

3 min to Spreed

```
In [61]: import numpy as np
import matplotlib.pyplot as plt
import os
import math
import scipy.io as spio
import pandas as pd
import itertools
from scipy.special import comb

plot_col = ['r', 'g', 'b', 'k', 'm']
plot_mark = ['o', '^', 'v', 'D', 'x', '+']

# Plots the rows in 'ymat' on the y-axis vs. 'xvec' on the x-axis
# with labels 'ylabels'
# and saves figure as pdf to 'dirname/filename'
def plotmatnsave(ymat, xvec, ylabels, dirname, filename):
    no_lines = len(ymat)
    fig = plt.figure(0)

    if len(ylabels) > 1:
        for i in range(no_lines):
            xs = np.array(xvec)
            ys = np.array(ymat[i])
            plt.plot(xs, ys, color = plot_col[i % len(plot_col)], lw=1,
label=ylabels[i])

        lgd = plt.legend(bbox_to_anchor=(0., 1.02, 1., .102), loc=3, ncol=3, mode="expand", borderaxespad=0.)

        savepath = os.path.join(dirname, filename)
        plt.xlabel('$x$', labelpad=10)
        plt.ylabel('$f(x)$', labelpad=10)
        plt.savefig(savepath, bbox_extra_artists=(lgd,), bbox_inches='tight')
    )
    plt.close()

# Sets the labels
labels = ['$e^x$', '1st order', '2nd order', '3rd order', '4th order']

# TODO: Given x values in "x_vec", save the respective function values ex,
# and its first to fourth degree Taylor approximations
# as rows in the matrix "y_mat"

#scipy.interpolate.approximate_taylor_polynomial(np.exp(x), 0, 0)
xvec = np.arange(-20, 8, 1)
y_mat = [[np.exp(x) for x in xvec],
          [1 + x for x in xvec],
          [1 + x + (x**2)/2 for x in xvec],
          [1 + x + (x**2)/2 + (x**3)/6 for x in xvec],
          [1 + x + (x**2)/2 + (x**3)/6 + (x**4)/24 for x in xvec]]
```

3 min to Spread

```
# Define filename, invoke plotmatnsave
filename = 'approx_plot_2.pdf'
plotmatnsave(y_mat, xvex, labels, '.', filename)
```

In [ ]:

```
In [135]: # There is numpy.linalg.lstsq, which you should use outside of this class
ss
def lstsq(A, b):
    return np.linalg.solve(A.T @ A, A.T @ b)

def fit(x, y, d):
    X = np.array([[x_t**i for i in range(d)] for x_t in x])
    return X, lstsq(X, y)

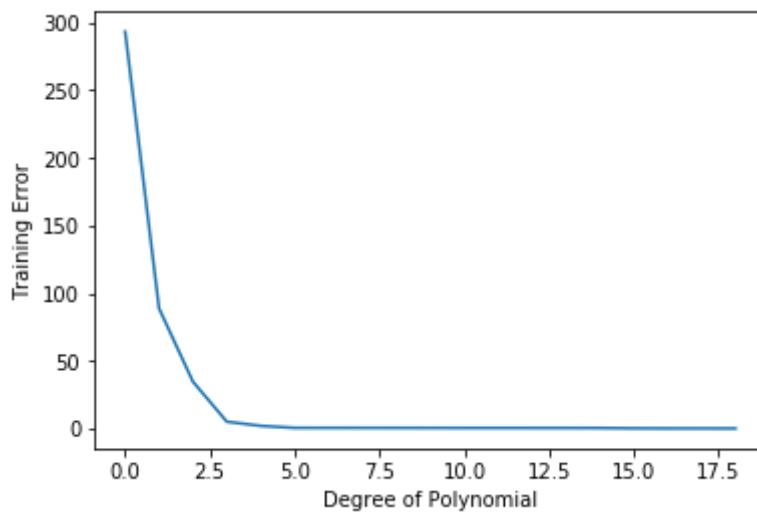
data = spio.loadmat('1D_poly.mat', squeeze_me=True)
x_train = np.array(data['x_train'])
y_train = np.array(data['y_train']).T

n = 20 # max degree

err = np.zeros(n - 1)

for i in range(n - 1):
    X, w = fit(x_train, y_train, i)
    err[i] = np.sum(np.square(y_train - X.dot(w))) / len(y_train)

plt.plot(err)
plt.xlabel('Degree of Polynomial')
plt.ylabel('Training Error')
plt.show()
```



## 5.c

3 min to Spreed

Average training error decreases as a function of D. Increasing the degree of the polynomial function used to fit the training data decreases the average error. Once the degree hits  $d-1$ , 1 - the number of data points used by our model, our model predicts the training values perfectly/exactly. We will not be able to fit a polynomial of degree n with the standard matrix inversion method as this will have a non-trivial null space.

## 5.d

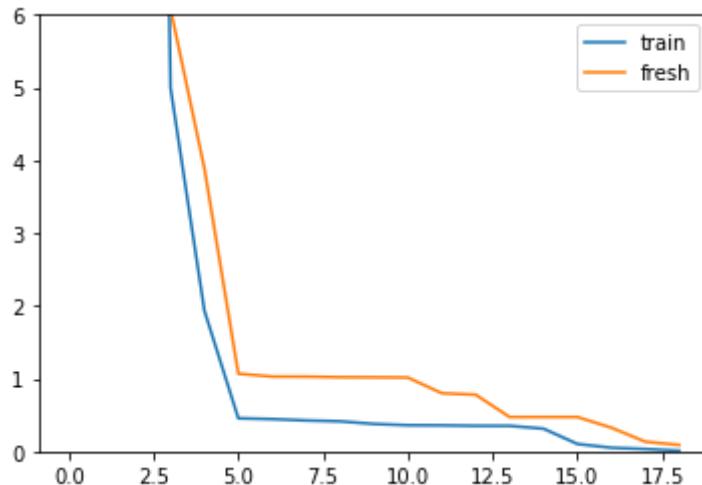
```
In [138]: data = spio.loadmat('1D_poly.mat', squeeze_me=True)
x_train = np.array(data['x_train'])
y_train = np.array(data['y_train']).T
y_fresh = np.array(data['y_fresh']).T

n = 20 # max degree
err_train = np.zeros(n - 1)
err_fresh = np.zeros(n - 1)

for i in range(n - 1):
    X, w = fit(x_train, y_train, i)
    err_train[i] = np.sum(np.square(y_train - X.dot(w))) / len(y_train)

for i in range(n - 1):
    X, w = fit(x_train, y_fresh, i)
    err_fresh[i] = np.sum(np.square(y_fresh - X.dot(w))) / len(y_fresh)

plt.figure()
plt.ylim([0, 6])
plt.plot(err_train, label='train')
plt.plot(err_fresh, label='fresh')
plt.legend()
plt.show()
```



3 min to Spreed

This plot shows how our model predictions for fresh are less accurate than they were for the training data. This is due to the fact that as the degree of the polynomial used to fit the data increases, the predictions become more complex and more tailored specifically to the training values. There is a tradeoff between reducing the error we get in the training data and the generalizability of the model to other data that is not part of the training set.

