

## 1. [0 points] Gradients and Hessians

Recall that a matrix  $A \in \mathbb{R}^{n \times n}$  is *symmetric* if  $A^T = A$ , that is,  $A_{ij} = A_{ji}$  for all  $i, j$ . Also recall the gradient  $\nabla f(x)$  of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , which is the  $n$ -vector of partial derivatives

$$\nabla f(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \vdots \\ \frac{\partial}{\partial x_n} f(x) \end{bmatrix} \quad \text{where } x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}.$$

The hessian  $\nabla^2 f(x)$  of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the  $n \times n$  symmetric matrix of twice partial derivatives,

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} f(x) & \frac{\partial^2}{\partial x_1 \partial x_2} f(x) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f(x) \\ \frac{\partial^2}{\partial x_2 \partial x_1} f(x) & \frac{\partial^2}{\partial x_2^2} f(x) & \cdots & \frac{\partial^2}{\partial x_2 \partial x_n} f(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f(x) & \frac{\partial^2}{\partial x_n \partial x_2} f(x) & \cdots & \frac{\partial^2}{\partial x_n^2} f(x) \end{bmatrix}.$$

- Let  $f(x) = \frac{1}{2}x^T A x + b^T x$ , where  $A$  is a symmetric matrix and  $b \in \mathbb{R}^n$  is a vector. What is  $\nabla f(x)$ ?
- Let  $f(x) = g(h(x))$ , where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is differentiable and  $h : \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable. What is  $\nabla f(x)$ ?
- Let  $f(x) = \frac{1}{2}x^T A x + b^T x$ , where  $A$  is symmetric and  $b \in \mathbb{R}^n$  is a vector. What is  $\nabla^2 f(x)$ ?
- Let  $f(x) = g(a^T x)$ , where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is continuously differentiable and  $a \in \mathbb{R}^n$  is a vector. What are  $\nabla f(x)$  and  $\nabla^2 f(x)$ ? (Hint: your expression for  $\nabla^2 f(x)$  may have as few as 11 symbols, including ' and parentheses.)

(a)  $Ax + b$

(b)  $h'(x) \cdot g'(h(x))$

(c)  $A$

(d)  ~~$\nabla f(x) = a$~~   $\nabla f(x) = a \cdot g'(a^T x)$

$\nabla^2 f(x) = \cancel{a^T a g''(a^T x)} \quad a^T a g''(a^T x)$

## 2. [0 points] Positive definite matrices

A matrix  $A \in \mathbb{R}^{n \times n}$  is *positive semi-definite* (PSD), denoted  $A \succeq 0$ , if  $A = A^T$  and  $x^T A x \geq 0$  for all  $x \in \mathbb{R}^n$ . A matrix  $A$  is *positive definite*, denoted  $A \succ 0$ , if  $A = A^T$  and  $x^T A x > 0$  for all  $x \neq 0$ , that is, all non-zero vectors  $x$ . The simplest example of a positive definite matrix is the identity  $I$  (the diagonal matrix with 1s on the diagonal and 0s elsewhere), which satisfies  $x^T I x = \|x\|_2^2 = \sum_{i=1}^n x_i^2$ .

- Let  $z \in \mathbb{R}^n$  be an  $n$ -vector. Show that  $A = zz^T$  is positive semidefinite.
- Let  $z \in \mathbb{R}^n$  be a non-zero  $n$ -vector. Let  $A = zz^T$ . What is the null-space of  $A$ ? What is the rank of  $A$ ?
- Let  $A \in \mathbb{R}^{n \times n}$  be positive semidefinite and  $B \in \mathbb{R}^{m \times n}$  be arbitrary, where  $m, n \in \mathbb{N}$ . Is  $BAB^T$  PSD? If so, prove it. If not, give a counterexample with explicit  $A, B$ .

(a)  $\text{rank}(A) = 1$  since every row/column is proportional to  $z$ .

The rank of  $A$  is equal to the number of non-zero eigenvalues of  $A$ .

$$Az = z(z^T z) = (z^T z)z = \lambda z$$

( $z^T z$  is a scalar and a scalar is commutative)

$$\lambda = z^T z = \|z\|^2 \geq 0$$

$A$  is a symmetric matrix and the definiteness of a symmetric matrix depends entirely on the sign of its eigenvalues.

Only non-zero eigenvalue  $\lambda > 0$

$A$  is positive semidefinite

$$(b) N(A) = \mathbb{R}^n - R(A)$$

$$R(A) = \{v \in \mathbb{R}^n : v = xz, x \in \mathbb{R}\}$$

$$\text{rank}(A) = 1$$

$$(c) \text{ Let } x \in \mathbb{R}^m$$

$$x^T BAB^T x = y^T A y \quad (\text{cancel } B^T x) \quad (B^T x = y) \geq 0$$

$BAB^T$  is PSD.

