

1. (a)

I did this homework with Luning Zhao.

We work on our homework separately, but we discuss when meeting problems.  
Homework is fine.

(b) 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.

2. (a) Rewrite  $\min_w \|y - Xw\|^2$  to unconstrained problem:  
 subject to  $\|w\|^2 \leq \beta^2$

$$L(w, \varepsilon) = \min_w (\|y - Xw\|^2 + \varepsilon (\|w\|^2 - \beta^2))$$

(b) Increasing  $\beta$  means  $\|w\|^2$  can go to larger value.

In the ridge regression context, this means  $\lambda$  can go smaller.

(c) change in the  $w^T x$  is  $w^T(x + \varepsilon) - w^T x = w^T \varepsilon \leq \|w\| \|\varepsilon\|$

$$\begin{aligned} (d) f(w) &= \|y - Xw\|^2 + \lambda \|w\|^2 \\ &= (y - Xw)^T (y - Xw) + \lambda w^T w \end{aligned}$$

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

$$\frac{\partial f(w)}{\partial w} = -y^T X - y^T X + w^T X^T X + w^T X^T X + \lambda w^T + \lambda w^T = 0$$

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

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

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

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

(e) If  $\sigma^2$  is eval for  $X^T X$ , then  $(X^T X) \vec{z} = \sigma^2 \vec{z}$ .

$$\text{For } (X^T X + \lambda I), \quad (X^T X + \lambda I) \vec{z} = X^T X \vec{z} + \lambda I \vec{z}$$

$$= (\sigma^2 + \lambda) \vec{z}.$$

so  $\sigma^2 + \lambda$  is eval for  $(X^T X + \lambda I)$ .

The evals for  $(X^T X + \lambda I)$  will be  $\sigma_1^2 + \lambda, \dots, \sigma_d^2 + \lambda$ .

If some eval for  $X^T X$   $\sigma_i^2 \ll 0$ , then it's numerically difficult to calculate  $(X^T X)^{-1}$ .  
 but  $\sigma_i^2 + \lambda$  will not be much smaller than 0, so it's much easier to calculate  $(X^T X + \lambda I)^{-1}$ .

(f) I think  $\lambda_2 = 0.5$  is a better choice.

First, with  $\lambda = 0.5$ , all eigenvalues for  $(X^T X + \lambda I)$  will not be closer to zero, so we can easily calculate  $(X^T X + \lambda I)^{-1}$ .

Second, with  $\lambda = 100$ ,  $\|w\|^2$  will be restricted to a small region, and  $(X^T X + \lambda I)^{-1}$  will be very small.

(g) because  $X \in \mathbb{R}^{4 \times 5}$ ,  $X$  must have nullspace  $W_0$  so that  $XW_0 = 0$ .

The reason is the 5 columns cannot be linear independent when they only have 4 numbers in each column.

So any vector  $\vec{w}(d)$  can be written as  $\vec{w}(d) = \vec{w}_0 + X^T d$ .

where  $\vec{w}_0$  is in  $X$  nullspace and  $X^T d$  is not. ( $\vec{w}_0^T X^T d = (X \vec{w}_0)^T d = 0$ )

$$\begin{aligned} \text{For OLS, } f(w) &= \|Xw - y\|^2 = (Xw - y)^T (Xw - y) \\ &= w^T X^T X w - w^T X^T y - y^T X w + y^T y \\ &= (w_0 + X^T d)^T X^T X (w_0 + X^T d) - (w_0 + X^T d)^T X^T y \\ &\quad - y^T X (w_0 + X^T d) + y^T y \\ &= \cancel{w_0^T X^T X w_0} + d^T X \cancel{X^T X w_0} + \cancel{w_0^T X^T X^T d} + d^T X X^T X X^T d \\ &\quad - \cancel{w_0^T X^T y} - d^T X X^T y - \cancel{y^T X w_0} - y^T X X^T d + y^T y \\ &= d^T X X^T X X^T d - d^T X X^T y - y^T X X^T d + y^T y. \end{aligned}$$

$f(w)$  is independent of  $w_0$ , so  $\frac{\partial f(w)}{\partial w_0} = 0 \Rightarrow w_0$  is arbitrary.

so OLS has infinite number of solutions.