## 5.e

"Peach lady" should use the model which best minimizes errors for the fresh data. This occurs when a 5 degree polynomial is used to fit the data. Beyond this degree the model becomes overfitted to the training data. The complexity serves to hamper the effectiveness of our model for fresh data.

## 5.F 5.G

In [ ]:

3 min to Spreed

```
In [57]: data = spio.loadmat('polynomial_regression_samples.mat', squeeze_me=True)
)
data_x = data['x']
data_y = data['y']
Kc = 4 # 4-fold cross validation
KD = 6 # max D = 6

def gen(D,k):
    " Creates the powers of all multivariate conditions with k variables
    and max degree n.\n",
    " Code from karakfa from\n",
    " https://stackoverflow.com/questions/37711817/generate-all-possible
-outcomes-of-k-balls-in-n-bins-sum-of-multinomial-catego\n",
    " \"\"\"\n
    if(k==1):
        return [[D]]
    if(D==0):
        return [[0]*k]
    return [ g2 for x in range(D+1) for g2 in [ u+[D-x] for u in gen(x,k
-1) ] ]
}

def multivariate(data_x, D):
    if D == 0:
        return [1]
    X = []
    for i in range(len(data_x)):
        X.append(multivariate_row(data_x[i], D))
    return X

def multivariate_row(row, D):
    data_row = [1]
    for x in row:
        data_row.append(x)
    data_row = np.array(data_row)

    powers = gen(D, len(data_row))
    result = []
    for j in range(len(powers)):
        result.append(np.prod(data_row**powers[j]))
    return result

#multivariate_features = multivariate(data_x, 2)
D_multivariate_features = []
for k in range(KD):
    D_multivariate_features.append(multivariate(data_x, k))
```

3 min to Spreed

```
In [146]: data = spio.loadmat('polynomial_regression_samples.mat', squeeze_me=True)
)
data_x = data['x']
data_y = data['y']

def lstsqMULTIVARIATE(A, b, lambda_):
    return np.linalg.solve((A.T @ A) + lambda_*np.identity(len(A[0])), A.T @ b)

X_train_1 = D_multivariate_features[k][:250]
X_test_1 = D_multivariate_features[k][250:]
Y_train_1 = data_y[:250]
Y_test_1 = data_y[250:]
#lambda_ = 0.1
#w = lstsqMULTIVARIATE(np.vstack(X_train_1), Y_train_1, lambda_)
#np.array(X_train_1).dot(w)
#D_multivariate_features[2][:50]
```

3 min to Spreed

```
In [148]: data = spio.loadmat('polynomial_regression_samples.mat', squeeze_me=True)
)
data_x = data['x']
data_y = data['y']
Kc = 4 # 4-fold cross validation
KD = 6 # max D = 6
LAMBDA = [0, 0.05, 0.1, 0.15, 0.2]

def lstsqMULTIVARIATE(A, b, lambda_):
    return np.linalg.solve((A.T @ A) + lambda_*np.eye(len(A[0])), A.T @ b)

#def fit(D, lambda_):

#    err_train = np.zeros(Kc)
#    err_test = np.zeros(Kc)
#    for k in range(Kc):
#
#        X_train = D_multivariate_features[k][k*250:250*(k+1)]
#        X_test = np.append(D_multivariate_features[k][(k+1)*250:], (D_m
#ultivariate_features[k][:250*k]))
#        Y_train = data_y[k*250:250*(k+1)]
#        Y_test = np.append(data_y[(k+1)*250:], (D_multivariate_features
#[k][:250*k]))
#
#        w = lstsqMULTIVARIATE(np.vstack(X_train), Y_train, lambda_)
#        print('hi mark')
#        err_train[k] = np.sum(np.square((Y_train - np.array(X_train).dot
#(w)))) / len(Y_train)
#        err_test[k] = np.sum(np.square((Y_test - np.array(X_test).dot
#(w)))) / len(Y_test)
#        return np.sum(err_train) / len(err_train), np.sum(err_test) / len(e
#rr_test)

def fit(D, lambda_):
```

3 min to Spreed

```

err_train = np.zeros(Kc)
err_test = np.zeros(Kc)
for k in range(Kc):
    np.random.shuffle(D_multivariate_features[k])
    np.random.shuffle(data_y)

    X_train = D_multivariate_features[k][250:]
    X_test = D_multivariate_features[k][:250]

    Y_train = data_y[250:]
    Y_test = data_y[:250]

    w = lstsqMULTIVARIATE(np.vstack(X_train), Y_train, lambda_)
    err_train[k] = np.sum(np.square((Y_train - np.array(X_train).dot(w)))) / len(Y_train)
    err_test[k] = np.sum(np.square((Y_test - np.array(X_test).dot(w)))) / len(Y_test)
return np.sum(err_train) / len(err_train), np.sum(err_test) / len(err_test)

np.set_printoptions(precision=11)
Etrain = np.zeros((KD, len(LAMBDA)))
Evalid = np.zeros((KD, len(LAMBDA)))
for D in range(KD):
    print(D)
    for i in range(len(LAMBDA)):
        Etrain[D, i], Evalid[D, i] = fit(D + 1, LAMBDA[i])

print('Average train error:', Etrain, sep='\n')
print('Average valid error:', Evalid, sep='\n')

# YOUR CODE to find best D and i

```

3 min to Spreed

```
0
hi mark
1
hi mark
2
hi mark
```

3 min to Spreed

```
hi mark
hi mark
hi mark
hi mark
hi mark
hi mark
3
hi mark
4
hi mark
5
hi mark
```

3 min to Spreed

```

hi mark
Average train error:
[[ 0.07182671679  0.07188244367  0.07187541048  0.07184186813
  0.07184415236]
 [ 0.071753036   0.07183388202  0.0718987976   0.07190746852
  0.07193515026]
 [ 0.07180763744  0.07180836585  0.07183598188  0.07190196048
  0.07182514924]
 [ 0.07175686784  0.07186490528  0.07184289915  0.07189173132
  0.07184330527]
 [ 0.0717986333   0.07186792507  0.07181433832  0.07184588244
  0.07182584588]
 [ 0.07177699539  0.07185205664  0.07185248816  0.07187020781
  0.07185754449]]
Average valid error:
[[ 0.0592034313   0.06190881825  0.06555698559  0.08045415212
  0.07802876645]
 [ 0.08709529159  0.082431425    0.05608399486  0.05306716989
  0.04105100448]
 [ 0.06841046751  0.09179487428  0.08164899329  0.0546876341   0.086535
 72942]
 [ 0.08877346328  0.06895134396  0.07905601196  0.05876199631
 0.07874661822]
 [ 0.06894647952  0.06831009089  0.08806512541  0.07703037398
 0.08551068358]
 [ 0.07683447185  0.07544122359  0.07545801487  0.0674781477   0.073897
 53221]]

```

Best D for 5.f for lambda = 0.1 appears to be D = 4. This is the case as the average error for the test/valid is smallest here.

