

Exercise List

Chapter 2,3

---

**Exercise 1.** Use *sklearn's* implementation of *k-Nearest Neighbors* for regression purposes, which is found in *sklearn.neighbors.KNeighborsRegressor*. You will find the best value of *k* using 10 fold CrossValidation (CV), which is found in *sklearn.model\_selection.KFold*.

- a) You will modify the python code below to generate 1000 data points, or alternatively you could use part of your semester project dataset if it is related to regression.

```
import numpy as np
from matplotlib import pyplot as plt

def genDataSet (N) :
    x = np . random . normal (0 , 1 , N)
    ytrue = (np.cos ( x ) + 2) / (np.cos ( x * . 4 ) + 2)
    noise = np . random . normal (0 , 0 . 2 , N)
    y = ytrue + noise
    return x , y , ytrue

x , y , ytrue = genDataSet(100)
plt.plot(x , y , '.' )
plt.plot(x , ytrue , 'rx')
plt.show( )
```

- b) Using 10-fold CV, you will report the three best values of *k-neighbors* that yield the best CV  $E_{out}$ . You will vary the values of *k* in the following range:  $k = 1, 3, 5, \dots, 2\lceil \frac{N+1}{2} \rceil - 1$ .
- c) You will report the best CV  $E_{out}$ .

*Solution. a)*

```
import numpy as np

from matplotlib import pyplot as plt


def genDataSet (N) :
```

```
x = np . random . normal (0 , 1 , N)

ytrue = (np.cos ( x ) + 2) / (np.cos ( x * . 4 ) + 2)

noise = np . random . normal (0 , 0 . 2 , N)

y = ytrue + noise

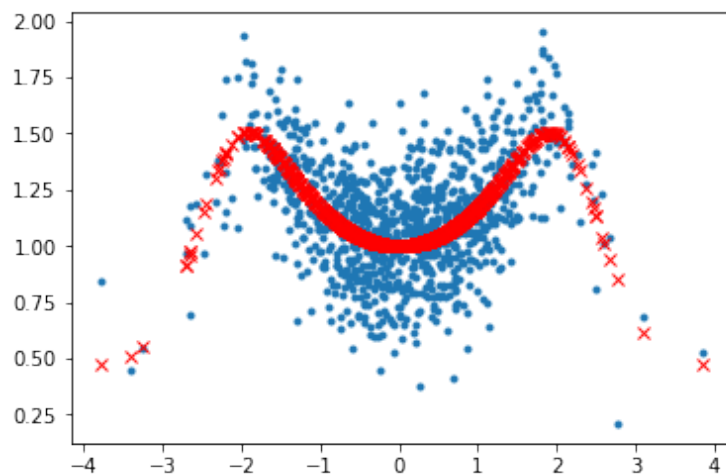
return x , y , ytrue

x , y , ytrue = genDataSet(1000) #changed this number to 1000

plt.plot(x , y , '.' )

plt.plot(x , ytrue , 'rx')

plt.show( )
```



