# HW 1 for CSCI 5525
**YangyangLi**
**li002252@umn.edu**

# 1 Problem 1

## 1.1 a

Show that $\langle A, B \rangle = \operatorname{tr}(A^\mathsf{T} B)$ and so $\|M\|_F = \sqrt{\operatorname{tr}(M^\mathsf{T} M)}$
**PROOF:**

$\because \langle A, B \rangle = \sum_{ij} a_{ij} b_{ij}$ and $A$ and $B$ have same size.

$\therefore \langle A, B \rangle = \sum_j a_j^\mathsf{T} b_j$, where $a_j$ and $b_j$ are column vector of $A$ and $B$ respectively.

$\because A^\mathsf{T} B = [a_1, \cdots, a_j]^\mathsf{T}[b_1, \cdots, b_j]$ and $\operatorname{tr}(M) = \sum_i m_{ii}$.

$$\therefore A^\mathsf{T} B = \begin{bmatrix} a_1^\mathsf{T} b_1 & \cdots & \cdots & \cdots & \cdots \\ \cdots & a_2^\mathsf{T} b_2 & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & a_j^\mathsf{T} b_j \end{bmatrix}$$

$\therefore \operatorname{tr}(A^\mathsf{T} B) = \sum_j a_j^\mathsf{T} b_j = \langle A, B \rangle$ and $\operatorname{tr}(M^\mathsf{T} M) = \langle M, M \rangle$

Summarizing:

$$\langle A, B \rangle = \operatorname{tr}(A^\mathsf{T} B) \tag{1}$$

$$\|M\|_F = \sqrt{\operatorname{tr}(M^\mathsf{T} M)} \tag{2}$$

## 1.2 b

Show that $\operatorname{tr}(A^\mathsf{T} B) = \operatorname{tr}(B^\mathsf{T} A)$
**PROOF:**

$\because \langle A, B \rangle = \operatorname{tr}(A^\mathsf{T} B)$ and $\langle A, B \rangle = \langle B, A \rangle$
$\therefore \operatorname{tr}(A^\mathsf{T} B) = \langle A, B \rangle = \langle B, A \rangle = \operatorname{tr}(B^\mathsf{T} A)$

Summarizing:

$$\operatorname{tr}(A^\mathsf{T} B) = \operatorname{tr}(B^\mathsf{T} A) \tag{3}$$

## 1.3 c

Assume $A$ and $B$ have the same size. In general, $AB^\mathsf{T}$ and $B^\mathsf{T} A$ have different sizes, but $\operatorname{tr}(AB^\mathsf{T}) = \operatorname{tr}(B^\mathsf{T} A)$. Show it.

**PROOF:**

$\because \operatorname{tr}(\boldsymbol{A}\boldsymbol{B}^\mathsf{T}) = \langle \boldsymbol{A}^\mathsf{T}, \boldsymbol{B}^\mathsf{T} \rangle = \sum_{ji} a_{ji} b_{ji}$, where $a_{ji}$ and $b_{ji}$ are $(j, i)$-th element of $\boldsymbol{A}^\mathsf{T}$ and $\boldsymbol{B}^\mathsf{T}$ respectively.

$\because \operatorname{tr}(\boldsymbol{B}^\mathsf{T}\boldsymbol{A}) = \langle \boldsymbol{A}, \boldsymbol{B} \rangle = \sum_{ij} a_{ij} b_{ij}$, where $a_{ij}$ and $b_{ij}$ are $(i, j)$-th element of $\boldsymbol{A}$ and $\boldsymbol{B}$ respectively.

$\therefore \operatorname{tr}(\boldsymbol{A}\boldsymbol{B}^\mathsf{T}) = \operatorname{tr}(\boldsymbol{B}^\mathsf{T}\boldsymbol{A})$

Summarizing:

$$\operatorname{tr}(\boldsymbol{A}\boldsymbol{B}^\mathsf{T}) = \operatorname{tr}(\boldsymbol{B}^\mathsf{T}\boldsymbol{A}) \tag{4}$$

## 1.4  d

Show that $\operatorname{tr}(\boldsymbol{M}_1\boldsymbol{M}_2\boldsymbol{M}_3) = \operatorname{tr}(\boldsymbol{M}_3\boldsymbol{M}_1\boldsymbol{M}_2) = \operatorname{tr}(\boldsymbol{M}_2\boldsymbol{M}_3\boldsymbol{M}_1)$, assuming that the sizes of $\boldsymbol{M}_1$, $\boldsymbol{M}_2$ and $\boldsymbol{M}_3$ are compatible with all the matrix multiplications. This is known as the *cyclic property* of matrix traces. (Hint: think of (c))
**PROOF:**