For ridge regression,  $f(w) = \|Xw - y\|^2 + \lambda \|w\|^2$

$$\begin{aligned} &= d^T X X^T X X^T d - d^T X X^T y - y^T X X^T d + y^T y \\ &\quad + \lambda [w_0^T w_0 + d^T \cancel{X} w_0 + \cancel{w_0^T X^T d} + d^T X X^T d] \end{aligned}$$

$$= \alpha^T X X^T X X^T \alpha - \alpha^T X X^T Y - Y^T X X^T \alpha + Y^T Y \\ + \lambda w_0^T w_0 + \lambda \alpha^T X X^T \alpha.$$

$$\frac{\partial f(w)}{\partial w_0} = \lambda w_0^T + \lambda w_0^T = 2\lambda w_0^T = 0 \Rightarrow w_0 = 0.$$

$\Rightarrow w = X^T \alpha$  is the unique solution for ridge regression.

(h) If  $\lambda \rightarrow 0$ , ridge regression will be the same as OLS.

$$(i) f(w) = \frac{1}{2} \|y - Xw\|^2 + \lambda \|Tw\|^2$$

$$= \frac{1}{2} [(y - Xw)^T (y - Xw)] + \lambda [(\Gamma w)^T \Gamma w]$$

$$= \frac{1}{2} (Y^T Y - Y^T X w - W^T X^T Y + W^T X^T X w) + \lambda (W^T \Gamma^T \Gamma w)$$

$$\frac{\partial f(w)}{\partial w} = \frac{1}{2} (-Y^T X - Y^T X + W^T X^T X + W^T X^T X) + \lambda (W^T \Gamma^T \Gamma + W^T \Gamma^T \Gamma) \\ = -Y^T X + W^T X^T X + 2\lambda W^T \Gamma^T \Gamma = 0$$

$$W^T (X^T X + 2\lambda \Gamma^T \Gamma) = Y^T X$$

$$W^T = Y^T X (X^T X + 2\lambda \Gamma^T \Gamma)^{-1}$$

$$W = (X^T X + 2\lambda \Gamma^T \Gamma)^{-1} X^T Y$$

$$3. (a) F = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix} \quad x_1 \neq x_2 \Rightarrow \text{rank}(F) = 2$$

$$(b) F' = \begin{bmatrix} 1 & x_1 & \dots & x_1^D \\ 0 & x_2 - x_1 & \dots & x_2^D - x_1^D \\ \vdots & & & \\ 0 & x_n - x_1 & \dots & x_n^D - x_1^D \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & \dots & 0 \\ -1 & 1 & \dots & 0 \\ -1 & 0 & 1 & \dots & 0 \\ -1 & 0 & 0 & \dots & 1 \end{bmatrix}}_A F$$

$$\det(F') = \det(AF) = \det(A) \det(F)$$

$$\det(A) = 1 \cdot 1 - 0 = 1 \Rightarrow \det(F') = \det(F)$$

$$(c) \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^D \\ 0 & x_2 - x_1 & x_2^2 - x_1^2 & \dots & x_2^D - x_1^D \\ \vdots & & & & \\ 0 & x_n - x_1 & x_n^2 - x_1^2 & \dots & x_n^D - x_1^D \end{bmatrix} \xrightarrow{i) P_0} \begin{bmatrix} 1 & x_1 & \dots & x_1^{D-1} & 0 \\ 0 & x_2 - x_1 & & x_2^{D-1} - x_1^{D-1} & x_2^{D-1}(x_2 - x_1) \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & x_n - x_1 & & x_n^{D-1} - x_1^{D-1} & x_n^{D-1}(x_n - x_1) \end{bmatrix}$$

$$\xrightarrow{ii) P_1} \begin{bmatrix} 1 & x_1 & \dots & 0 & 0 \\ 0 & x_2 - x_1 & & x_2^{D-2}(x_2 - x_1) & x_2^{D-1}(x_2 - x_1) \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & x_n - x_1 & & x_n^{D-2}(x_n - x_1) & x_n^{D-1}(x_n - x_1) \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & x_2 - x_1 & x_2^2 - x_2 x_1 & & x_2^{D-1}(x_2 - x_1) \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & x_n - x_1 & x_n^2 - x_n x_1 & & x_n^{D-1}(x_n - x_1) \end{bmatrix}$$