## 3. [0 points] Eigenvectors, eigenvalues, and the spectral theorem

The eigenvalues of an  $n \times n$  matrix  $A \in \mathbb{R}^{n \times n}$  are the roots of the characteristic polynomial  $p_A(\lambda) = \det(\lambda I - A)$ , which may (in general) be complex. They are also defined as the values  $\lambda \in \mathbb{C}$  for which there exists a vector  $x \in \mathbb{C}^n$  such that  $Ax = \lambda x$ . We call such a pair  $(x, \lambda)$  an *eigenvector, eigenvalue* pair. In this question, we use the notation  $\text{diag}(\lambda_1, \dots, \lambda_n)$  to denote the diagonal matrix with diagonal entries  $\lambda_1, \dots, \lambda_n$ , that is,

$$\text{diag}(\lambda_1, \dots, \lambda_n) = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ 0 & 0 & \lambda_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \lambda_n \end{bmatrix}.$$

- (a) Suppose that the matrix  $A \in \mathbb{R}^{n \times n}$  is diagonalizable, that is,  $A = T\Lambda T^{-1}$  for an invertible matrix  $T \in \mathbb{R}^{n \times n}$ , where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  is diagonal. Use the notation  $t^{(i)}$  for the columns of  $T$ , so that  $T = [t^{(1)} \dots t^{(n)}]$ , where  $t^{(i)} \in \mathbb{R}^n$ . Show that  $At^{(i)} = \lambda_i t^{(i)}$ , so that the eigenvalues/eigenvector pairs of  $A$  are  $(t^{(i)}, \lambda_i)$ .

A matrix  $U \in \mathbb{R}^{n \times n}$  is orthogonal if  $U^T U = I$ . The spectral theorem, perhaps one of the most important theorems in linear algebra, states that if  $A \in \mathbb{R}^{n \times n}$  is symmetric, that is,  $A = A^T$ , then  $A$  is *diagonalizable by a real orthogonal matrix*. That is, there are a diagonal matrix  $\Lambda \in \mathbb{R}^{n \times n}$  and orthogonal matrix  $U \in \mathbb{R}^{n \times n}$  such that  $U^T A U = \Lambda$ , or, equivalently,

$$A = U \Lambda U^T.$$

Let  $\lambda_i = \lambda_i(A)$  denote the  $i$ th eigenvalue of  $A$ .

- (b) Let  $A$  be symmetric. Show that if  $U = [u^{(1)} \dots u^{(n)}]$  is orthogonal, where  $u^{(i)} \in \mathbb{R}^n$  and  $A = U \Lambda U^T$ , then  $u^{(i)}$  is an eigenvector of  $A$  and  $Au^{(i)} = \lambda_i u^{(i)}$ , where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ .

- (c) Show that if  $A$  is PSD, then  $\lambda_i(A) \geq 0$  for each  $i$ .

(a)  $A = T \Lambda T^{-1} \Rightarrow A T = T \Lambda$   
 $A T = [A t^{(1)} \dots A t^{(n)}] \quad T \Lambda = [T \lambda^{(1)} \dots T \lambda^{(n)}] = [t^{(1)} \lambda_1 \dots t^{(n)} \lambda_n]$   
 $[A t^{(1)} \dots A t^{(n)}] = [t^{(1)} \lambda_1 \dots t^{(n)} \lambda_n]$

(b)  $A = U \Lambda U^T \Rightarrow A U = U \Lambda U^T U = U \Lambda$   
 $A U = [A u^{(1)} \dots A u^{(n)}] \quad U \Lambda = [U \lambda^{(1)} \dots U \lambda^{(n)}] = [u^{(1)} \lambda_1 \dots u^{(n)} \lambda_n]$   
 $[A u^{(1)} \dots A u^{(n)}] = [u^{(1)} \lambda_1 \dots u^{(n)} \lambda_n]$



(C) Show that if  $A$  is PSD, then  $\lambda_i(A) \geq 0$  for each  $i$ .

$$X^T A X = X^T U \Lambda U^T X = \hat{X}^T \Lambda \hat{X} = \sum_{i=1}^n \lambda_i \hat{x}_i^2 \geq 0$$

To  $\sum_{i=1}^n \lambda_i \hat{x}_i^2 \geq 0$  be true for all  $\hat{x}$ ,

$\lambda_i(A) \geq 0$  must be true.

## 4. [0 points] Probability and multivariate Gaussians

Suppose  $X = (X_1, \dots, X_n)$  is sampled from a multivariate Gaussian distribution with mean  $\mu$  in  $\mathbb{R}^n$  and covariance  $\Sigma$  in  $S_+^n$  (i.e.  $\Sigma$  is positive semidefinite). This is commonly also written as  $X \sim \mathcal{N}(\mu, \Sigma)$ .

(a) Describe the random variable  $Y = X_1 + X_2 + \dots + X_n$ . What is the mean and variance? Is this a well known distribution, and if so, which?

(b) Now, further suppose that  $\Sigma$  is invertible. Find  $\mathbb{E}[X^T \Sigma^{-1} X]$ . (Hint: use the property of trace that  $x^T A x = \text{tr}(x^T A x)$ ).

(a)  $Y$  is the linear combination of the elements of a multivariate normal random variable.

mean:  $\mu$

variance:  $\Sigma$

This is a multivariate Gaussian distribution.

$$\begin{aligned}
 (b) \mathbb{E}[X^T \Sigma^{-1} X] &= \mathbb{E}[\text{tr}(X^T \Sigma^{-1} X)] \\
 &= \mathbb{E}[\text{tr}(X X^T \Sigma^{-1})] \\
 &= \text{tr}(\mathbb{E}[X X^T \Sigma^{-1}]) \\
 &= \text{tr}(\Sigma^{-1} \mathbb{E}[X X^T]) \\
 &= \text{tr}(\Sigma^{-1} (\Sigma + \mu \mu^T)) \\
 &= \text{tr}(I + \Sigma^{-1} \mu \mu^T) \\
 &= \text{tr}(I) + \text{tr}(\Sigma^{-1} \mu \mu^T) \\
 &= n + \text{tr}(\mu^T \Sigma^{-1} \mu) \\
 &= n + \mu^T \Sigma^{-1} \mu
 \end{aligned}$$