b) and c)

```
from sklearn.neighbors import KNeighborsRegressor

from sklearn.metrics import mean_squared_error

from sklearn.model_selection import KFold

kf = KFold(n_splits=10)
```

```
X = np.array([x,y]).T
kf.get_n_splits(X)

msek = {}
bestk= 0
bestmse = 100000
for k in np.arange(1,2*np.floor(((len(ytrue)*0.9)+1)/2)-1,2):

    print(k)

    mse = []

    for train_index , test_index in kf.split(X):

        X_train , X_test =      X[train_index] ,      X[test_index]

        y_train , y_test = ytrue[train_index] , ytrue[test_index]

        neigh = KNeighborsRegressor(n_neighbors = int(k))

        neigh.fit(X_train , y_train)

        predict = neigh.predict(X_test)

        mse.append(mean_squared_error(y_test , predict))

    print(np.mean(mse))

    msek[k] = np.mean(mse)

#part c

if bestmse > msek[k]:

    bestmse = msek[k]

    bestk = k
```

```

    print(k)

    print(bestmse)

plt.plot(k,msek[k], 'r. ') #mse

plt.show()

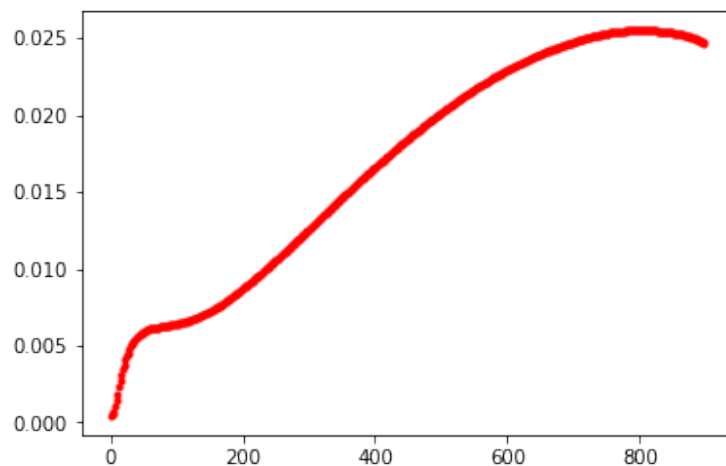
print(bestk)

```

The best three best  $k$  using mse as an evaluation metric are:

$k$	1.0	3.0	5.0
MSE	0.00045241177636707794	0.000512687006146458	0.0006638327509173002

This can also be seen by plotting the MSE by  $k$ .



The smallest MSE is the best, making smaller  $k$ s better than larger  $k$ s.

**Exercise 2.** Using the same dataset you just tried in the previous problem, repeat the experiment 100 times storing the best three  $k$  number of neighbors in every single trial, and at the end of all trials plot a histogram of all the values of  $k$  that you saved.

*Solution.*

```

from sklearn.neighbors import KNeighborsRegressor

```

```
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import KFold

kf = KFold(n_splits=10)
bestkall = []
for i in range(100):
    x, y, ytrue = genDataSet(1000)
    X = np.array([x,y]).T
    kf.get_n_splits(X)

    msek = {}
    bestk= 0
    bestmse = 100000
    for k in np.arange(1,2*np.floor(((len(ytrue)*0.9)+1)/2)-1,2):
        mse = []

        for train_index, test_index in kf.split(X):
            X_train, X_test = X[train_index], X[test_index]
            y_train, y_test = ytrue[train_index], ytrue[test_index]

            neigh = KNeighborsRegressor(n_neighbors = int(k))
            neigh.fit(X_train, y_train)

            predict = neigh.predict(X_test)
```

```
mse.append(mean_squared_error(y_test , predict))

mse[k] = np.mean(mse)

if bestmse > mse[k]:

    bestmse = mse[k]

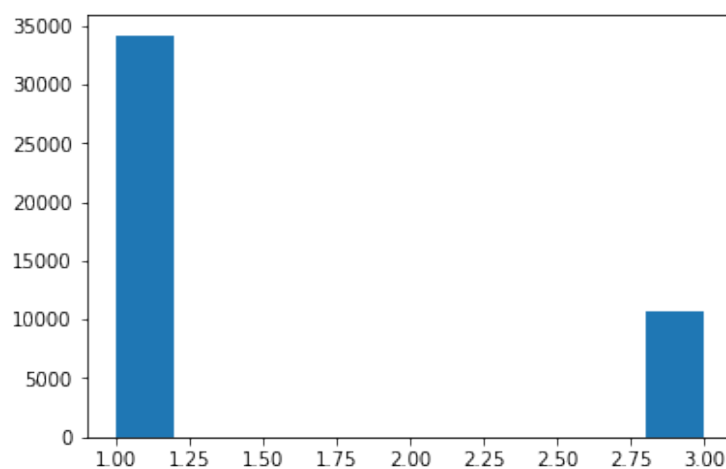
    bestk = k

    print(k)

    print(bestmse)

bestkall.append(bestk)

plt.hist(bestkall)
```