In [ ]:

In [ ]:

3 min to Spreed

). a.. Worked on HW with: Elimare Okoyoma,  
Prashanth Gareth, Daniel Mockaitis

b. I certify that all solutions are  
entirely in my own words & that I  
have not looked at another student's solutions.  
I have credited all extend sources in this  
write up.

Nich Lomo (26089160)

Q2. Settigate 2 for noise

2.a.

$$x \in \mathbb{R}^{n \times d}, y \in \mathbb{R}^{n \times 1}$$

$$\text{Min } \|y - Xw\|_2^2$$

$$\text{s.t. } \|w\|_2^2 \leq \beta^2$$

Lagrangian unconstrained form

$$\text{Min } \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

b. as we increase  $\beta$  our weight, used to penalize, is increased as well. With the lesser weight our lagrangian multiplier should be smaller.

- Big  $\beta \rightarrow$  lesser weights  $\rightarrow$  less  $\lambda$  req to penalize.

Conclusion:  $\beta$  inversely proportional to  $\lambda$

c.

$$w^T x = \hat{y} = w_1(x_1) + w_2 x_2 + \dots + w_n x_n \quad \text{for } Xw=y$$

adding error to feature vector  $x$

$$\hat{y} = w_1(x_1 + \epsilon_1) + \dots + w_n(x_n + \epsilon_n)$$

$$\hat{\hat{y}} - \hat{y} = w_1 \epsilon_1 + \dots + w_n \epsilon_n = w^T \epsilon$$

upper bound  
"how much could increase above  $\hat{y}$  given noise"

2.c.

$$\hat{w}_r = (X^T X + \lambda I)^{-1} X^T y \quad \text{as } \lim_{\lambda \rightarrow \infty} \hat{w}_r$$

$$\arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

- the weights estimate become super small.  $\rightarrow$  to zero.

- our weights are guaranteed to be @ most  $\frac{1}{2}$  as  $\lambda \rightarrow \infty$  our weights go to zero.

$$(y - Xw)^T (y - Xw) + \lambda w^T w$$

$$y^T y - y^T Xw - w^T X^T y + w^T X^T Xw + \lambda w^T w$$

$$y^T y - 2w^T X^T y + w^T X^T Xw + \lambda w^T w$$

Take  
 $\nabla($

$$)/\delta w$$

$$= 0 - 2y^T X + 2X^T Xw + 2\lambda w = 0$$

$$w(2X^T X + 2\lambda) = 2y^T X$$

$$w = (X^T X + \lambda I)^{-1} y^T X$$

2.e.

Part I.  $X^T X + \lambda I$

SVD of  $X^T X + \lambda I$

$$A = U \Sigma V^T$$

$$\text{let } X = U \Sigma V^T$$

$$= \sqrt{\begin{bmatrix} \sigma_1 & 0 \\ 0 & 0 \end{bmatrix}} V^T \begin{bmatrix} \sigma_1 & 0 \\ 0 & 0 \end{bmatrix} V^T + \lambda I$$

$$= \sqrt{\begin{bmatrix} \sigma_1^2 & 0 \\ 0 & 0 \end{bmatrix}} V^T + \lambda I$$

$\uparrow \downarrow \rightarrow V$  is orthogonal  
 $\therefore \{V^T = V^{-1}\}$

$$= \sqrt{\begin{bmatrix} \sigma_1^2 & 0 \\ 0 & 0 \end{bmatrix}} V^T + \sqrt{\lambda} I V^T \quad \{V^T V = I\}$$

$$= \sqrt{\left[ \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & 0 \end{bmatrix} V^T + \lambda I \right] V^T}$$

$$= \sqrt{\begin{bmatrix} \sigma_1^2 + \lambda & 0 \\ 0 & \lambda \end{bmatrix}} V^T \quad \begin{array}{l} \text{eigenvalues are along} \\ \text{the diagonal} \\ \sigma_1^2 + \lambda, \dots, \lambda \end{array}$$

thus the eigenvalues for  $X^T X + \lambda I$  are

$\sigma_1^2 + \lambda, \dots, \sigma_d^2 + \lambda, \lambda$  for regular eigenvalues but do not correspond to eigenvectors for the  $X^T X$  matrix.

Part II. regularizer term makes  $X^T X$  full rank.

Therefore  $X^T X + \lambda I$  now has good numerical inversion operations.

2. T.F

$X^T X = (X^T X)^{-1}$  will reduce sensitivity issues. For OLS fit this will result in a sensitive situation, where small changes in the features matter. Result in large / drastic changes to the weight estimate.

$\lambda_1 = 100$  is a better option for the eigenvalues of  $X$ . That is because the eigenvalues of  $X^T X + \lambda I$  will then be larger & thus  $(X^T X + \lambda I)^{-1}$  will have singular values which can be easily adapted to the resulting weight

$\sqrt{\lambda_i}$ . That is become true

Eigenvalues of  $X^T X + \sigma I$  will then be  
larger & thus  $(X^T X + \sigma I)^{-1}$  will have

Ex 489 Singular vectors which can be  
Note 2. at most  $\frac{1}{\sigma}$ , Not sensitive or  
highly amplified to the resulting weight  
estimate.

$$X^T X = \sqrt{[1000 \quad \sqrt{1000 + \sigma^2}]} \quad \text{some small eigenvectors}$$

$$\text{"along } \sigma I"$$

$$X^T X + \sigma I = \sqrt{[c \quad 1+\sigma \quad c \quad c]} \quad \sqrt{\sigma}$$

• If  $(X^T X)^{-1}$  is big then the estimate is  
very sensitive "little changes result in big changes"  
this occurs if  $X^T X$  has small eigenvalues. Only  
a large  $\sigma$  (prior w.r.t respect to the eigenvectors)  
helps overcome this sensitive nature of  $(X^T X + \sigma I)^{-1}$ .

2.h.

$$\lim_{\lambda \rightarrow 0}$$

•  $X^T X + \lambda I$  is full rank & invertible.

• The singular values become

$$\frac{1}{\sqrt{\lambda}} \neq \frac{1}{\lambda}, \text{ Many values are}$$

@ most  $\frac{1}{\lambda}$  ∵ as  $\lim_{\lambda \rightarrow 0}$  than our singular values approach zero. This effectively removes the effect of ridge regression by making our hypothesis approach zero.