Provided that $\boldsymbol{M}_1\boldsymbol{M}_2 = \boldsymbol{A}$, $\boldsymbol{M}_3 = \boldsymbol{B}^\mathsf{T}$, according to the Eq. (4):

$$\operatorname{tr}(\boldsymbol{M}_1\boldsymbol{M}_2\boldsymbol{M}_3) = \operatorname{tr}(\boldsymbol{A}\boldsymbol{B}^\mathsf{T}) = \operatorname{tr}(\boldsymbol{B}^\mathsf{T}\boldsymbol{A}) = \operatorname{tr}(\boldsymbol{M}_3\boldsymbol{M}_1\boldsymbol{M}_2)$$

Similarly,

$$\operatorname{tr}(\boldsymbol{M}_3\boldsymbol{M}_1\boldsymbol{M}_2) = \operatorname{tr}(\boldsymbol{M}_2\boldsymbol{M}_3\boldsymbol{M}_1)$$

Summarizing:

$$\operatorname{tr}(\boldsymbol{M}_1\boldsymbol{M}_2\boldsymbol{M}_3) = \operatorname{tr}(\boldsymbol{M}_3\boldsymbol{M}_1\boldsymbol{M}_2) = \operatorname{tr}(\boldsymbol{M}_2\boldsymbol{M}_3\boldsymbol{M}_1) \tag{5}$$

## 1.5  e

For any matrices $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}, \boldsymbol{D}$ of compatible sizes, we always have $\langle \boldsymbol{A}\boldsymbol{C}\boldsymbol{B}, \boldsymbol{D} \rangle = \langle \boldsymbol{C}\boldsymbol{B}, \boldsymbol{A}^\mathsf{T}\boldsymbol{D} \rangle = \langle \boldsymbol{A}\boldsymbol{C}, \boldsymbol{D}\boldsymbol{B}^\mathsf{T} \rangle$, i.e., we can always move the **leading** matrix of one side of the inner product to the other side as **leading** matrix **once transposed** (if these matrices are complex-valued, should be conjugate transposed), and similarly the **trailing** matrix to the other side as **trailing** matrix once **transposed**. Why? (Hint: think of the above results and also try to remember this important property that will be useful for calculation later)
**PROOF:**

In terms of the Eq. (1):

$$\langle \boldsymbol{A}\boldsymbol{C}\boldsymbol{B}, \boldsymbol{D} \rangle = \langle \boldsymbol{D}, \boldsymbol{A}\boldsymbol{C}\boldsymbol{B} \rangle = \operatorname{tr}(\boldsymbol{D}^\mathsf{T}\boldsymbol{A}\boldsymbol{C}\boldsymbol{B})$$

Provided that $\boldsymbol{D}^\mathsf{T}\boldsymbol{A} = \boldsymbol{M}_1$ and $\boldsymbol{C}\boldsymbol{B} = \boldsymbol{M}_2$. In terms of the Eq. (1):

$$\operatorname{tr}(\boldsymbol{D}^\mathsf{T}\boldsymbol{A}\boldsymbol{C}\boldsymbol{B}) = \operatorname{tr}(\boldsymbol{M}_1\boldsymbol{M}_2) = \langle \boldsymbol{M}_1^\mathsf{T}, \boldsymbol{M}_2 \rangle = \langle \boldsymbol{A}^\mathsf{T}\boldsymbol{D}, \boldsymbol{C}\boldsymbol{B} \rangle = \langle \boldsymbol{C}\boldsymbol{B}, \boldsymbol{A}^\mathsf{T}\boldsymbol{D} \rangle$$

Similarly, in terms of the Eq. (5):

$$\langle ACB, D \rangle = \langle D, ACB \rangle = \operatorname{tr}(D^\mathsf{T}ACB) = \operatorname{tr}(BD^\mathsf{T}AC) = \langle DB^\mathsf{T}, AC \rangle$$

Summarizing:

$$\langle ACB, D \rangle = \langle CB, A^\mathsf{T}D \rangle = \langle AC, DB^\mathsf{T} \rangle \tag{6}$$