$$P_1 = P_0 \underbrace{\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & & \ddots & \dots & 1 - x_1 \\ 0 & 0 & \dots & 1 & 1 \end{bmatrix}}_A \quad \det(A) = 1 \Rightarrow \det(F'') = \det(F')$$

$$(d) \text{ If } A\vec{x}_1 = \lambda_1 \vec{x}_1, \text{ then } \begin{bmatrix} 1 & 0 \\ 0 & A \end{bmatrix} \begin{bmatrix} 0 \\ \vec{x}_1 \end{bmatrix} = \begin{bmatrix} 0 \\ A\vec{x}_1 \end{bmatrix} = \begin{bmatrix} 0 \\ \lambda_1 \vec{x}_1 \end{bmatrix} = \lambda_1 \begin{bmatrix} 0 \\ \vec{x}_1 \end{bmatrix}$$

So  $\lambda_1, \dots, \lambda_d$  is evals for  $B$ .

$$\begin{bmatrix} 1 & 0^T \\ 0 & A \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \text{so } 1 \text{ is the other eval for } B.$$

$$\det(B) = 1 \det(A) - 0 = \det(A).$$

$$\begin{aligned} (e) \det(F') &= \det \begin{bmatrix} x_2 - x_1 & x_2(x_2 - x_1) & \dots & x_2^{D-1}(x_2 - x_1) \\ \vdots & & & \\ x_n - x_1 & x_n(x_n - x_1) & \dots & x_n^{D-1}(x_n - x_1) \end{bmatrix} \\ &= \prod_{i=2}^n (x_i - x_1) \det \begin{bmatrix} 1 & x_2 & \dots & x_2^{D-1} \\ \vdots & & & \\ 1 & x_n & \dots & x_n^{D-1} \end{bmatrix} \\ &= \prod_{i=2}^n (x_i - x_1) \cdot \det[P_{D-1}(x_2), P_{D-1}(x_3), \dots, P_{D-1}(x_n)]^T \\ &= \prod_{1 \leq i < j \leq n} (x_j - x_i) \neq 0 \Rightarrow \det(F') \neq 0 \Rightarrow \text{full rank.} \end{aligned}$$

(f) Using stars and bars argument, we first make  $D+l$  \*s.

Then we mark  $l$  bars, the number of \*s after each bar is the exponents for  $x_i$ .

For example,  $l=2$ ,  $D=3$ .

$$x \boxed{x} x \boxed{x} x \Rightarrow x_1 x_2$$

$$\boxed{x} x \boxed{x} x x \Rightarrow x_1 x_2^2$$

So size of  $P_D(x)$  is  $\binom{D+l}{l}$ .

$$(g) F_d = \left[ \begin{array}{c|c|c|c|c|c|c|c|c|c|c} 1 & \overline{d_1} & \overline{d_1} & \cdots & d_1 & d_1^2 & \cdots & d_1^D & d_1^D & \cdots & d_1^D \\ 1 & \overline{d_2} & \overline{d_2} & \cdots & d_2 & d_2^2 & \cdots & d_2^D & d_2^D & \cdots & d_2^D \\ \vdots & & & & & & & & & & \end{array} \right]$$

$\text{Col}_2 = \text{Col}_3 \Rightarrow F_d \text{ always has linearly dependent columns.}$

When  $D=1$ ,  $F_d$  is more likely linear independent compared with large  $D$ .