---

**Exercise 3. Experiment with k-means for color quantization.**

Using *sklearn*'s implementation of k-means, find the best color clustering for an image of your choice. A good portrait of yourself could be fun (just sayin'). Color quantization is the science behind compression of images. The idea is to represent an image with fewer colors than the original. The experiment consists of the following steps:

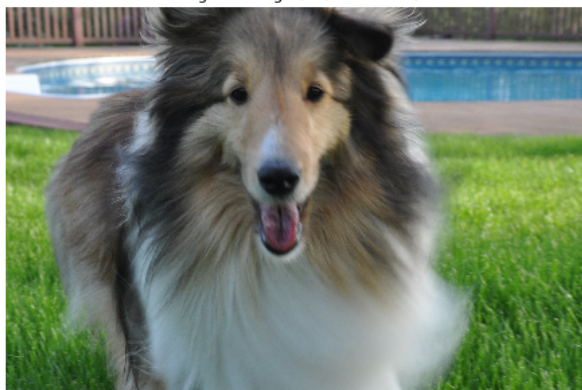
- Download to your computer the file `hw3.kmeans.img.py`
- In the same folder where you downloaded the program, save a copy of your picture for experimentation.
- Go to line 23 and set `n_colors` with your choice of a number of colors between 2 and 64. This number is actually the number of clusters (or  $k$ ) we are searching for in an

unsupervised manner using *kmeans*.

- d) Then go to line 25 and replace the image file name with the name of your image file. This is where your image is read into a numpy array.
- e) Run the program. Observe the result. If the result does not look funny to you, repeat from 3.(c) until it does. Then, report your result image along with your answers to the following questions:
  - i) Explain what happens when you increase or decrease the value of *n\_colors*?
  - ii) Explain in what other possible applications do you think this can be useful?
  - iii) Why do you think the resulting picture was funny at the end?

*Solution.* This is my puppy Penny

Original image (96,615 colors)



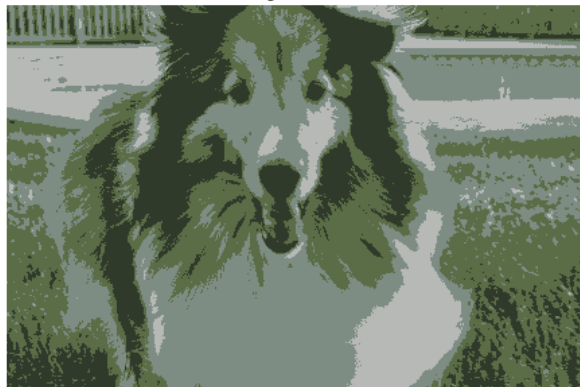
Quantized image (20 colors, K-Means)



Quantized image (7 colors, K-Means)



Quantized image (4 colors, K-Means)



- i) When I increase the *n\_color* there are more colors in the picture, when I decrease it there are less colors. When there is only 2 *n\_colors* then the picture has only two colors which end up being the averages of two groups of colors related to each other

ii) This would be helpful in displaying images on devices that do not have the ability to represent a lot of different colors. It decreases the complexity and it only takes away some of the detail. One place where this is helpful is on old video game handheld devices.

iii) The picture looks cool to me. It is more cartoon-ish but less bright with color. With less group the color turns more grayish green and more washed out.

#### Exercise 4. **Neural networks: The MLP**

- a) Download the python program `hw3.MLP.sol.py` which implements a 10 fold cross validation approach to find the best number of neurons and the best learning parameter  $\eta$  (eta) in an MLP for regression. For learning purposes you could download `hw3.MLP.py` first, which is a simple implementation of the MLP.
- b) Download the python program `hw3_4_a_gendata.py` which generates random data.
- c) Run the program in 4.(a) for 1,000 samples, and then take note of the best number of neurons and  $\eta$ . Go here <https://goo.gl/forms/QFmaNWYaFLcPWdim2> and report your results. You can do it as many times as you want, but at least one is required.
- d) Explain your results. What do you think is happening? What is your interpretation of the number of neurons with respect to the performance of the network?
- e) (extra credit) Repeat 4.(c)-(d) but for 10,000 samples.