## 1.6 f

For $M$, let's perform a *compact SVD* (if not sure, check up Wikipedia! https://en.wikipedia.org/wiki/Singular_value_decomposition#Compact_SVD) to obtain $M = U\Sigma V^\mathsf{T}$, so that $U$ and $V$ are orthonormal (not necessarily square) matrices, i.e., $U^\mathsf{T}U = I$ and $V^\mathsf{T}V = I$. Use the cyclic property of trace and that $\|M\|_F = \sqrt{\operatorname{tr}(M^\mathsf{T}M)}$ to show that

$$\|M\|_F = \sqrt{\sum_{i=1}^{r} \sigma_i^2},$$

assuming the rank of $M$ is $r$. Here $\sigma_i$'s are the singular values of $M$.
**PROOF**:

since $M = U\Sigma V^\mathsf{T}$, $\|M\|_F = \sqrt{\operatorname{tr}(M^\mathsf{T}M)}$, $U^\mathsf{T}U = I$ and $V^\mathsf{T}V = I$ so,

$$\|M\|_F = \sqrt{\operatorname{tr}(V\Sigma U^\mathsf{T}U\Sigma V^\mathsf{T})} = \sqrt{\operatorname{tr}(V\Sigma\Sigma V^\mathsf{T})} = \sqrt{\operatorname{tr}(V^\mathsf{T}V\Sigma\Sigma)} = \sqrt{\operatorname{tr}(\Sigma_r\Sigma_r)} = \sqrt{\sum_{i=1}^{r}\sigma_i^2}$$

# 2 Problem 2

## 2.1 a

Let $A$ be a square matrix. Deriving the gradient and Hessian of the quadratic function $f(x) = x^\mathsf{T}Ax + b^\mathsf{T}x$. Please include your calculation details. (Hint: note that Hessian must be a symmetric matrix.
**SOLUTION**:

$$\begin{aligned}
f(x+\sigma) &= (x+\sigma)^\mathsf{T}A(x+\sigma) + b^\mathsf{T}(x+\sigma) \\
&= (x^\mathsf{T}+\sigma^\mathsf{T})(Ax+A\sigma) + b^\mathsf{T}x + b^\mathsf{T}\sigma \\
&= x^\mathsf{T}Ax + x^\mathsf{T}A\sigma + \sigma^\mathsf{T}Ax + \sigma^\mathsf{T}A\sigma + b^\mathsf{T}x + b^\mathsf{T}\sigma \\
&= f(x) + \sigma^\mathsf{T}Ax + x^\mathsf{T}A\sigma + b^\mathsf{T}\sigma + \sigma^\mathsf{T}A\sigma
\end{aligned} \tag{7}$$

$\because \sigma^\mathsf{T}Ax$ is a scalar

$\therefore \sigma^\mathsf{T}Ax = (\sigma^\mathsf{T}Ax)^\mathsf{T} = x^\mathsf{T}A^\mathsf{T}\sigma$

$$\begin{aligned}
f(x+\sigma) &= f(x) + (x^\mathsf{T}A^\mathsf{T} + x^\mathsf{T}A + b^\mathsf{T})\sigma + \sigma^\mathsf{T}A\sigma \\
&= f(x) + \langle Ax + A^\mathsf{T}x + b, \sigma \rangle + \langle \sigma, A\sigma \rangle
\end{aligned} \tag{8}$$

Consequently,

$$\nabla f\left(\boldsymbol{x}\right) = \boldsymbol{A}\boldsymbol{x} + \boldsymbol{A}^{\mathsf{T}}\boldsymbol{x} + \boldsymbol{b}$$

$$\nabla^2 f\left(\boldsymbol{x}\right) = 2\boldsymbol{A}$$

## 2.2  b

Let $p\left(\boldsymbol{x}; \boldsymbol{\beta}\right) = \frac{e^{\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x}}}{1+e^{\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x}}}$. The log-likelihood for logistic regression with two classes is (assuming $N$ samples of the form $\left(\boldsymbol{x}_i, y_i\right)$)

$$
\begin{aligned}
f\left(\boldsymbol{\beta}\right) &= \sum_{i=1}^{N} \left[y_i \log p\left(\boldsymbol{x}_i; \boldsymbol{\beta}\right) + (1-y_i) \log\left(1 - p\left(\boldsymbol{x}_i; \boldsymbol{\beta}\right)\right)\right] \\
&= \sum_{i=1}^{N} \left[y_i \boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x}_i - \log\left(1 + e^{\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x}_i}\right)\right].
\end{aligned}
$$

