

① a) i) Let $A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, then the characteristic polynomial is
 $\det(A - \lambda I) = \det \begin{pmatrix} -\lambda & 1 \\ 1 & -\lambda \end{pmatrix} = \lambda^2 - 1$. Note $AA^T = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$.

Setting the characteristic polynomial to 0, we get $\lambda^2 - 1 = 0 \rightarrow \lambda = \pm 1$.
 Hence, the 2 eigenvalues are ± 1 .

Let \vec{v} and \vec{w} be the eigenvectors corresponding to $1, -1$ respectively.
 Then,

$$(A - (1)I_2) \vec{v} = 0$$

$$\left(\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} - \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right) \vec{v} = 0$$

$$\begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} \vec{v} = 0 \rightarrow \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \rightarrow \begin{aligned} -v_1 + v_2 &= 0 \\ \rightarrow v_2 &= v_1 \end{aligned}$$

$$\text{Then, } \vec{v} = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \begin{pmatrix} v_1 \\ v_1 \end{pmatrix} = v_1 \begin{pmatrix} 1 \\ 1 \end{pmatrix} = k \begin{pmatrix} 1 \\ 1 \end{pmatrix}, k \in \mathbb{R}.$$

Similarly, for $\lambda = -1$,

$$(A - (-1)I_2) \vec{w} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$w_1 + w_2 = 0 \rightarrow w_1 = -w_2$$

$$\vec{w} = \begin{pmatrix} w_1 \\ -w_1 \end{pmatrix} = w_1 \begin{pmatrix} 1 \\ -1 \end{pmatrix}.$$

$$\text{Hence, } \vec{v} = k_1 \begin{pmatrix} 1 \\ 1 \end{pmatrix}, k_1 \in \mathbb{R}, \text{ and } \vec{w} = k_2 \begin{pmatrix} 1 \\ -1 \end{pmatrix}, k_2 \in \mathbb{R}$$

We notice that each eigenvalue has norm 1, and the eigenvectors are orthogonal.

ii) Since $AB = I \rightarrow BA = I$, $AA^T = I \rightarrow A^T A = I$

Let \vec{v} be an eigenvector with eigenvalue λ such that $\lambda \vec{v} = A \vec{v}$. Then,

$$\begin{aligned} \vec{v}^T A^T A \vec{v} &= (A \vec{v})^T (A \vec{v}) = (\lambda \vec{v})^T (\lambda \vec{v}) = (\lambda \vec{v}) \cdot (\lambda \vec{v}) \\ &= \lambda^2 \vec{v} \cdot \vec{v} = \lambda^2 |\vec{v}|^2. \end{aligned}$$

$$\text{Also, } \vec{v}^T A^T A \vec{v} = \vec{v}^T I \vec{v} = \vec{v} \cdot \vec{v} = |\vec{v}|^2 \text{ since } A^T A = I.$$

$$\text{Hence, } \lambda^2 |\vec{v}|^2 = |\vec{v}|^2, \text{ or } \lambda^2 = 1 \rightarrow \lambda = \pm 1.$$

Hence all eigenvalues must be ± 1 , so all eigenvalues must have norm 1.

iii) From Piazza, we only need to focus on real eigenvalues, so we can ignore the adjoint case.

From ii), we know we can only have eigenvalues ± 1 .

Let \vec{x} be an eigenvector with eigenvalue λ_1 , and \vec{y} be an eigenvector with λ_2 . Then $A\vec{x} = \lambda_1\vec{x}$, $A\vec{y} = \lambda_2\vec{y}$, with $\lambda_1 \neq \lambda_2$.

$$\begin{aligned}\text{Then, } \vec{x} \cdot \vec{y} &= \vec{x}^T \vec{y} = \vec{x}^T I \vec{y} = \vec{x}^T A^T A \vec{y} = (A\vec{x})^T (A\vec{y}) \\ &= (\lambda_1 \vec{x})^T (\lambda_2 \vec{y}) = \lambda_1 \lambda_2 (\vec{x} \cdot \vec{y}).\end{aligned}$$

Hence either $\lambda_1 \lambda_2 = 1$ or $\vec{x} \cdot \vec{y} = 0$. Since $\lambda_1 \cdot \lambda_2 = 1 \cdot -1 = -1 \neq 1$, $\vec{x} \cdot \vec{y} = 0$, and hence \vec{x} and \vec{y} are orthogonal, as needed.

iv) Since the eigenvalues are ± 1 , the transformation does not scale the input in any way. Since length and angles must be preserved, $A\vec{x}$ will simply be \vec{x} after a rotation (around the origin), or a flip/reflection across some subspace.