*Solution. a)*

```
import hw3_4_a_gendata

import numpy as np

import matplotlib.pyplot as plt

from sklearn.neural_network import MLPRegressor

from sklearn.model_selection import KFold

# generates data & split it into X (training input) and y (target output)
X, y = hw3_4_a_gendata.genDataSet(100)

neurons = 1000 # <- number of neurons in the hidden layer
```



```

eta = 0.01          # <- the learning rate parameter

# here we create the MLP regressor

mlp = MLPRegressor(hidden_layer_sizes=(neurons,), verbose=True, learning_rate=eta)

# here we train the MLP

mlp.fit(X, y)

# E-out in training

print("Training_set_score: %f" % mlp.score(X, y))

# now we generate new data as testing set and get E-out for testing set

X, y = hw3_4_a_gendata.genDataSet(100)

print("Testing_set_score: %f" % mlp.score(X, y))

ypred = mlp.predict(X)

plt.plot(X[:, 0], X[:, 1], '-')

plt.plot(X[:, 0], y, '-r')

plt.plot(X[:, 0], ypred, '-k')

plt.show()

```

b)

```

import hw3_4_a_gendata

import numpy as np

import matplotlib.pyplot as plt

from sklearn.neural_network import MLPRegressor

from sklearn.model_selection import KFold