Derive the gradient and Hessian of $f\left(\boldsymbol{\beta}\right)$. Please include your calculation details. (1/12) For logistic regression, we are going to maximize $f\left(\boldsymbol{\beta}\right)$, which is equivalent to minimize $-f\left(\boldsymbol{\beta}\right)$. Does the minimization problem has a unique minimizer or not?
**SOLUTION:**

denote $f\left(\boldsymbol{\beta}\right)$ as matrix form,

$$
\begin{aligned}
f\left(\boldsymbol{\beta}\right) &= \sum_{i=1}^{N} \left[y_i \boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x}_i - \log\left(1 + e^{\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{x}_i}\right)\right] \\
&= \boldsymbol{y}^{\mathsf{T}}\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{1}^{\mathsf{T}} \log\left(1 + e^{\boldsymbol{X}\boldsymbol{\beta}}\right)
\end{aligned}
\tag{9}
$$

where $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^{\mathsf{T}} \\ \boldsymbol{x}_2^{\mathsf{T}} \\ \vdots \\ \boldsymbol{x}_N^{\mathsf{T}} \end{bmatrix}$, $\boldsymbol{y} = [y_1, \ldots, y_N]^{\mathsf{T}}$, $\boldsymbol{1} = [1, \ldots, 1]$.

Since $f\left(\boldsymbol{\beta}\right)$ map vector to scalar, $df = \sum_i \frac{\partial f}{\partial \beta_i} d\beta_i = \left(\frac{\partial f}{\partial \boldsymbol{\beta}}\right)^{\mathsf{T}} d\boldsymbol{\beta} = \nabla_{\boldsymbol{\beta}} f^{\mathsf{T}} d\boldsymbol{\beta}$

$$
\begin{aligned}
df &= \boldsymbol{y}^{\mathsf{T}}\boldsymbol{X} d\boldsymbol{\beta} - \frac{\boldsymbol{1}^{\mathsf{T}}\left(e^{\boldsymbol{X}\boldsymbol{\beta}} \circ d\boldsymbol{X}\boldsymbol{\beta}\right)}{1 + e^{\boldsymbol{X}\boldsymbol{\beta}}} \\
&= \boldsymbol{y}^{\mathsf{T}}\boldsymbol{X} d\boldsymbol{\beta} - \frac{\left(e^{\boldsymbol{X}\boldsymbol{\beta}}\right)^{\mathsf{T}} d\boldsymbol{X}\boldsymbol{\beta}}{1 + e^{\boldsymbol{X}\boldsymbol{\beta}}}
\end{aligned}
\tag{10}
$$

Let $\sigma\left(a\right) = \frac{e^a}{1+e^a}$, $\sigma'\left(a\right) = \frac{e^a}{\left(1+e^a\right)^2}$,

$$
\begin{aligned}
df &= \boldsymbol{y}^\mathsf{T}\boldsymbol{X}d\boldsymbol{\beta} - \frac{\mathbf{1}^\mathsf{T}\left(e^{\boldsymbol{X}\boldsymbol{\beta}} \circ d\boldsymbol{X}\boldsymbol{\beta}\right)}{1+e^{\boldsymbol{X}\boldsymbol{\beta}}} \\
&= \boldsymbol{y}^\mathsf{T}\boldsymbol{X}d\boldsymbol{\beta} - \frac{\left(e^{\boldsymbol{X}\boldsymbol{\beta}}\right)^\mathsf{T}\boldsymbol{X}d\boldsymbol{\beta}}{1+e^{\boldsymbol{X}\boldsymbol{\beta}}} \\
&= \boldsymbol{y}^\mathsf{T}\boldsymbol{X}d\boldsymbol{\beta} - \sigma\left(\boldsymbol{X}\boldsymbol{\beta}\right)^\mathsf{T}\boldsymbol{X}d\boldsymbol{\beta} \\
&= \left(\boldsymbol{y}^\mathsf{T} - \sigma\left(\boldsymbol{X}\boldsymbol{\beta}\right)^\mathsf{T}\right)\boldsymbol{X}d\boldsymbol{\beta}
\end{aligned}
\tag{11}
$$

So $\nabla f = \boldsymbol{X}^T\left(\boldsymbol{y} - \sigma\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right)$. $\nabla^2 f = \frac{\partial^2 f}{\partial\boldsymbol{\beta}\partial\boldsymbol{\beta}}$, so,

$$
\begin{aligned}
d\nabla f &= d\left(\boldsymbol{X}^\mathsf{T}\left(\boldsymbol{y} - \sigma\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right)\right) \\
&= -\boldsymbol{X}^\mathsf{T}d\sigma\left(\boldsymbol{X}\boldsymbol{\beta}\right) \\
&= -\boldsymbol{X}^\mathsf{T}\left(\sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right) \circ d\boldsymbol{X}\boldsymbol{\beta}\right) \\
&= -\boldsymbol{X}^\mathsf{T}\operatorname{diag}\left(\sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right)\boldsymbol{X}d\boldsymbol{\beta}
\end{aligned}
\tag{12}
$$