(h) To make  $F_d$  full rank, we should try to make  $x_i$  different from each other.

$$4. (a) i) e^x = e^0 + f'(x)x = 1 + e^0 \cdot x = 1+x$$

$$e^x = e^0 + \frac{f'(x)}{1} x + \frac{f''(x)}{2!} x^2 = 1 + x + \frac{1}{2} x^2$$

$$e^x = e^0 + \frac{f'(x)}{1} x + \frac{f''(x)}{2!} x^2 + \frac{f'''(x)}{3!} x^3 = 1 + x + \frac{1}{2} x^2 + \frac{1}{6} x^3$$

$$e^x = 1 + x + \frac{1}{2} x^2 + \frac{1}{6} x^3 + \frac{1}{24} x^4$$

$$ii) \sin x = \sin 0 + \cos 0 \cdot x = x$$

$$\sin x = x + (-\sin 0)/2! x^2 = x$$

$$\sin x = x - \frac{\cos 0}{3!} x^3 = x - \frac{1}{6} x^3$$

$$\sin x = x - \frac{1}{6} x^3 + \frac{\sin 0}{4!} x^4 = x - \frac{1}{6} x^3$$

(b) see plot in next page.

- When  $m=2$ ,  $\|f - \phi_m\|_\infty = \sup_{x \in I} |e^x - 1 - x - \frac{1}{2}x^2|$

From the plot we can see that  $\|f - \phi_m\|$  gets maximum when  $x=3$

$$\text{So } \|f - \phi_m\|_\infty = |e^3 - 1 - 3 - \frac{1}{2} \cdot 9| = |e^3 - 8.5| \approx 11.6$$

- For non-zero integer  $m$ ,

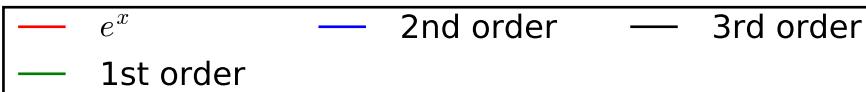
$$|f - \phi_m| = |e^x - \sum_{i=0}^m \frac{1}{i!} x^i| \leq \frac{T|x|^{m+1}}{(m+1)!}$$

where  $T = \max |f^{(m)}x| = \max |e^x| = e^3$ .

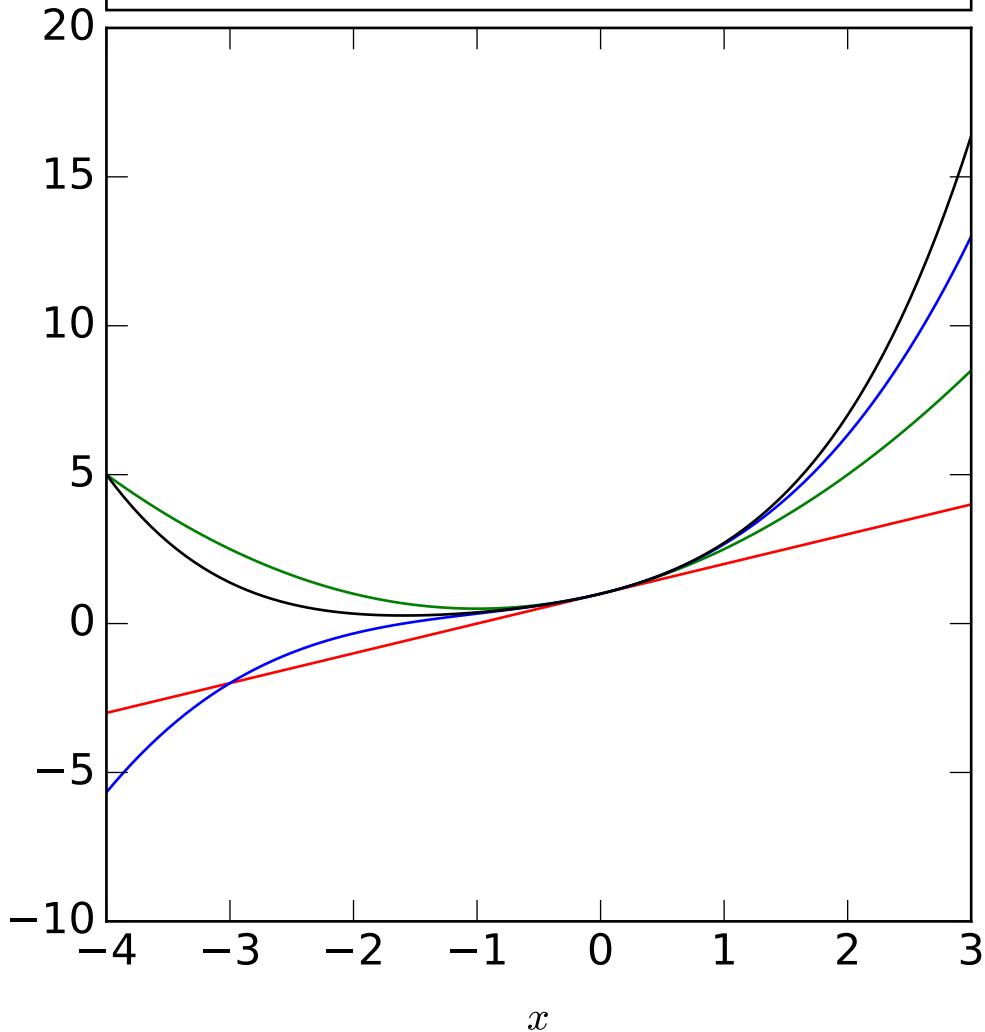
$$\text{So } |f - \phi_m| \leq \frac{e^3 \cdot 4^{m+1}}{(m+1)!}$$