```

```
# number of samples
N = 1000

# generate data & split it into X (training input) and y (target output)
X, y = hw3_4_a_gendata.genDataSet(N)

# linear regression solution
w=np.linalg.pinv(X.T.dot(X)).dot(X.T).dot(y)

#neurons <- number of neurons in the hidden layer
#eta <- the learning rate parameter

bestNeurons=0
bestEta=0
bestScore=float('-inf')
score=0
for neurons in range(1,101,1):
    for eta in range(1,11,1):
        eta=eta/10.0
        kf = KFold(n_splits=10)
```

```

cvscore=[]

for train, validation in kf.split(X):

    X_train, X_validation, y_train, y_validation = X[train, :], X[validation

    # here we create the MLP regressor

    mlp = MLPRegressor(hidden_layer_sizes=(neurons,), verbose=False, learn

    # here we train the MLP

    mlp.fit(X_train, y_train)

    # now we get E_out for validation set

    score=mlp.score(X_validation, y_validation)

    cvscore.append(score)


# average CV score

score=sum(cvscore)/len(cvscore)

if (score > bestScore):

    bestScore=score

    bestNeurons=neurons

    bestEta=eta

    print("Neurons_" + str(neurons) + ",_eta_" + str(eta) + "._Testing_set_

# here we get a new training dataset

X, y = hw3_4_a_gendata.genDataSet(N)

# here we create the final MLP regressor

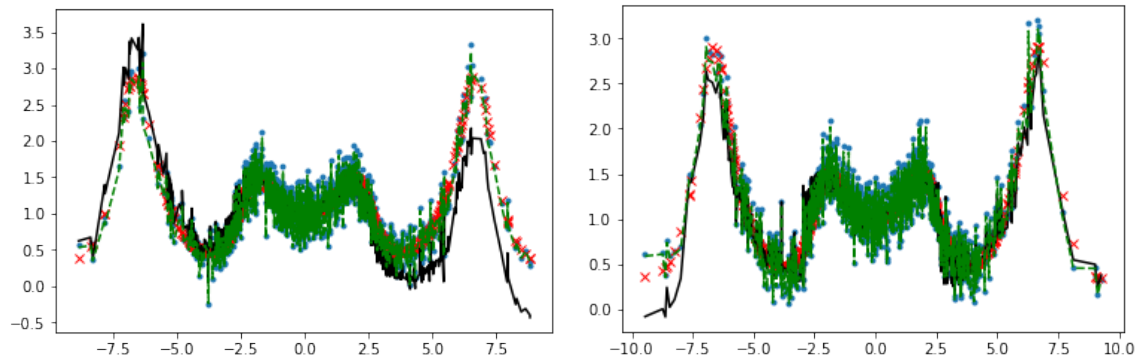
```

```

mlp = MLPRegressor(hidden_layer_sizes=(bestNeurons, ), verbose=True, learning_rate_milestones=(100000, 200000, 400000, 800000, 1600000, 3200000, 6400000, 12800000, 25600000, 51200000, 102400000, 204800000, 409600000, 819200000, 1638400000, 3276800000, 6553600000, 13107200000, 26214400000, 52428800000, 104857600000, 209715200000, 419430400000, 838860800000, 1677721600000, 3355443200000, 6710886400000, 13421772800000, 26843544000000, 53687088000000, 107374176000000, 214748352000000, 429496704000000, 858993408000000, 1717986816000000, 3435973632000000, 6871947264000000, 13743894528000000, 27487789056000000, 54975578112000000, 109951156224000000, 219902312448000000, 439804624896000000, 879609249792000000, 1759218499584000000, 3518436999168000000, 7036873998336000000, 14073747996672000000, 28147495993344000000, 56294991986688000000, 112589983973376000000, 225179967946752000000, 450359935893504000000, 900719871787008000000, 1801439743574016000000, 3602879487148032000000, 7205758974296064000000, 14411517948592128000000, 28823035897184256000000, 57646071794368512000000, 115292143588737024000000, 230584287177474048000000, 461168574354948096000000, 922337148709896192000000, 1844674297419792384000000, 3689348594839584768000000, 7378697189679169536000000, 14757394379358339072000000, 29514788758716678144000000, 59029577517433356288000000, 118059155034866712576000000, 236118310069733425152000000, 472236620139466850304000000, 944473240278933700608000000, 1888946480557867401216000000, 3777892961115734802432000000, 7555785922231469604864000000, 15111571844462939209728000000, 30223143688925878419456000000, 60446287377851756838912000000, 120892574755703513677824000000, 241785149511407027355648000000, 483570299022814054711296000000, 967140598045628109422592000000, 1934281196091256218845184000000, 3868562392182512437690368000000, 7737124784365024875380736000000, 15474249568730049750761472000000, 30948499137460099501522944000000, 61896998274920199003045888000000, 123793996549840398006091776000000, 247587993099680796012183552000000, 495175986199361592024367104000000, 990351972398723184048734208000000, 1980703944797446368097468416000000, 3961407889594892736194936832000000, 7922815779189785472389873664000000, 15845631558379570944779747328000000, 31691263116759141889559494656000000, 63382526233518283779118989312000000, 126765052467036567558237978624000000, 253530104934073135116475957248000000, 507060209868146270232951914496000000, 1014120419736292540465903828992000000, 2028240839472585080931807657984000000, 4056481678945170161863615315968000000, 8112963357890340323727230631936000000, 16225926715780680647454461263872000000, 32451853431561361294908922527744000000, 64903706863122722589817845055488000000, 129807413726245445179635690110976000000, 259614827452490890359271380221952000000, 519229654904981780718542760443904000000, 1038459309809963561437085520887808000000, 2076918619619927122874171041775616000000, 4153837239239854245748342083551232000000, 8307674478479708491496684167102464000000, 16615348956959416982993368334204928000000, 33230697913918833965986736668409856000000, 66461395827837667931973473336819712000000, 132922791655675335863946946673639424000000, 265845583311350671727893893347278848000000, 531691166622701343455787786694557696000000, 1063382333245402686911575573389115392000000, 2126764666490805373823151146778230784000000, 4253529332981610747646302293556461568000000, 