•  $\lambda$  is no longer  $> 0$  ∵ we no longer have a unique solution for when the Feature Matrix & target samples make up an undetermined system.

$$\left[ \lim_{\lambda \rightarrow 0} \tilde{w}(x) = \lim_{\lambda \rightarrow 0} (X^T X)^{-1} X^T y + (\lambda I)^{-1} X^T y \right]$$

<sup>"From g"</sup>

$$\Rightarrow \text{approaches OLS} = (X^T X)^{-1} X^T y + \dots$$

2.I.  $\vec{w} = \arg \min_{\vec{w}} \frac{1}{2} \|y - \vec{x}\vec{w}\|_2^2 + \lambda \|\Gamma\vec{w}\|_2^2$   
 $= \min \left( \frac{1}{2} (\vec{y} - \vec{x}\vec{w})^T (\vec{y} - \vec{x}\vec{w}) + \lambda (\Gamma\vec{w})^T (\Gamma\vec{w}) \right)$   
 $\frac{1}{2} (\vec{y}^T - \vec{w}^T \vec{x}^T) (\vec{y} - \vec{x}\vec{w}) + \lambda \vec{w}^T \Gamma^T \Gamma \vec{w}$   
 $\frac{1}{2} (\vec{y}^T \vec{y} - \vec{y}^T \vec{x}\vec{w} - \vec{w}^T \vec{x}^T \vec{y} + \vec{w}^T \vec{x}^T \vec{x}\vec{w}) + \lambda \vec{w}^T \Gamma^T \Gamma \vec{w}$   
 are  
gradt  
resp/w  
0 =  $\frac{\partial}{\partial \vec{w}} 0 = \frac{1}{2} \cdot 0 - \frac{1}{2} \cdot 2 \vec{y}^T \vec{x} \cdot 1 + \frac{1}{2} \cdot 2 \cdot \vec{x}^T \vec{x} \vec{w} + \lambda \cdot 2 \Gamma^T \Gamma \vec{w} = 0$   
 $- \vec{y}^T \vec{x} + \vec{x}^T \vec{x}\vec{w} + 2\lambda \Gamma^T \Gamma \vec{w} = 0$  4  
 $F(\vec{w}) = \vec{w}^T \vec{x}^T \vec{y}$        $F(\vec{w} + \Delta) = (\vec{w} + \Delta)^T \vec{x}^T \vec{x} (\vec{w} + \Delta)$   
 $\quad \quad \quad (\vec{w}^T + \Delta^T) \vec{x}^T \vec{x} (\vec{w} + \Delta)$   
 $\frac{\delta F(\vec{w})}{\delta \Delta} = \frac{F(\vec{w} + \Delta) - F(\vec{w})}{\Delta}$   
 gradient  
derivative of  
transpose  
 $\therefore$  transpose of  
derivative is gradient  
 $((\vec{w}^T \vec{x}^T \vec{x}) + \Delta^T \vec{x}^T \vec{x})(\vec{w} + \Delta) - \vec{w}^T \vec{x}^T \vec{x} \vec{w}$   
 $= F(\vec{w}) + \vec{w}^T (\vec{x}^T \vec{x} + \vec{x} \vec{x}^T) \Delta + \vec{x}^T \vec{x} \Delta$   
 $F(\vec{w}) + ((\vec{x}^T \vec{x} + \vec{x} \vec{x}^T) \vec{w})^T \Delta + 0$   
 $\nabla F(\vec{w}) = 2\vec{x}^T \vec{x}$

4  $(\vec{x}^T \vec{x} + 2\lambda \Gamma^T \Gamma) \vec{w} = \vec{y}^T \vec{x}$

$\hat{\vec{w}} = (\vec{x}^T \vec{x} + 2\lambda \Gamma^T \Gamma)^{-1} \vec{y}^T \vec{x}$

closed form general solution.

(CONT.)

2g. Show  $(X^T X)^{-1} X^T y = \vec{w}$  has  $\infty = \# \text{ of solutions.}$

If a system is underdetermined with  $\lambda < 0$  then the matrix  $X$  (where  $X$  is the Features Matrix  $\in \mathbb{R}^{4 \times 5}$ ) has a nontrivial Null Space. Thus it is never invertible.

$$X\vec{w} = \vec{y}$$

$$\vec{w} = \vec{w}_0 + X^T \vec{x} \quad \vec{w}_0 \in \text{Nullspace}(X) \quad \vec{x} \in \mathbb{R}^5$$

- the vectors spanning the Null Space can have any linear combination of the basis vectors for the Null Space added to them & still get a solution

$$\vec{w} = \vec{w}_0 + \cancel{X^T \vec{x}}$$

any linear combination of basis of null space gives us a solution as well  $\therefore$  we have  $\infty$  solutions for  $X\vec{w} = \vec{y}$

$$\vec{w} = (X^T X)^{-1} X^T \vec{y}$$

① theorem from Math 54

2g  $(X^T X + \lambda I)$   $\Rightarrow$   $X^T y = b$  has unique solutions for  $\lambda > 0$

• PSD, Positive semi definite

- ① If  $A$  is invertible  $Ax=b$  has unique soln
- ② A square matrix  $A$  is invertible iff determinant  $\neq 0$

① From part 2e we determine that  $(X^T X + \lambda I)$  is invertible with singular values of at least  $\frac{1}{\sqrt{\lambda}}$ . In particular from part c:

$$\begin{aligned}(X^T X + \lambda I)^{-1} &= \left( \sqrt{\begin{bmatrix} \lambda + \sigma_r^2 & 0 \\ 0 & \lambda \end{bmatrix}} V^T \right)^{-1} \\ &= (\sqrt{\lambda})^{-1} \left[ (\lambda + \sigma_r^2)^{-1} \quad 0 \right] V^{-1} \\ &\quad \left[ \begin{array}{cc} 0 & (\lambda + \sigma_r^2)^{-1} \end{array} \right]\end{aligned}$$

Singular values of  $\lambda$

last  $\frac{1}{\sqrt{\lambda}}$  means that for  $\lambda > 0$  none of our singular values for  $(X^T X + \lambda I)^{-1}$  are zero.

②  $(X^T X + \lambda I)$  is a square matrix, due to its precisely stated. None of its eigenvalues or singular values are zero  $\therefore$  the det, which equals the product of the eigenvalues  $\neq 0$ . Thus  $(X^T X + \lambda I)$  is invertible.

" $A$  is invertible then  $Ax=b$  has unique soln"

$$Ax=b \quad A^{-1}A \mathbf{x} = A^{-1}b \quad \mathbf{x} = A^{-1}b \quad \text{unique solution}$$

Do taught

3 h for New Party

3. F.

Stars & Bars       $\ell$ -Dimension: bins      we have  $\ell+1$  bins  
Bars       $d$ -degree: balls

$$\ell=2 \quad D=3 \quad \binom{5}{2}$$

degree 0

$$1, x_1, x_{12}, x_{12}^2, x_{12}^3, x_1 x_2, x_1^2 x_2, x_1 x_2^2, x_1^3, x_2^3$$

$x_{12}, x_1$

$$\ell=2, D=2 \quad \binom{4}{2}$$

$$1, x_{11}, x_{12}, x_1 x_{12}, x_{12} x_{11}, x_{11}^2, x_{12}^2$$

Want: # of ways to distribute coeff. [to add up  
 $\ell$  variables, that sum to  $D$ ]  $\binom{\ell+D-1}{\ell-1}$   
(with exponents)