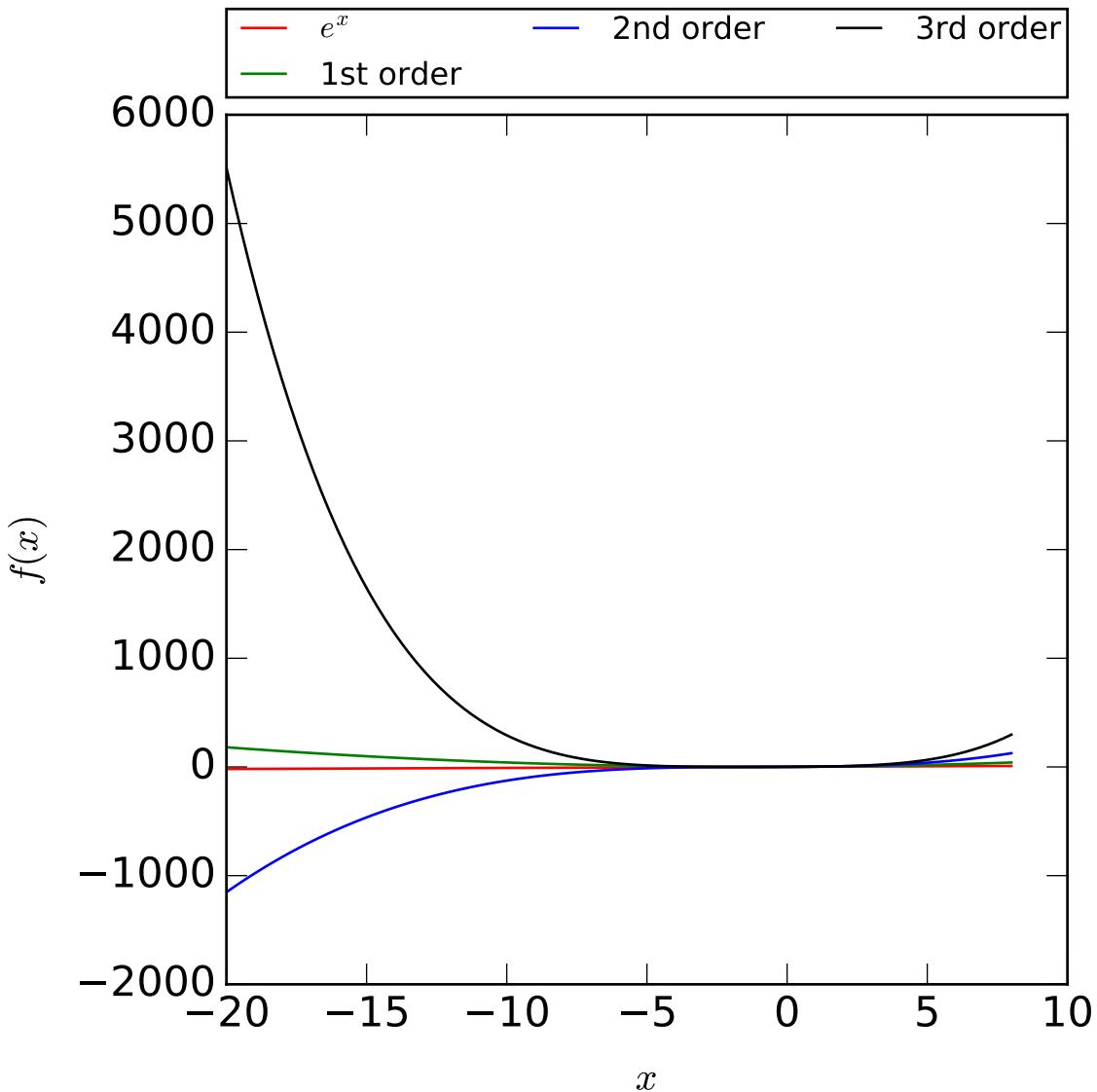
- $\lim_{m \rightarrow \infty} \frac{e^3 \cdot 4^{m+1}}{(m+1)!} = \lim_{m \rightarrow \infty} \frac{e^3 \cdot 4^{m+1}}{\sqrt{2\pi(m+1)} \left(\frac{m+1}{e}\right)^{m+1}} = \lim_{m \rightarrow \infty} \frac{e^3}{\sqrt{2\pi(m+1)}} \left(\frac{4e}{m+1}\right)^{m+1} = 0$