If $\det(A) = 1$, we will not have a reflection, while if $\det(A) = -1$, we will (note we can still rotate before/after a rotation).

i) We utilize the Spectral Theorem, which tells us that if A is $n \times n$ and symmetric, $A = PDP^{-1} = PDP^T$, where P is columns are an orthonormal eigenbasis of \mathbb{R}^n , and D is a diagonal matrix with entries $\lambda_1, \dots, \lambda_n$.

Let $A = U \Sigma V^T$, where $A \in \mathbb{R}^{m \times n}$. Then,

$$AA^T = (U \Sigma V^T)(U \Sigma V^T)^T = U \Sigma V^T V \Sigma^T U^T.$$

Since V is orthogonal, $V^T V = I$, so

$AA^T = U \Sigma \Sigma^T U^T$. Since $(AA^T)^T = A^T A = AA^T$, AA^T is symmetric, so we can use Spectral decomposition to conclude the left singular vectors of A are the eigenvectors of AA^T .

Similarly, since $A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T U^T U \Sigma V^T$
 $= V \Sigma^T \Sigma V^T$,

the right singular vectors of A are the eigenvectors of $A^T A$.

ii) From i), we know $AA^T = U \Sigma \Sigma^T U^T$, so by the spectral theorem $\Sigma \Sigma^T$ is the eigenvalue matrix for AA^T .

If A is $m \times n$, then Σ is also $m \times n$, so $\Sigma \Sigma^T$ is $m \times m$. However, A only has $\min(m, n)$ singular values.

Hence, the first $\min(m, n)$ eigenvalues of AA^T correspond to the singular values of A , and any remaining eigenvalues (any remaining diagonal entries in $\Sigma \Sigma^T$) will be 0.

The eigenvalues of $A^T A$ are the same, except we are filling in $\Sigma^T \Sigma$ instead (so we might not need to pad with 0s depending on the dimension of $\Sigma^T \Sigma$).

- c) i) False
ii) False
iii) True
iv) True
v) True

② a) i) Using the probability chain rule, we get

$$P(H50 | \text{Tails}) = \frac{P(H50 \wedge \text{Tails})}{P(\text{Tails})} = \frac{0.5 \cdot 0.5}{P(\text{Tails} | H50) P(H50) + P(\text{Tails} | H60) P(H60)}$$
$$= \frac{0.25}{0.4 \cdot 0.5 + 0.5 \cdot 0.5} = \boxed{\frac{5}{9}}$$

ii) We want to find $P(H50 | THHH)$. Using Bayes' Thm, we can rewrite as

$$P(H50 | THHH) = \frac{P(THHH | H50) P(H50)}{P(THHH)}$$

Then using the law of probability we can rewrite $P(THHH)$ as

$$P(THHH) = P(THHH | H50) P(H50) + P(THHH | H60) P(H60)$$
$$= 0.5^4 \cdot 0.5 + 0.4 \cdot 0.6^3 \cdot 0.5$$

$$\text{Hence, } P(H50 | THHH) = \frac{0.5^5}{0.5^5 + 0.5 \cdot 0.4 \cdot 0.6^3} \approx \boxed{42\%}$$

iii) We first need to find $P(9 \text{ heads}) = P(H50) P(9 \text{ heads} | H50) + P(H55) P(9 \text{ heads} | H55) + P(H60) P(9 \text{ heads} | H60)$.

We can do this with the binomial distribution. We use

$$\binom{n}{k} p^k (1-p)^{n-k} \quad \text{with } n=10, k=9, \text{ and } p = 0.5, 0.55, \text{ and } 0.6.$$

$$\text{Then, } P(9 \text{ heads} | H50) = \binom{10}{9} 0.5^9 0.5^1 = 10 \cdot 0.5^{10} = \frac{5}{512}$$

$$P(9 \text{ heads} | H55) = \binom{10}{9} 0.55^9 0.45^1 \approx 0.020724, \text{ and}$$

$$P(9 \text{ heads} | H60) = \binom{10}{9} 0.6^9 0.4^1 \approx 0.04031$$

$$\text{Hence, } P(9 \text{ heads}) = 0.0236.$$

$$\text{Then, } P(H50 | 9 \text{ heads}) = \frac{P(9 \text{ heads} | H50) P(H50)}{P(9 \text{ heads})} = \frac{\frac{5}{512} \cdot \frac{1}{3}}{0.0236} \approx \boxed{14\%}$$

$$P(H55 | 9 \text{ heads}) = \frac{0.020724 \cdot \frac{1}{3}}{0.0236} \approx \boxed{29\%}$$

$$P(H60 | 9 \text{ heads}) = \frac{0.04031 \cdot \frac{1}{3}}{0.0236} \approx \boxed{57\%}$$

b) We once again use Bayes' Thm.