8507058665963221495292604587112923136000000, 17014117331926442990585209174225846272000000, 34028234663852885981170418348451692544000000, 68056469327705771962340836696903385088000000, 136112938655411543924681673393806770176000000, 272225877310823087849363346787613540352000000, 544451754621646175698726693575227080704000000, 1088903509243292351397453387150454161408000000, 2177807018486584702794906774300908322816000000, 4355614036973169405589813548601816645632000000, 8711228073946338811179627097203633291264000000, 17422456147892677622359254194407266582528000000, 34844912295785355244718508388814533165056000000, 69689824591570710489437016777629066330112000000, 139379649183141420978874033555258132660224000000, 278759298366282841957748067110516265320448000000, 557518596732565683915496134221032530640896000000, 1115037193465131367830992268442065061281792000000, 2230074386930262735661984536884130122563584000000, 4460148773860525471323969073768260245127168000000, 8920297547721050942647938147536520490254336000000, 17840595095442101885295876295073040980508672000000, 35681190190884203770591752590146081961017344000000, 71362380381768407541183505180292163922034688000000, 142724760763536815082367010360584327844069376000000, 285449521527073630164734020721168655688138752000000, 570899043054147260329468041442337311376277504000000, 1141798086108294520658936082884674622752555008000000, 2283596172216589041317872165769349245505110016000000, 4567192344433178082635744331538698491010220032000000, 9134384688866356165271488663077396982020440064000000, 18268769377732712330542977326154793964040880128000000, 36537538755465424661085954652309587928081760256000000, 73075077510930849322171909304619175856163520512000000, 146150155021861698644343818609238351712327041024000000, 292300310043723397288687637218476703424654082048000000, 584600620087446794577375274436953406849308164096000000, 1169201240174893589154750548873906813698616328192000000, 2338402480349787178309501097747813627397232656384000000, 4676804960699574356619002195495627254794465312768000000, 9353609921399148713238004390991254509588930625536000000, 18707219842798297426476008781982509019177861251072000000, 37414439685596594852952017563965018038355722502144000000, 74828879371193189705904035127930036076711445004288000000, 149657758742386379411808070255860072153422890008576000000, 299315517484772758823616140511720144306845780017152000000, 598631034969545517647232281023440288613691560034304000000, 1197262069939091035294464562046880577227383120068608000000, 2394524139878182070588929124093761154454766240137216000000, 4789048279756364141177858248187522308909532480274432000000, 9578096559512728282355716496375044617819064960548864000000, 19156193119025456564711432992750089235638129921097728000000, 38312386238050913129422865985500178471276259842195456000000, 76624772476101826258845731971000356942552519684390912000000, 153249544952203652517691463942000713885105039368781824000000, 306499089904407305035382927884001427770210078737563648000000, 612998179808814610070765855768002855540420157475127296000000, 1225996359617629220141531711536005711080840314950254592000000, 2451992719235258440283063423072011422161680629900509184000000, 4903985438470516880566126846144022844323361259801018368000000, 9807970876941033761132253692288045688646722519602036736000000, 19615941753882067522264507384576091377293445039204073472000000, 39231883507764135044529014769152182754586890078408146944000000, 78463767015528270089058029538304365509173780156816293888000000, 156927534031056540178116059076608731018347560313632587776000000, 313855068062113080356232118153217462036695120627265175552000000, 627710136124226160712464236306434924073390241254530351104000000, 1255420272248452321424928472612869848146780482509060702208000000, 2510840544496904642849856945225739696293560965018121404416000000, 5021681088993809285699713890451479392587121930036242808832000000, 10043362177987618571399427780902958785174243860072485617664000000, 20086724355975237142798855561805917570348487720144971235328000000, 40173448711950474285597711123611835140696975440289942470656000000, 80346897423900948571195422247223670281393950880579884941312000000, 160693794847801897142390844494447340562787901761159769882624000000, 321387589695603794284781688988894681125575803522319539765248000000, 642775179391207588569563377977789362251151607044639079530496000000, 1285550358782415177139126755955578724502303214089278159060992000000, 2571100717564830354278253511911157449004606428178556318121984000000, 5142201435129660708556507023822314898009212856357112636243968000000, 10284402870259321417113014047644629796018425712714225272487936000000, 20568805740518642834226028095289259592036851425428450544975872000000, 