- approximation error becomes very large outside the bounded interval  $I$ .



$f(x)$





$$\begin{aligned}
 (c) \text{ For } f(x) = e^x, |f(x) - \phi_D(x)| &\leq \frac{e^3 \cdot 4^{D+1}}{(D+1)!} \leq \varepsilon \\
 \Rightarrow 3(D+1) \ln 4 - \ln \sqrt{2\pi(D+1)} \left(\frac{D+1}{e}\right)^{D+1} &\leq \ln \varepsilon \\
 \Rightarrow -3(D+1) \ln 4 + \ln \sqrt{2\pi(D+1)} + (D+1) \ln(D+1) - D+1 &\geq \ln \frac{1}{\varepsilon} \\
 \Rightarrow (D+1) \ln(D+1) - (\ln 4 + 1)(D+1) + \frac{1}{2} \ln(D+1) &\geq O(\ln \frac{1}{\varepsilon}) \\
 \Rightarrow D \ln D + \frac{1}{2} \ln D - 2D &\geq O(\ln \frac{1}{\varepsilon})
 \end{aligned}$$

when  $\varepsilon \ll 1$ ,  $\ln(\frac{1}{\varepsilon}) \gg 1$ ,

If  $D > O(\ln \frac{1}{\varepsilon})$ , then  $D \ln D \geq O(\ln \frac{1}{\varepsilon})$ .

$$\text{For } f(x) = \sin x, |f(x) - \phi_D(x)| \leq \frac{1 \cdot 4^{D+1}}{(D+1)!} \leq \varepsilon,$$

Similar as  $f(x) = e^x$ , we can prove if  $D > O(\log \frac{1}{\varepsilon})$ ,  $|f(x) - \phi_D(x)| \leq \varepsilon$ .

$$(d) \|f - \phi_D\|_{\infty} \leq \frac{T|x-x_0|^{D+1}}{(D+1)!} \leq \varepsilon.$$

$$\Rightarrow \frac{T|x-x_0|^{D+1} \cdot e^{D+1}}{\sqrt{2\pi(D+1)} (D+1)^{D+1}} = \frac{T}{\sqrt{2\pi}} \frac{(e|x-x_0|)^{D+1}}{(D+1)^{D+\frac{3}{2}}} \leq \varepsilon.$$

because  $T, x, x_0$  are bounded, when  $D \rightarrow \infty$ ,  $\frac{T}{\sqrt{2\pi}} \frac{(e|x-x_0|)^{D+1}}{(D+1)^{D+\frac{3}{2}}} \rightarrow 0$ .