$$P(\text{pregnant} | \text{positive}) = \frac{P(\text{positive} | \text{pregnant}) P(\text{pregnant})}{P(\text{positive})}, \text{ and then}$$

$$P(\text{positive}) = P(\text{positive} | \text{pregnant}) P(\text{pregnant}) + P(\text{positive} | \text{not pregnant}) P(\text{not pregnant})$$
$$P(\text{positive}) = 0.99 (0.01) + 0.1 (0.99) = 0.1089$$

$$\text{Hence } P(\text{pregnant} | \text{positive}) = \frac{0.99 (0.01)}{0.1089} \approx \boxed{9.1\%}$$

The number is low due to the high rate of false positives; since only 1% of the population is pregnant, most positive tests will come from false positives instead of true positives.

c) Since expected value is a linear operator, we have

$$E(A\vec{x} + \vec{b}) = E(A\vec{x}) + E(\vec{b}).$$

Then since \vec{b} is deterministic (not random), $E(\vec{b}) = \vec{b}$.

Finally, since A is deterministic, we have

$$E(A\vec{x}) = E \left[\begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \right] = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} E \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$
$$= A E(\vec{x})$$

$$\text{Hence } \boxed{E(A\vec{x} + \vec{b}) = A E(\vec{x}) + \vec{b}}$$

$$\text{d) We get } \text{cov}(A\vec{x} + \vec{b}) = E \left((A\vec{x} + \vec{b} - E(A\vec{x} + \vec{b})) * (A\vec{x} + \vec{b} - E(A\vec{x} + \vec{b}))^T \right)$$

$$= E \left((A\vec{x} + \vec{b} - A E(\vec{x}) - \vec{b}) * (A\vec{x} + \vec{b} - A E(\vec{x}) - \vec{b})^T \right)$$

$$= E \left((A\vec{x} - A E(\vec{x})) (A\vec{x} - A E(\vec{x}))^T \right)$$

$$= E \left(A (\vec{x} - E(\vec{x})) (\vec{x} - E(\vec{x}))^T A^T \right)$$

$$= A E((\vec{x} - E(\vec{x})) (\vec{x} - E(\vec{x}))^T) A^T$$

$$\boxed{= A \text{cov}(\vec{x}) A^T}$$

we are able to cancel \vec{b} and factor A since both are deterministic.

③ a) $\nabla_{\vec{x}} \vec{x}^T A \vec{y} = \nabla_{\vec{x}} \vec{x} \cdot A \vec{y}$. Let $\vec{d} = A \vec{y}$, then

$$\nabla_{\vec{x}} (\vec{x} \cdot \vec{d}) = \begin{pmatrix} \partial/\partial x_1 (\vec{x} \cdot \vec{d}) \\ \partial/\partial x_2 (\vec{x} \cdot \vec{d}) \\ \vdots \\ \partial/\partial x_n (\vec{x} \cdot \vec{d}) \end{pmatrix} = \begin{pmatrix} d_1 \\ \vdots \\ d_n \end{pmatrix} = \vec{d} = \boxed{A \vec{y}}$$

b) $\nabla_{\vec{y}} \vec{x}^T A \vec{y} = \nabla_{\vec{y}} (A^T \vec{x})^T \vec{y}$. Let $A^T \vec{x} = \vec{d}$, then

$$\nabla_{\vec{y}} \vec{x}^T A \vec{y} = \nabla_{\vec{y}} \vec{d}^T \vec{y} = \nabla_{\vec{y}} (d_1 y_1 + \dots + d_m y_m) = \vec{d} = \boxed{A^T \vec{x}}$$

c) $\vec{x}^T A \vec{y} = \sum_i^n \sum_j^m a_{ij} x_i y_j$. Then,

$$\nabla_A (\vec{x}^T A \vec{y}) = \begin{bmatrix} \frac{\partial f}{\partial a_{11}} x_1 y_1 & & & \\ & \ddots & & \\ & & \frac{\partial f}{\partial a_{nm}} x_n y_m & \\ & & & \ddots \end{bmatrix} = \begin{bmatrix} x_1 y_1 & & & \\ & \ddots & & \\ & & x_n y_m & \\ & & & \ddots \end{bmatrix}$$

$$= \boxed{\vec{x} \vec{y}^T}$$

d) $\nabla_{\vec{x}} f = \nabla_{\vec{x}} (\vec{x}^T A \vec{x}) + \nabla_{\vec{x}} (\vec{b}^T \vec{x})$.

From class, we know the gradient of the first term is $A \vec{x} + A^T \vec{x}$.

From (a), we know $\nabla_{\vec{x}} (\vec{b}^T \vec{x}) = \vec{b}$. Since gradient is a linear operator,

$$\boxed{\nabla_{\vec{x}} f = A \vec{x} + A^T \vec{x} + \vec{b}}$$

$$e) \operatorname{tr}(AB) = \sum_i \sum_j (ab)_{ij} = \sum_i \sum_j b_{ji} a_{ij}$$

$$\bullet \text{ Then } \frac{\partial f}{\partial a_{11}} = \frac{\partial}{\partial a_{11}} (a_{11}b_{11} + a_{12}b_{21} + \dots + a_{nm}b_{mn}) = b_{11},$$

$$\frac{\partial f}{\partial a_{12}} = \frac{\partial}{\partial a_{12}} (a_{11}b_{11} + a_{12}b_{21} + \dots + a_{nm}b_{mn}) = b_{21}, \text{ and}$$

$$\nabla_A f = \begin{bmatrix} b_{11} & b_{21} & b_{31} & \dots & b_{m1} \\ b_{12} & & & & \\ \vdots & & & & \\ b_{1m} & & & & \end{bmatrix}$$

these are swapped in B

As we can see, the matrix is flipped from B.

$$\text{Hence, } \boxed{\nabla_A f = B^T}$$

$$\textcircled{4} \text{ Since } \|\vec{x}\|^2 = \operatorname{tr}(\vec{x}^T \vec{x}),$$

$$f = \frac{1}{2} \sum_{i=1}^n \|y^i - W x^i\|^2 = \frac{1}{2} \sum_{i=1}^n \operatorname{tr}((y^i - W x^i)^T (y^i - W x^i)).$$

Then, using FOIL / distributing,

$$f = \frac{1}{2} \sum_{i=1}^n \operatorname{tr}(y^i{}^T y^i - y^i{}^T W x^i - (W x^i)^T y^i + (W x^i)^T (W x^i)).$$

Then, since $\operatorname{tr}(A+B) = \operatorname{tr}(A) + \operatorname{tr}(B)$, and $\operatorname{tr}(AB) = \operatorname{tr}(BA)$,

$$f = \frac{1}{2} \sum_{i=1}^n (\operatorname{tr}(y^i{}^T y^i) - \operatorname{tr}(y^i{}^T W x^i) - \operatorname{tr}((W x^i)^T y^i) + \operatorname{tr}((W x^i)^T (W x^i)))$$

Now we apply $\partial/\partial W$, noting $\frac{\partial}{\partial X} \operatorname{tr}(AXB) = A^T B^T$ and $\frac{\partial}{\partial X} \operatorname{tr}(AX^T B) = BA$.

$$\text{Then } \frac{\partial f}{\partial W} = \frac{1}{2} \sum_{i=1}^n (0 - y^i{}^T x^i{}^T - y^i{}^T x^i{}^T + \frac{\partial}{\partial W} \operatorname{tr}((W x^i)^T (W x^i)))$$

$$\text{Using the hint, we get } \frac{\partial}{\partial W} \operatorname{tr}((W x^i)^T (W x^i)) = \frac{\partial}{\partial W} \operatorname{tr}((W x^i)(W x^i)^T)$$

$$= \frac{\partial}{\partial W} \operatorname{tr}(W x^i x^i{}^T W) = W (x^i x^i{}^T)^T + W (x^i x^i{}^T)$$

Putting it all together, we get

$$D_i = -2 y^i x^i{}^T + W (x^i x^i{}^T)^T + W (x^i x^i{}^T)$$

$$D_i = -2 y^i x^i{}^T + 2W x^i x^i{}^T, \text{ where}$$

$$\frac{\partial f}{\partial W} = \frac{1}{2} \sum_{i=1}^n D_i = 0$$

Now we solve for W . We have

$$\frac{1}{2} (-2 y^i x^i{}^T + 2W x^i x^i{}^T) = 0$$

$$W x^i x^i{}^T = y^i x^i{}^T$$

$$W = y^i x^i{}^T (x^i x^i{}^T)^{-1} \text{ for } i \in [1, n], \text{ since } (x^i x^i{}^T) \text{ is symmetric so is invertible.}$$

Since this holds for all $i \in [1, n]$, we can collapse the x^i and y^i to get:

$$\boxed{W = y x^T (x x^T)^{-1}}$$