41137611481037285668452056190578519184073702850856901089951744000000, 82275222962074571336904112381157038368147405701713802179903488000000, 164550445924149142673808224762314076736294811403427604359806976000000, 329100891848298285347616449524628153472589622806855208719613952000000, 658201783696596570695232899049256306945179245613710417439227904000000, 1316403567393193141390465798098512613890358491227420834878455808000000, 2632807134786386282780931596197025227780716982454841669756911616000000, 5265614269572772565561863192394050455561433964909683339513823232000000, 10531228539145545131123726384788100911122867929819366679027646464000000, 21062457078291090262247452769576201822245735859638733358055292928000000, 42124914156582180524494905539152403644491471719277466716110585856000000, 84249828313164361048989811078304807288982943438554933432221171712000000, 168499656626328722097979622156609614577965886877109866864442343424000000, 336999313252657444195959244313219229155931773754219733728884686848000000, 673998626505314888391918488626438458311863547508439467457769373696000000, 1347997253010629776783836977252876916623727095016878934915538747392000000, 2695994506021259553567673954505753833247454190033757869831077494784000000, 5391989012042519107135347909011507666494908380067515739662154989568000000, 10783978024085038214270695818023015332989816760135031479324309979136000000, 21567956048170076428541391636046030665979633520270062958648619958272000000, 43135912096340152857082783272092061331959267040540125917297239916544000000, 86271824192680305714165566544184122663918534081080251834594479833088000000, 172543648385360611428331133088368245327837068162160503669188959666176000000, 345087296770721222856662266176736490655674136324321007338377919332352000000, 690174593541442445713324532353472981311348272648642014676755838664704000000, 1380349187082884891426649064706945962622696545297284029353511677329408000000, 2760698374165769782853298129413891925245393090594568058707023354658816000000, 5521396748331539565706596258827783850490786181189136117414046709317632000000, 11042793496663079131413192517655567700981572362378272234828093418635264000000, 22085586993326158262826385035311135401963144724756544469656186837270528000000, 44171173986652316525652770070622270803926289449513088939312373674541056000000, 88342347973304633051305540141244541607852578899026177878624747349082112000000, 176684695946609266102611080282489083215705157798052355757249494698164224000000, 353369391893218532205222160564978166431410315596104711514498989396328448000000, 706738783786437064410444321129956332862820631192209423028997978792656896000000, 1413477567572874128820888642259912665725641262384418846057995957585313792000000, 2826955135145748257641777284519825331451282524768837692115991915170627584000000, 5653910270291496515283554569039650662902565049537675384231983830341255168000000, 11307820540582993030567109138079301325805130099075350768463967660682510336000000, 22615641081165986061134218276158602651610260198150701536927935321365020672000000, 45231282162331972122268436552317205303220520396301403073855870642730041344000000, 90462564324663944244536873104634410606441040792602806147711741285460082688000000, 180925128649327888489073746209268821212882081585205612295423482570920165376000000, 361850257298655776978147492418537642425764163170411224590846965141840330752000000, 723700514597311553956294984837075284851528326340822449181693930283680661504000000, 1447401029194623107912589969674150569703056652681644898363387860567361323008000000, 2894802058389246215825179939348301139406113305363289796726775721134722646016000000, 5789604116778492431650359878696602278812226610726579593453551442269445292032000000, 11579208233556984863300719757393204557624453221453159186907102884538890584064000000, 23158416467113969726601439514786409115248906442906318373814205769077781168128000000, 46316832934227939453202879029572818230497812885812636747628411538155562336256000000, 92633665868455878906405758059145636460995625771625273495256823076311124672512000000, 185267331736911757812811516118291272921991251543250546990513646152622249345024000000, 370534663473823515625623032236582545843982503086501093
```

plato. The second run had 36 Neurons and at eta of 0.2.

Training set score: 0.900475, Testing set score: 0.869687



e) Neurons 20, eta 0.1. Testing set CV score: -1.505540.

This is also a relatively small number of Neurons comparatively. Even though it seems small it looks to do a okay job looking at the graph below. The green line seems to represent the data to be able to use to predict it. It is just completely green in the center which means all the data values are accounted for.

Training set score: 0.935958

Testing set score: 0.933730