so for any  $\varepsilon > 0$ , there is a  $D \geq 1$  s.t.  $\|f - \phi_D\|_{\infty} < \varepsilon$ .

$$(e) f(\vec{v}) = f(\vec{v}_0) + \left[ \frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \right] \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix} + \dots$$

$$+ \frac{1}{2!} [x - x_0 \quad y - y_0] \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix} \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix}$$

$$= f(\vec{v}_0) + \nabla f^T(\vec{v} - \vec{v}_0) + \frac{1}{2!} (\vec{v} - \vec{v}_0)^T \nabla^2 f (\vec{v} - \vec{v}_0)$$

$$\text{Here, } \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} \quad \nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$

For  $f(\vec{v}) = e^x y^2$ ,  $\frac{\partial f}{\partial x} = e^x y^2$ ,  $\frac{\partial f}{\partial y} = 2e^x y$

$$\frac{\partial^2 f}{\partial x^2} = e^x y^2 \quad \frac{\partial^2 f}{\partial x \partial y} = 2e^x y \quad \frac{\partial^2 f}{\partial y \partial x} = 2e^x y \quad \frac{\partial^2 f}{\partial y^2} = 2e^x$$

$$\Rightarrow f(\vec{v}) = e^{x_0} y_0^2 + [e^{x_0} y_0^2 \quad 2e^{x_0} y_0] \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix}$$

$$+ \frac{1}{2} [x - x_0 \quad y - y_0] \begin{bmatrix} e^{x_0} y_0^2 & 2e^{x_0} y_0 \\ 2e^{x_0} y_0 & 2e^{x_0} \end{bmatrix} \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix}$$

(f)  $\vec{v}(t) = t \vec{v} = t \begin{bmatrix} x \\ y \end{bmatrix}$

$$y(t) = f(\vec{v}(t)) = e^{tx_0} (\cancel{ty_0})^2 + [e^{tx_0} (\cancel{ty_0})^2 \quad 2e^{tx_0} \cancel{ty_0}] \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix}$$

$$+ \frac{1}{2} [x - x_0 \quad y - y_0] \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix}$$

$$= \frac{1}{2} [0 \quad y - y_0] \begin{bmatrix} x - x_0 \\ y - y_0 \end{bmatrix} = (y - y_0)^2 = t^2 y_0^2$$

5. (a) Jaina's problem can be written as:

$$\underbrace{\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^D \\ 1 & x_2 & x_2^2 & \dots & x_2^D \\ \vdots & & & & \\ 1 & x_n & x_n^2 & \dots & x_n^D \end{bmatrix}}_X \underbrace{\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}}_w = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_y.$$

$$\arg \min_w \|xw - y\|^2$$

(b) see output of b.py

(c) The average training error decreases with D, because with higher poly order, we have more model parameters, so we can fit data better.

If our model has  $n+1$  parameters,  $X \in R^{n \times n+1}$ , so we can't find the inverse of  $(X^T X)^{-1}$ . In the end, we are over fitting with  $n$  poly order.

(d) Test error is always larger than training error. [plot see output of c.py]  
Because w is minimized using training data, so bias will be minimized for training data.

(e) We should use the order when testing error becomes flat, which is poly order = 4.  
(plot from d)

(f) See code

(g) See code.

7.

Q: Prove Sherman-Morrison formula.

Suppose  $A \in R^{n \times n}$  is an invertible square matrix, and  $u, v \in R^n$  are column vectors. If  $A + uv^T$  is invertible, then

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u}$$

$$\begin{aligned} A: & (A + uv^T)(A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u}) \\ &= AA^{-1} - \frac{uv^TA^{-1}}{1 + v^TA^{-1}u} + uv^TA^{-1} - \frac{uv^TA^{-1}v^TA^{-1}}{1 + v^TA^{-1}u} \\ &= I + uv^TA^{-1} - \frac{u(1 + v^TA^{-1}u)v^TA^{-1}}{1 + v^TA^{-1}u} \\ &= I + uv^TA^{-1} - uv^TA^{-1} \\ &= I \end{aligned}$$