$\ell=2, D=2$  categories [for ex.]

$$x_{11}^0 x_{12}^0, x_{12}^2 | x_{11} x_{12}^0, x_{12}^1 | x_{11}^0 x_{12}^1 x_{12}^1 | x_{11}^2 x_{12}^0 x_{12}^0 | \dots$$

we need to count for the terms where the sum of  
the exponents is  $D$ ,  $\therefore$  we add an extra dummy  
variable to fulfill this in each of our categories  
 $\therefore$  the updated problem is

[# ways to add up  $\ell+1$  variables with exponents  
that sum to degree  $D$ .] =  $\binom{\ell+D}{\ell}$

3.g. A sample points  $\{x_i\}_{i=1}^n$   $x_{i,1} = x_{i,2} = \dots = x_{i,l} = \alpha_i$

Want  $F_l \in \mathbb{R}^{n \times (D+l)}$

distinct scalars

$\alpha \in \{1, 2, \dots, n\}$

Show  $F_l$  linearly dependent columns  $n \geq 2$ . Compare w/ e.

$$F = \begin{bmatrix} P_D(x_1) & \dots & P_D(x_n) \end{bmatrix}^T = \begin{bmatrix} P_D(x_1) \\ \vdots \\ P_D(x_n) \end{bmatrix} = \begin{array}{l} \text{(lets look at because} \\ \text{explore in 2.F)} \\ \text{w/ } l=2, D=2 \text{ & } n \geq 2 \end{array}$$

$$= \begin{bmatrix} 1 & \alpha_1 & \alpha_1 & \alpha_1^2 & \alpha_1^2 & \alpha_1^2 \\ 1 & \alpha_2 & \alpha_2 & \alpha_2^2 & \alpha_2^2 & \alpha_2^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \alpha_n & \alpha_n & \alpha_n^2 & \alpha_n^2 & \alpha_n^2 \end{bmatrix} \quad \begin{bmatrix} \alpha_1^3 & \alpha_1^3 & \alpha_1^3 & \alpha_1^3 \\ \alpha_2^3 & \alpha_2^3 & \alpha_2^3 & \alpha_2^3 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$\wedge = 2 \text{ for easier ex.}$

• 3 k+1 identical columns for degree k terms in  $P_D(x_i)$

∴ 3 linearly dependent columns for a Feature Matrix  
w/  $n \geq 2$  &  $l \geq 2$ .

Compare w/ • "the det is zero iff column vectors are linearly dependent" ∵ we know that  $F$  has a determinant of zero if we know that  $F$  is not full rank as at least two of four input points are equal (stated in part e), causing at least one of our eigenvalues for the Feature matrix to be zero.

3.b. Set of Sample point  $x_i$  s.t.  $F_L$  is full column  
 (informal) Rank (all column vectors are independent)  
 (Idea) Sample method: obtain sample points from sample  
 until have  $\lambda$  distinct points.  
 (concept) By (e) we know  $F$  full rank unless two  
 input data points are equal.  
 $\{x_i : i \in \{1, \dots, \lambda\}\}$  s.t.  $x_i \neq x_j \quad j=1, \dots, \lambda$   
 want our method to be functional for arbitrary  $D, l$

Multivariate case:  $F_L$  full column Rank: independent columns /  $l$ -dimension

let  $x_{i,1}, x_{i,2}, \dots, x_{i,l}$  be distinct scalars.

• for each  $j=1, \dots, l$  w/in each sample  
 point  $i$  select a distinct scalar (  
 as opposed to having one distinct scalar  $x_i$   
 for each).

let each

$$x_{i,j} = \alpha_i^{(j-1)}$$

$$x_{i,1} = \alpha_i^0$$

$$x_{i,2} = \alpha_i^1$$

$$\vdots$$

$$x_{i,l} = \alpha_i^{l-1}$$

4.a. i)  $e^x$

@  $x_0=0$

$$1^{\text{st}}: f(x) \approx 1 + x$$

$$2^{\text{nd}}: f(x) \approx 1 + x + \frac{x^2}{2!}$$

$$3^{\text{rd}}: f(x) \approx 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!}$$

$$4^{\text{th}}: f(x) \approx 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!}$$

ii)  $\sin(x)$

$$1^{\text{st}}: 0 + x$$

$$2^{\text{nd}}: 0 + x - \frac{x^2}{2!}$$

$$3^{\text{rd}}: 0 + x - \frac{x^3}{3!}$$

$$4^{\text{th}}: 0 + x - \frac{x^3}{3!}$$

b.  $\|f - \phi_M\|_\infty = \text{For } M=2 \quad I := [-4, 3]$

$$\sup_{\substack{x \in I \\ M \text{ th derivative}}} |f(x) - \phi_M(x)| \leq T \cdot \frac{|x-x_0|^{M+2}}{(M+1)!}$$

$$|f^{(m)}(x)| \leq T \quad x \in (x_0-s, x_0+s)$$

gutte Taylor approximation

$$\|f - \phi_2\|_\infty := \sup_{x \in I} |f(x) - \phi_2(x)| \leq \frac{T \cdot |x-x_0|^3}{3!}$$

$$\approx \sqrt{2\pi \lambda} \left(\frac{\lambda}{e}\right)^e$$

$$\geq \frac{|f''(x)| \cdot |x-x_0|^3}{3!} = \frac{|e^3| |x-0|^3}{3!}$$

largest value within  $[-4, 3]$ ,  
the interval, for the second  
derivative of the function  
from def  $\Delta T \geq |f''(x)|$

$$= \frac{|e^3| |x|^3}{3!}$$

4.b Compute  $\|f - g_M\|_\infty$  for arbitrary non-zero integers  $M$ .

$$I = [-4, 3] \quad T \geq |f^M(x)| \quad f^{(M)}(x) = f'(x) = e^x$$

• on domain

$(-4, 3]$ , largest

$M^{\text{th}}$  bound is

w/  $x=3$  for

$e^x$ .

$$\leq \frac{T|x-x_0|^{M+1}}{(M+1)!}$$

$$\leq \frac{C^3|x|^{M+2}}{(M+1)!}$$

$$N! \sim \sqrt{2\pi N} \cdot \left(\frac{N}{e}\right)^N \quad \left(\frac{M+1}{e}\right)^{M+1}$$

on interval  $I = [-4, 3]$

$$\leq \frac{C^3|x|^{M+2}}{\sqrt{2\pi(M+2)}} \cdot \left(\frac{M+1}{e}\right)^{M+2} = \frac{e^3|x|^{M+1}}{\sqrt{2\pi(M+1)}} \cdot \left(\frac{e}{M+1}\right)^{M+1}$$

By wolfram alpha.

$$\lim_{M \rightarrow \infty} \|f - g_M\|_\infty \leq \lim_{M \rightarrow \infty} \left( \frac{e^3|x|^{M+1}}{\sqrt{2\pi(M+1)}} \cdot \left(\frac{M+1}{e}\right)^{M+1} \right) = 0 \quad \text{logic}$$

• as the order of our Taylor approx approaches  $\infty$  the errors approach 0.

