of ydly (5988, 2)
Shape of Matrix X before transformation: (6000, 52)
Shape of Matrix y before transformation: (6000, 2)
Shape of Matrix X after transformation: (5987, 728)
Supppose shape of Xdly (5987, 728)
Shape of Matrix y after transformation: (5987, 2)
Supppose shape of ydly (5987, 2)
Shape of Matrix X before transformation: (6000, 52)
Shape of Matrix y before transformation: (6000, 2)
Shape of Matrix X after transformation: (5986, 780)
Supppose shape of Xdly (5986, 780)
Shape of Matrix y after transformation: (5986, 2)
Supppose shape of ydly (5986, 2)
Shape of Matrix X before transformation: (6000, 52)
Shape of Matrix y before transformation: (6000, 2)
Shape of Matrix X after transformation: (5985, 832)
Supppose shape of Xdly (5985, 832)
Shape of Matrix y after transformation: (5985, 2)
Supppose shape of ydly (5985, 2)

```

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.

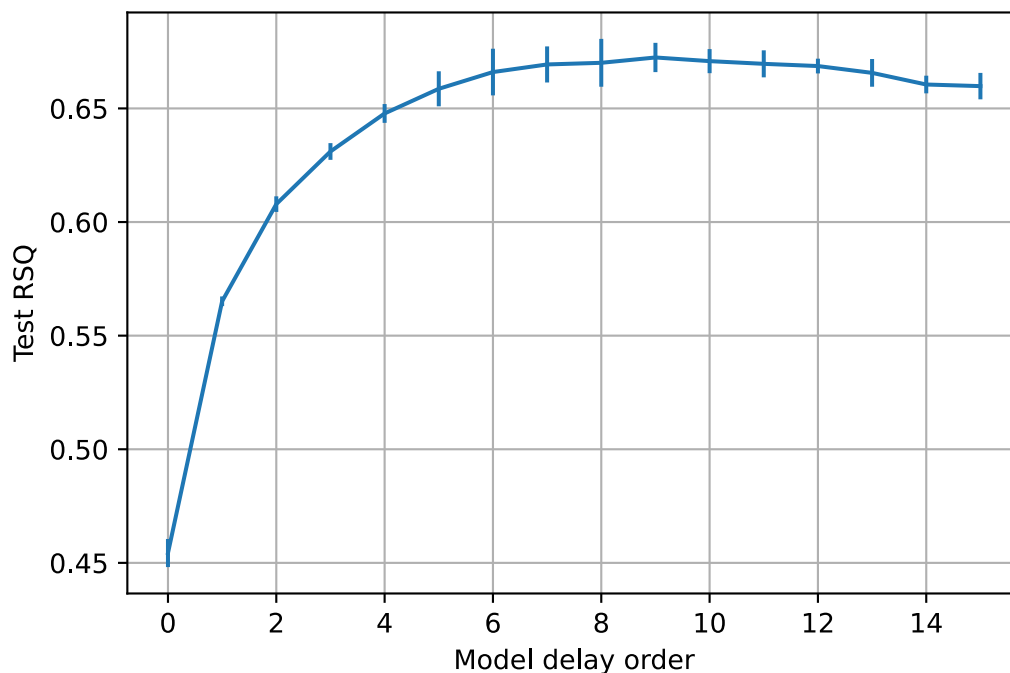
In [199...

```

# TODO
RSQ_mean = np.mean(RSQ_kfold_test, axis=1)
RSS_std = np.std(RSQ_kfold_test, axis=1) / np.sqrt(nfold-1)

plt.errorbar(dtest, RSQ_mean, yerr=RSS_std, fmt='-.')
plt.xlabel("Model delay order")
plt.ylabel("Test RSQ")
plt.grid(True)

```



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

In [200...

```

# TODO
imax = np.argmax(RSQ_mean)
print("The estimated model order with normal rule is {}".format(imax))

```

The estimated model order with normal rule is 9

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

In [201...

```
# TODO
RSQ_tgt = RSQ_mean[imax] + RSS_std[imax]
print("The Target RSQ is ", RSQ_tgt)
I = np.where(RSQ_mean <= RSQ_tgt)[0]
iopt = np.argmax(RSQ_mean[I])
dopt = dtest[iopt]

# Print results
print("The estimated model order is {}".format(dopt))
print("The test RSQ for this model is {}".format(RSQ_mean[iopt]))
```

The Target RSQ is 0.6788661608250948

The estimated model order is 9

The test RSQ for this model is 0.6724243819747516

In [ ]: