

Instructions

- The homework is due on Tuesday 4/25 at 5pm ET before the lecture starts.
- There are 15 points available for extra credit.
- No extension will be provided, unless for serious documented reasons.
- Start early!
- Study the material taught in class, and feel free to do so in small groups, but the solutions should be a product of your own work.
- This is not a multiple choice homework; reasoning, and mathematical proofs are required before giving your final answer.

1 SVD again [20 points]

1. (5 pts) Find the SVD of $A = [1, 1]$ without the use of computing devices/software.
2. (15 pts) Let $A \in \mathbb{R}^{m \times n}$ and let σ_1 be the maximum singular value of A . For $x \in \mathbb{R}^n \setminus \{0\}$ the spectral norm of A is defined as $\|A\|_2 = \max_x \frac{\|Ax\|_2}{\|x\|_2}$. Prove that

$$\|A\|_2 = \sigma_1.$$

2 Taylor polynomial approximation [10 points]

1. (5 pts) Let $f(x) = \sin(x) + \cos(x)$. Compute the degree 5 Taylor polynomial for f at $x = 0$.
2. (5 pts) Compute the quadratic approximation of the function $f(x, y) = x^2 + y^2 + 2xy - 3x + 2y + 5$ at the point $x = 5, y = 10$.

3 Derivatives [35 points]

Compute the derivative $\frac{df}{dx}$ for the following functions. It will be helpful to identify n, m where $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, and the dimensions of the derivative first.

(a) [5pts] $f(x) = \frac{1}{1+e^{-x}}, x \in \mathbb{R}$

(b) [5 pts] $f(x) = \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right), x \in \mathbb{R}$

(c) [5 pts] $f(x) = \sin(x_1) \cos(x_2), x \in \mathbb{R}^2$. $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

(d) [5 pts] $f(x) = xx^T, x \in \mathbb{R}^n$. $f: \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$

- (e) [5 pts] $f(x) = \sin(\log(x^T x)), x \in \mathbb{R}^n$.
- (f) [5 pts] $f(z) = \log(1 + z)$ where $z = x^T x, x \in \mathbb{R}^n$
- (g) [5 pts] $f(x) = x^T A x$ where $x \in \mathbb{R}^n, A \in \mathbb{R}^{n \times n}$.

4 Optimization [15 points]

1. (7.5 pts) Consider the univariate function $f(x) = x^3 + 6x^2 - 3x - 5$. Find its stationary points and indicate whether they are maximum, minimum or saddle points.
2. (7.5 pts) Explain how to solve the least squares loss in a linear model using (i) gradient descent and (ii) SVD. Discuss the pros and cons.

5 Coding [35 points]

Check the Jupyter notebook on our Git repo.

1.

1. (5 pts) Find the SVD of $A = [1, 1]$ without the use of computing devices/software.

$$A^T A = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\lambda = 2$$

$$\lambda = 0$$

$$\det(A^T A - \lambda I) = 0$$

$$\begin{bmatrix} 1-2 & 1 \\ 1 & 1-2 \end{bmatrix} \sim \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \sim \begin{bmatrix} 1 & -1 \\ 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \sim \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1-\lambda & 1 \\ 1 & 1-\lambda \end{bmatrix} = 0$$

$$\begin{bmatrix} 1 & -1 & | & 0 \\ 0 & 0 & | & 0 \end{bmatrix} \Rightarrow \begin{matrix} x_1 = x_2 \\ x_2 \text{ free} \end{matrix}$$

$$\begin{bmatrix} 1 & 1 & | & 0 \\ 0 & 0 & | & 0 \end{bmatrix} \Rightarrow \begin{matrix} x_1 + x_2 = 0 \\ x_2 \text{ free} \end{matrix}$$

$$(1-\lambda)^2 = 1$$

$$x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} x_2$$

$$x_1 = -x_2$$

$$1-\lambda = \pm 1$$

$$v_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}$$

$$x = \begin{bmatrix} -1 \\ 1 \end{bmatrix} x_2$$

$$\lambda = 0, 2$$

$$v_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}$$

$$V = \begin{bmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}, \Sigma = [\sqrt{2}, 0]$$

$$u_i = \frac{1}{\sigma} A v_i$$

$$u_1 = \frac{1}{\sqrt{2}} [1, 1] \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix} = 1$$

$$A = [1] \cdot [\sqrt{2}, 0] \cdot \begin{bmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix}$$

2. (15 pts) Let $A \in \mathbb{R}^{m \times n}$ and let σ_1 be the maximum singular value of A . For $x \in \mathbb{R}^n \setminus \{0\}$ the spectral norm of A is defined as $\|A\|_2 = \max_x \frac{\|Ax\|_2}{\|x\|_2}$. Prove that

$$\|A\|_2 = \sigma_1.$$

Proof:

$$i) \|Ax\| = \sqrt{(Ax)^T Ax} = \sqrt{x^T A^T A x} = \sqrt{x^T \lambda x}, \quad \sigma_{\max} = \sigma_1 = \sqrt{\lambda_{\max}}$$

$$ii) \max \frac{\|Ax\|}{\|x\|} = \sqrt{\frac{x^T \lambda_{\max} x}{x^T x}} = \sqrt{\frac{\lambda_{\max} x^T x}{x^T x}} = \sqrt{\lambda_{\max}} = \sigma_1$$



2.

1. (5 pts) Let $f(x) = \sin(x) + \cos(x)$. Compute the degree 5 Taylor polynomial for f at $x = 0$.

$$f(0) = \sin(0) + \cos(0) = 1$$

$$f'(0) = \cos(0) - \sin(0) = 1$$

$$f''(0) = -\sin(0) - \cos(0) = -1$$

$$f'''(0) = -\cos(0) + \sin(0) = -1$$

$$f^{(4)}(0) = \sin(0) + \cos(0) = 1$$

$$f^{(5)}(0) = 1$$

$$\sin' = \cos$$

$$\cos' = -\sin$$



Since $x=0$, we have the special case of Taylor polynomial = Maclaurin.

$$\begin{aligned} \Rightarrow P_5(x) &= f(0) + f'(0) \cdot x + \frac{f''(0)}{2!} x^2 + \frac{f'''(0)}{3!} x^3 + \frac{f^{(4)}(0)}{4!} x^4 + \frac{f^{(5)}(0)}{5!} x^5 \\ &= 1 + x - \frac{1}{2} x^2 - \frac{1}{6} x^3 + \frac{1}{24} x^4 + \frac{1}{120} x^5 \end{aligned}$$

2. (5 pts) Compute the quadratic approximation of the function $f(x, y) = x^2 + y^2 + 2xy - 3x + 2y + 5$ at the point $x = 5, y = 10$.

$$\begin{aligned} Q(x, y) &= f(x, y) + \frac{df}{dx} (x - x_0) + \frac{df}{dy} (y - y_0) \\ &\quad + \frac{1}{2} \frac{d^2f}{dx^2} (x - x_0)^2 + \frac{d^2f}{dxdy} (x - x_0)(y - y_0) + \frac{1}{2} \frac{d^2f}{dy^2} (y - y_0)^2 \end{aligned}$$

$$f(x, y) = x^2 + y^2 + 2xy - 3x + 2y + 5 = 25 + 100 + 100 - 15 + 20 + 5 = 235$$

$$\frac{df}{dx} = 2x + 2y - 3 = 10 + 20 - 3 = 27$$

$$\frac{df}{dy} = 2y + 2x + 2 = 20 + 10 + 2 = 32$$

$$\frac{d^2f}{dx^2} = 2 \quad \frac{d^2f}{dy^2} = 2 \quad \frac{d^2f}{dxdy} = 2$$

$$\begin{aligned} \Rightarrow Q(x, y) &= 235 + 27(x - 5) + 32(y - 10) + (x - 5)^2 \\ &\quad + 2(x - 5)(y - 10) + (y - 10)^2 \end{aligned}$$

3.

$$\begin{aligned} \text{a)} \quad \frac{df}{dx} &= \left(\frac{1}{1+e^{-x}} \right)' = \left[(1+e^{-x})^{-1} \right]' \\ &= -1 (1+e^{-x})^{-2} \cdot (-1) e^{-x} \\ &= (1+e^{-x})^{-2} \cdot e^{-x} \end{aligned}$$

$$\begin{aligned} \text{b)} \quad f(x) &= e^{-u} \quad u = \frac{1}{2\sigma^2}(x-\mu)^2 \\ \frac{df}{du} &= -e^{-u} \\ \frac{du}{dx} &= \frac{2}{2\sigma^2}(x-\mu) = \frac{1}{\sigma^2}(x-\mu) \\ \Rightarrow \frac{df}{dx} &= -e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \cdot \frac{1}{\sigma^2}(x-\mu) \end{aligned}$$

$$\text{c)} \quad f(x) = \sin(x_1) \cos(x_2), \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\frac{df}{dx} = \cos(x_1) \cos(x_2) - \sin(x_1) \sin(x_2)$$

$$\text{d)} \quad f(x) = x x^T \quad x \in \mathbb{R}^n$$

$$x x^T = \begin{bmatrix} x_1 x_1 & x_1 x_2 & \dots & x_1 x_n \\ x_2 x_1 & & & \\ \vdots & & \ddots & \\ x_n x_1 & & & x_n x_n \end{bmatrix}$$

$$\frac{df}{dx_1} \hookrightarrow \lim_{h \rightarrow 0} \frac{f(x_1+h, x_2, x_3, \dots) - f(x_1, x_2, \dots)}{h} = \lim \begin{bmatrix} (x_1+h)^2 - x_1^2 & x_2 \dots x_n \\ \vdots & 0 \\ x_n & \end{bmatrix}$$

Similarly,

$$\frac{df}{dx_2} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2+h, x_3, \dots) - f(x_1, x_2, x_3, \dots)}{h} = \lim \begin{bmatrix} 0 & -x_1 & \dots & 0 \\ x_1 & (x_2+h)^2 - x_2^2 & x_3 \dots x_n \\ \vdots & x_n & 0 \end{bmatrix}$$

Thus, each layer x_i in (x_1, x_2, \dots, x_n) has $\frac{df}{dx_i} = \begin{bmatrix} \dots & x_i & \dots \\ x_1 & 2x_i & \dots x_n \\ \vdots & \vdots & \end{bmatrix}$, other cells = 0.

$$\frac{df}{dx} = \begin{bmatrix} \frac{df}{dx_1} \\ \frac{df}{dx_2} \\ \vdots \\ \frac{df}{dx_n} \end{bmatrix}$$

e) $f(x) = \sin(\log(x^T x))$

$f: \mathbb{R}^n \rightarrow \mathbb{R}, \quad x^T x = 2x \in \mathbb{R}$

$\frac{df}{du} = (\sin u)' = \cos u$

$$\boxed{\frac{df}{dx} = \cos(\log(x^2)) \cdot \frac{2}{x}}$$

$$\begin{aligned} \frac{du}{dx} &= [\log(x^T x)]' = \frac{1}{x^T x} \cdot (x^T x)' \\ &= \frac{1}{x^T x} \cdot 2x = \frac{1}{x^2} \cdot 2x = \frac{2}{x} \end{aligned}$$

f) $f(z) = \log(1+z) \quad z = x^T x = x^2$

$$\frac{df}{dx} = \frac{df}{dz} \cdot \frac{dz}{dx} = \frac{1}{1+z} \cdot 2x$$

$$= \frac{1}{1+x^2} \cdot 2x$$

g) $f(x) = x^T A x \quad x \in \mathbb{R}^n \quad A \in \mathbb{R}^{n \times n}$

$x^T A x \in \mathbb{R}$

$$\begin{aligned} f(x) &= \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= \sum_{j=1}^n \sum_{i=1}^n a_{ij} x_i x_j \end{aligned}$$

$$\frac{df}{dx} = \sum_{k=1}^n \left(\underbrace{\sum_{j=1}^n a_{kj}}_{\text{row}} x_j + \sum_{i=1}^n \underbrace{a_{ik}}_{\text{col}} x_i \right) = x^T A^T + x \cdot A = 2Ax$$

4 Optimization [15 points]

- (7.5 pts) Consider the univariate function $f(x) = x^3 + 6x^2 - 3x - 5$. Find its stationary points and indicate whether they are maximum, minimum or saddle points.
- (7.5 pts) Explain how to solve the least squares loss in a linear model using (i) gradient descent and (ii) SVD. Discuss the pros and cons.

$$\begin{aligned}
 1) \quad f'(x) &= 3x^2 + 12x - 3 = 0 \\
 3(x^2 + 4x - 1) &= 0 \\
 x^2 + 4x + 4 &= 5 \\
 (x+2)^2 &= 5 \\
 x &= \pm\sqrt{5} - 2
 \end{aligned}$$

$$f''(x) = 6x + 12$$

$$f''(x_1) = 6\sqrt{5} - 12 + 12 = 6\sqrt{5} > 0$$

$$f''(x_2) = -6\sqrt{5} - 12 + 12 = -6\sqrt{5} < 0$$

$$\begin{aligned}
 x_1 &= \sqrt{5} - 2 \quad x_2 = -\sqrt{5} - 2 \\
 f(x) &\xrightarrow{x \rightarrow \infty} \infty \quad f(x) \xrightarrow{x \rightarrow -\infty} -\infty \\
 \therefore x_1, x_2 &\text{ are not global.}
 \end{aligned}$$

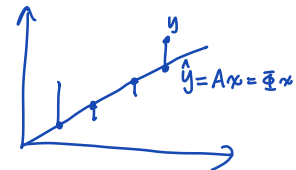
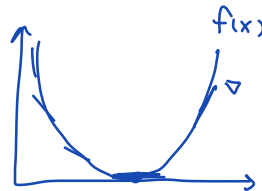
Since $f''(x_1) > 0$, it's a concave-up graph at x_1 . $\Rightarrow x_1$ is local minimum

Similarly, $f''(x_2) < 0$, so x_2 is local maximum.

2)

i) gradient.

$$\hat{x} = x - \eta (\nabla f(x))^T, \eta > 0$$



We'll choose a proper η , η is the gap between each selection of x . We'll continue calculating the gradient of x we selected until we get a $\nabla_x f(x^*) = 0$, which is a minimum.

ii) SVD is basically doing calculation.

$$\min \|Ax - b\|^2, A = U\Sigma V^T$$

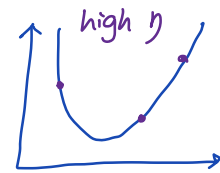
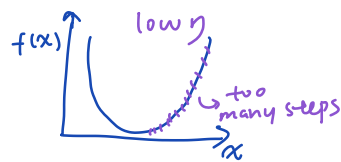
$$\Rightarrow A^T A \hat{x} = A^T b$$

$$V \Sigma^T U^T A \hat{x} = V \Sigma^T U^T b$$

$\hat{x} = V \Sigma^T U^T b$ is the optimal solution.

Pro & Con
Next Page

- | | Pro | con |
|-------------------|---|---|
| Gradient descent. | <ul style="list-style-type: none"> If the sample is huge, it's faster than SVD, since SVD needs to compute an inverse matrix which takes time. | <ul style="list-style-type: none"> It might get local minimum of x rather than global minimum. We need to choose proper η, if η is too high, it'll cause drastic changes and we might miss the minimum. If η is too low, it'll take too many steps to reach a minimum. |



- | | | |
|-----|--|--|
| SVD | <ul style="list-style-type: none"> If the sample is not too huge, SVD is faster than Gradient descent. It's more accurate than using Gradient descent. | <ul style="list-style-type: none"> When sample is too large, it'll take a long time computing the inverse matrix. |
|-----|--|--|