• Outside of  $I: [-4, 3]$ : the approximation error grows. From  $(-20, -4]$  the approx is diverse from the true function & from  $[3, 8]$  the true value diverges from the Taylor approximations.

$$I = [-4, 3], \quad \varepsilon \ll 1$$

$x$  or 3, constant?

4.C.

$$\text{1. Upper bound: } \frac{e^3 |x|^{(M+1)}}{(M+1)!} = \frac{e^3 |x|^{M+2}}{\sqrt{2\pi M + 2\pi}} \cdot \left(\frac{e}{M+1}\right)^{M+1}$$

on approx error.

$$\text{Show IF } D > O(\log(1/\varepsilon)) \Rightarrow |F(x) - f_D(x)| \leq \varepsilon$$

(Bounded by)

$$\text{For degree } D, \text{ approx error is: } \frac{e^3}{\sqrt{2\pi D + 2\pi}} \cdot \left(\frac{|x|e}{D+1}\right)^{D+2} \leq \varepsilon$$

$\sqrt{3 \cdot 3} = \sqrt{3 \cdot 3}$  • Constant factors not relevant so can drop:  $\sqrt{2\pi}, \frac{x}{(3)}$ ,

$$\frac{1}{\varepsilon} \leq \frac{\sqrt{2\pi D + 2\pi}}{e^3} \cdot \left(\frac{D+1}{|x|e}\right)^{D+2} \cdot \sqrt{2\pi(D+1)} =$$

$$\log(1/\varepsilon) \leq \log\left(\sqrt{D+2} \cdot e^{-3} \cdot \left(\frac{D+1}{|x|e}\right)^{D+2}\right)$$

$$\log(1/\varepsilon) \leq \log\left(\sqrt{D+2} - \underset{\text{constant}}{3 \cdot \log(e)} + (D+1) \log\left(\frac{D+1}{|x|e}\right)\right)$$

$$\leq \frac{1}{2} \log(D+1) + (D+1) \log(D+1) - (D+1) \log(|x|e) - (D+1) \log(|x|) - (D+1) \log(e)$$

$$\log(1/\varepsilon) \leq \frac{1}{2} \log(D+1) + (D+1) \log(D+1) - (D+1) \log(|x|) - (D+1)$$

$$\log(1/\varepsilon) \leq \frac{-(D+1)(1 + \log(|x|))}{(D+2)} - \underset{\text{constant}}{\log(D+1)}$$

$$\log(1/\varepsilon) \leq \underbrace{\frac{D+1}{2} \cdot \log(D+1)}_{\text{Argument}} - (D+1)$$

$$\leadsto \log(1/\varepsilon) \leq D \cdot \log(D) \quad \therefore \quad (?) \text{ Come Back}$$

4.C.

ii) Upper bound:  $\frac{1 \cdot |x|^{D+1}}{\sqrt{2\pi(D+2)}} \left(\frac{e}{D+1}\right)^{D+1}$

$$\varepsilon \geq (2\pi(D+1))^{-1/2} \left(\frac{|x|e}{D+1}\right)^{D+1}$$

$$\frac{1}{\varepsilon} \leq (2\pi(D+1))^{1/2} \left(\frac{D+1}{|x|e}\right)^{D+1}$$

$$\log(1/\varepsilon) \leq \frac{1}{2} \log(D+1) + (D+1) \log(D+1) - (D+1)(\log|x|e) - (D+1)$$

$$\log(1/\varepsilon) \leq \frac{D+1}{2} \log(D+1) - (D+1)$$

$$\Rightarrow \underbrace{\log(1/\varepsilon)}_{\text{arbitrarily small}} \leq D \log(D) \quad \textcircled{B}$$

Argument for part i) & ii)

$\Rightarrow$  given that  $D > \log(1/\varepsilon) \therefore D > \text{Big } *$  as  $\varepsilon \ll 1$  (Small Sufficiently)

thus we know that  $\log(D) > 2$  & we

$\therefore$  know  $D \cdot \log(D) > \log(1/\varepsilon)$ .

- we can now reverse the arguments shown above for ① & ② to show this implies that  $|F(x) - \phi_D(x)| \leq \varepsilon$ . (applicable for both  $\sin(x)$  &  $\cos(x)$ )

$$\text{u.d. } \exists \varepsilon > 0 \quad \exists \quad D \geq 1 \quad \text{s.t.} \quad \|F - P_D\|_{\infty} < \varepsilon$$

w.r.t respect to interval  $I$  (any closed  $I$ )

$$|F^M(x)| \leq T \quad M \geq 2 \quad x \in I \quad \begin{matrix} \text{"any function"} \\ \text{Univariate polynomial} \end{matrix}$$

$$g(x) \leq \frac{1}{(M+1)!} T (x-x_0)^{M+1} < \varepsilon$$

$$\lim_{M \rightarrow \infty} g(x) = T \cdot \lim_{M \rightarrow \infty} \frac{(x-x_0)^{M+1}}{(M+1)!} = 0$$

- as  $M$  approaches infinity (we get a large high degree polynomial) that we can use to approximate any  $F$  on a closed interval.

$$\bullet (M+1)! > (x-x_0)^{M+1} \quad \text{as} \quad M \rightarrow \infty \quad \therefore \lim_{M \rightarrow \infty} \left( \frac{1}{M+1} \right) = 0$$

4.2

I) Multivariate Taylor Series expansion in Matrix form

$$F(v) = F(v_0) \quad v = [x, y]^T \quad v_0 = [x_0, y_0]^T$$

$$F(v) = F(v_0) + \nabla F(v_0)^T \Delta v + \frac{1}{2} \Delta v^T \nabla^2 F(v_0) \Delta v + \text{Higher order terms}$$

$$\Delta v = \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix}$$

$$\nabla^2 F(v) = H(F(x)) = \text{Hess. of } F(x) = \begin{bmatrix} f_{xx} & f_{xy} \\ f_{yx} & f_{yy} \end{bmatrix}$$

$$\nabla F(v) = \begin{bmatrix} f_x \\ f_y \end{bmatrix}$$

II)  $F(v) = e^x y^2 \quad v = [x, y]^T \quad @ \quad v_0 \text{ makes } z^2 \text{ term zero}$

$$F(v) \approx e^{x_0} y_0^2 + \left[ e^{x_0} y_0^2 \right] \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} x - x_0 & y - y_0 \end{bmatrix} \begin{bmatrix} e^{x_0} y_0^2 & 2y_0 x_0 \\ 2y_0 e^{x_0} & 2e^{x_0} \end{bmatrix} \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix}$$

$$\nabla F(v) = \begin{bmatrix} e^x y^2 \\ 2y e^x \end{bmatrix} \quad \nabla^2 F(v) = \begin{bmatrix} e^x y^2 & 2y e^x \\ 2y e^x & 2e^x \end{bmatrix}$$

$$4.F. \quad F(v) = e^x y^2 \quad v_0 = 0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad v - v_0 = v \quad "directa" \\ \vec{v}(t) = \vec{0} + t(\vec{v} - \vec{v}_0)$$

$$g(t) = F(v(t)) \quad \text{Scalar to a scalar}$$

$$\vec{v}(t) = t \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} tx \\ ty \end{bmatrix} = [tx \ ty]^T$$

$$F(tx \ ty)$$

$$F(v(t)) = e^{tx} (yt)^2$$

$$\frac{1}{2}[tx \ ty] \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} tx \\ ty \end{bmatrix}$$

$$g(t) = \left\{ F(v(t)) = 0 + t \begin{bmatrix} 0 \\ 0 \end{bmatrix}^T \begin{bmatrix} x \\ y \end{bmatrix} + t^2 \frac{1}{2} \begin{bmatrix} x \ y \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \right\}$$

$$= t \cdot \frac{1}{2} \cdot 2y^2 = t \cdot y^2$$

$$(x+y) \begin{bmatrix} 0 \\ 2y \end{bmatrix} = \begin{bmatrix} 2y^2 \end{bmatrix}$$

$$g(t) = t^2 \cdot y^2$$

$$g(t) = g(t_0) + g'(t_0)(t-t_0) + \frac{1}{2}g''(t_0)(t-t_0)^2$$

$$g(t_0) = e^{\frac{tx}{(yt_0)^2}}$$

$$g'(t_0) = (yt_0)^2 \times e^{\frac{tx}{(yt_0)^2}} + 2y^2 t_0 e^{\frac{tx}{(yt_0)^2}}$$

$$g''(t_0) = 2t_0 y^2 x e^{\frac{tx}{(yt_0)^2}} + y^2 t_0^2 \times 2e^{\frac{tx}{(yt_0)^2}} + 2y^2 e^{\frac{tx}{(yt_0)^2}} + 2y^2 x t_0 e^{\frac{tx}{(yt_0)^2}}$$

$$g(t) =$$

$$\text{as } g(t_0) = F(v(t_0))$$

$$g'(t_0) = D F'(v(t_0))$$

$$g''(t_0) = D^2 F(v(t_0))$$

$$t_0 = 0 \quad \therefore$$

$$g(t) = g(0) + g'(0) \cdot t + \frac{1}{2} g''(0) \cdot t^2 \quad \therefore$$

5. John E has big ass peaches

a. A day peaches  $x_1 \dots x_n$  (Mars) predictors

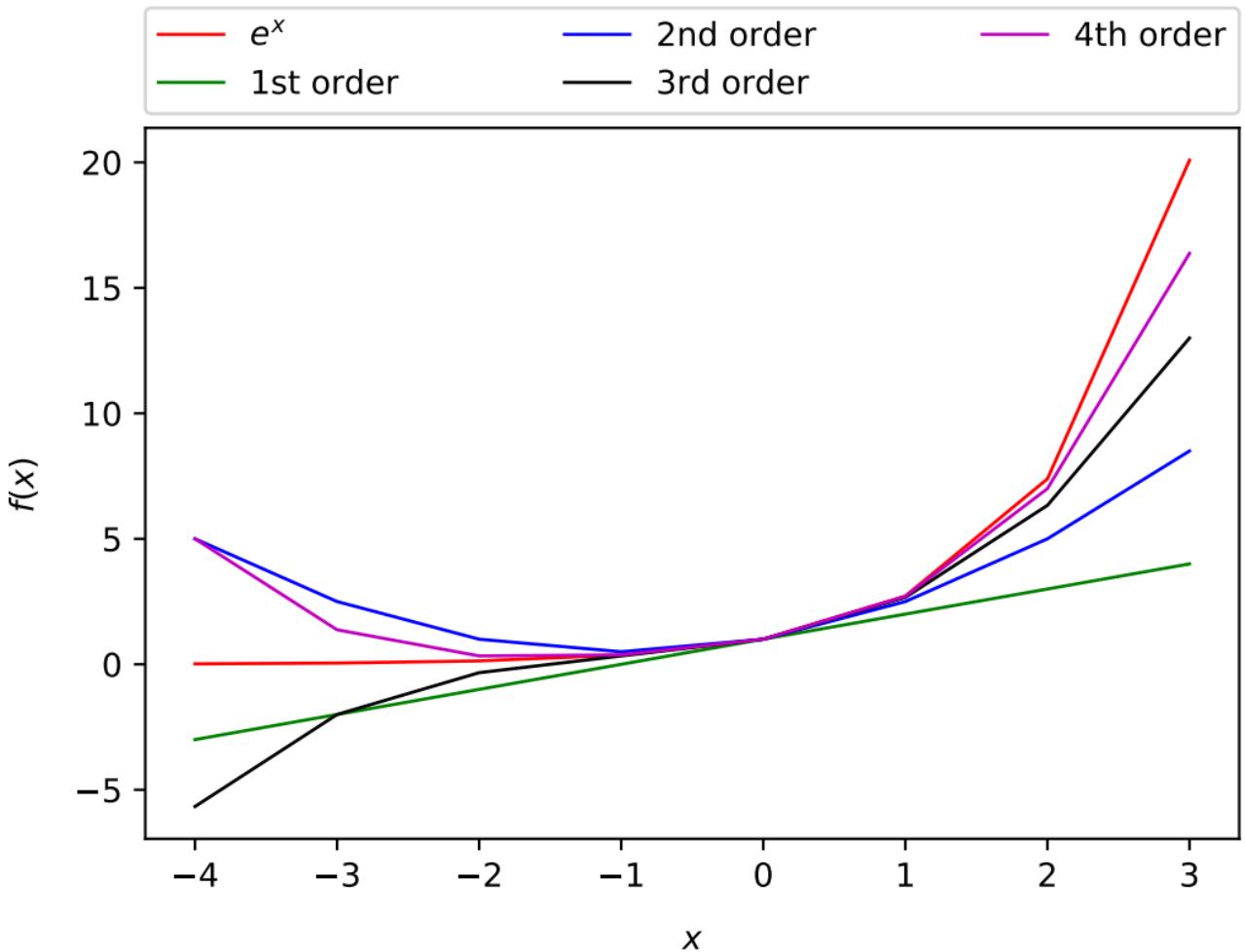
- $x_i$  for  $i=1, \dots, n$  are our predictors.
- $y_i$  represents our response
- we can use  $(x_i, y_i)$  as our day data  
 $\hat{y}_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_D x_i^D$

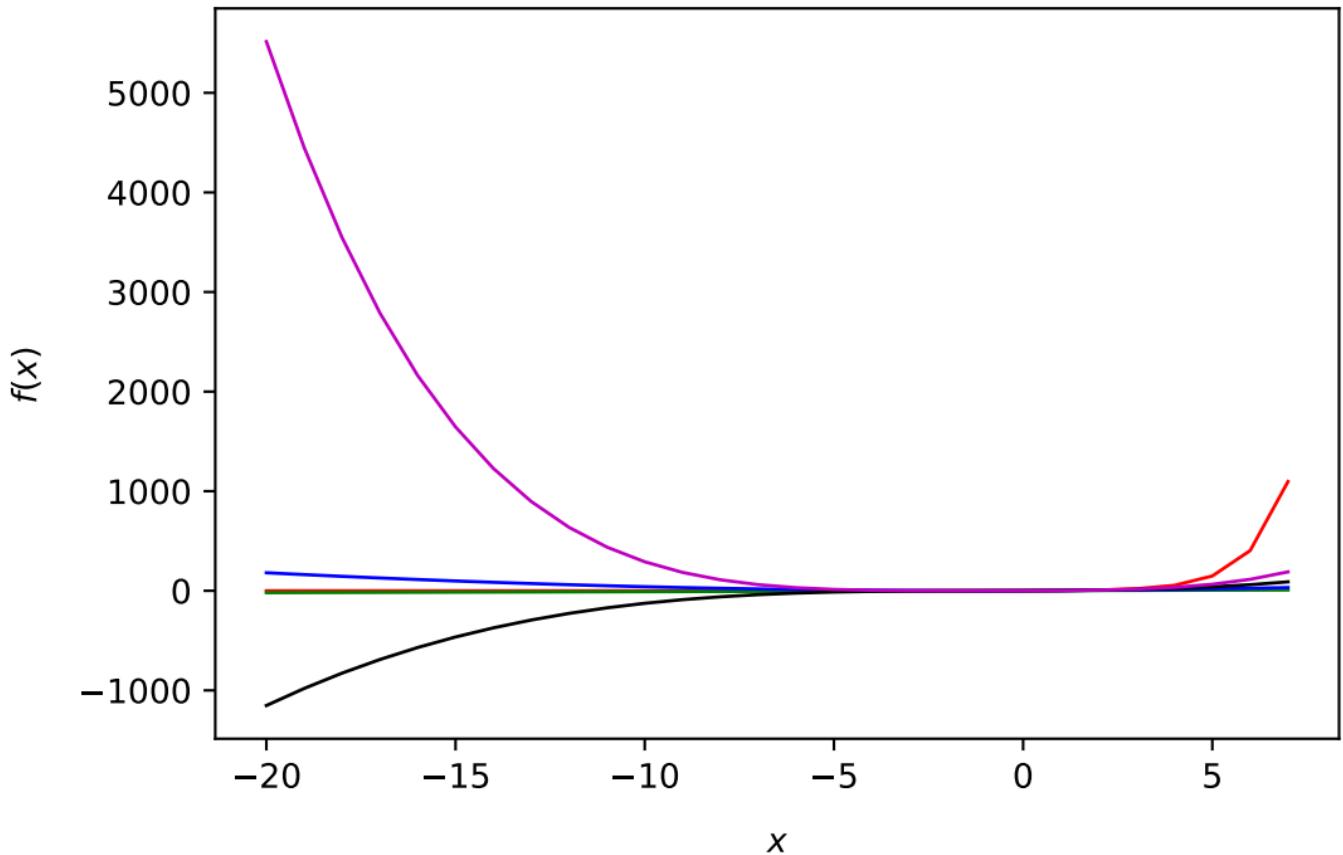
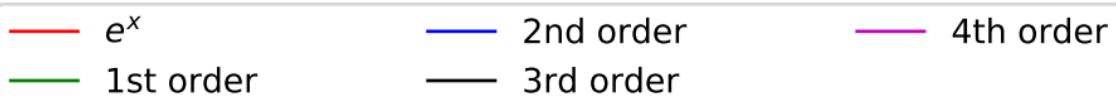
$$\text{min } \tilde{w} = \arg \min_w \|Xw - y\|_2^2$$

let  $X = \begin{bmatrix} 1 & x_1 & x_1 & \dots & x_1^D \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \end{bmatrix} \quad i=1, \dots, n$  be the feature matrix

$$X \begin{bmatrix} w_0 \\ \vdots \\ w_n \\ \tilde{w} \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \tilde{w} = (X^T X)^{-1} X^T y$$







## HW02 - CS 189

1.

Who else did you work with on this homework? In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

I worked on this homework with Ehimare Okoyomon, Prashanth Ganeth, and Daniel Mockaitis. We worked by getting together throughout the week and communicating on facebook.

*I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up}*

*Nicholas Lorio, 26089160*

### Question 8.

**What implication does the squared nature of the bias found through bias-variance decomposition have on its behavior (as N approaches infinity,  $f(x)$  in  $h(x, D)$ , etc). More specifically, if  $f$  is in the set of hypothesis functions, why does bias decrease as N increases?**

To be clear, bias will only decrease as the sample size increases **if** the samples are representative and without inherent biases themselves. Bias arises from omitted variable bias (non representative data), data collection bias, data cleaning bias, etcetera. High bias is related to under-fitting of the, it represents the ability of the model function to approximate the data. Therefore, a higher sample size N (with good/unbiased data points) will decrease the bias of the model. This however is only true given that our function estimate is within the true value function. E.G. that we're not trying to use a linear model to fit a quadratic relationship.

### Is hypothesis testing extrapolation or interpolation?

Hypothesis testing is a theoretical method of measuring the effectiveness of a hypothesis function  $h$  in relation to our real function  $f$ . We are comparing the data amongst known true values, as such it is more akin to interpolation than extrapolation. For real world data, the model is rarely known and the noise model may be entirely wrong. As such we cannot calculate bias and variance. We can however test algorithms over synthetic data.

### Source Note 5

<http://www.eecs189.org/static/notes/n5.pdf>

*“4. If  $f$  is in the set of hypothesis functions, bias will decrease with more data. If  $f$  is not in the set of hypothesis functions, then there is an underfitting problem and more data won’t help.”*

## Sources

[https://en.wikipedia.org/wiki/Taylor\\_series](https://en.wikipedia.org/wiki/Taylor_series)

<https://math.stackexchange.com/questions/499321/why-is-the-determinant-zero-iff-the-column-vectors-are-linearly-dependent>

[https://drive.google.com/file/d/0BzpkjIP\\_hA7\\_anFRU0FjQkIBTWs/view](https://drive.google.com/file/d/0BzpkjIP_hA7_anFRU0FjQkIBTWs/view)

<http://www.eecs189.org/static/notes/n3.pdf>

<http://www.eecs189.org/static/notes/n2.pdf>

<http://snasiriany.me/files/ml-book.pdf>

<https://math.stackexchange.com/questions/1985144/relationship-between-null-space-and-determinant>

[https://en.wikipedia.org/wiki/Taylor%27s\\_theorem#Taylor's\\_theorem\\_for\\_multivariate\\_functions](https://en.wikipedia.org/wiki/Taylor%27s_theorem#Taylor's_theorem_for_multivariate_functions)

<https://stackoverflow.com/questions/37711817/generate-all-possible-outcomes-of-k-balls-in-n-bins-sum-of-multinomial-catego>