Similarly,

$$
\nabla^2 f = -\boldsymbol{X}^\mathsf{T}\operatorname{diag}\left(\sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right)\boldsymbol{X}, \text{ where } \sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right) = \frac{e^{\boldsymbol{X}\boldsymbol{\beta}}}{\left(1+e^{\boldsymbol{X}\boldsymbol{\beta}}\right)^2}
$$

Consequently,

$$
\nabla f = \boldsymbol{X}^T\left(\boldsymbol{y} - \sigma\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right),
$$

$$
\nabla^2 f = -\boldsymbol{X}^\mathsf{T}\operatorname{diag}\left(\sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right)\boldsymbol{X}
$$

where $\sigma\left(\boldsymbol{X}\boldsymbol{\beta}\right) = \frac{e^{\boldsymbol{X}\boldsymbol{\beta}}}{1+e^{\boldsymbol{X}\boldsymbol{\beta}}}$, $\sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right) = \frac{e^{\boldsymbol{X}\boldsymbol{\beta}}}{\left(1+e^{\boldsymbol{X}\boldsymbol{\beta}}\right)^2}$.

Also,

$$
\begin{aligned}
\nabla^2\text{-}f &= \boldsymbol{X}^T\operatorname{diag}\left(\sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right)\boldsymbol{X} \\
&= \begin{bmatrix} \sum_i^N \lambda_i x_{i,1}^2 & \cdots & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & \cdots & \sum_i^N \lambda_i x_{i,M}^2 \end{bmatrix}
\end{aligned}
\tag{13}
$$

Where $\lambda_i$ is diagonal entries of $\operatorname{diag}\left(\sigma'\left(\boldsymbol{X}\boldsymbol{\beta}\right)\right)$, and $M$ is the number of column vector of $\boldsymbol{X}$. As $\sigma'\left(a\right) > 0$, diagonal entries of $\nabla^2\text{-}f$ more than $0$. So, $\nabla^2\text{-}f$ is positive definite.

Therefore, only one unique minimizer exists.

## 2.3  c

Let $A \in \mathbb{R}^{m \times n}$ with $m \leq n$, then given a $y \in \mathbb{R}^m$, the least-squares problem $\min_x \; \|y - Ax\|_2^2$ has infinitely many solutions. Now let's say we want a solution with a small $\ell_2$ norm, then it is reasonable to put a penalty on the $\ell_2$ norm:

$$\min_x \; \|y - Ax\|_2^2 + \lambda \|x\|_2^2$$

with a chosen $\lambda > 0$. This is *ridge regression*. Now we know that for an unconstrained first-order differentiable function $g(x)$, any of its local minimizer $x_*$ must satisfy the *first-order optimality condition*: $\nabla g(x_*) = 0$. Use this to derive $x_*$ (1/12). Is the $x_*$ unique? Why?

**SOLUTION:**

$$
\begin{aligned}
g(x + \sigma) &= \|y - A(x + \sigma)\|_2^2 + \lambda \|x + \sigma\|_2^2 \\
&= \|y - Ax\|_2^2 + \lambda \|x\|_2^2 + \lambda \|\sigma\|_2^2 + 2\lambda \langle x, \sigma \rangle - 2 \langle A^\mathsf{T}(y - Ax), \sigma \rangle + \|A\sigma\|_2^2 \\
&= g(x) + \sigma^\mathsf{T}(\lambda I + A^\mathsf{T}A)\sigma + 2(\lambda x^\mathsf{T} - (y - Ax)^\mathsf{T}A)\sigma \\
&= g(x) + 2\langle \lambda x - A^\mathsf{T}(y - Ax), \sigma \rangle + \langle \sigma, (\lambda I + A^\mathsf{T}A)\sigma \rangle
\end{aligned}
\tag{14}
$$

So,

$$\nabla g = 2(\lambda x - A^\mathsf{T}y + A^\mathsf{T}Ax)$$

$$\nabla^2 g = 2(\lambda I + A^\mathsf{T}A)$$

Let $\nabla g(x_*)) = 0$,

$$x_* = (\lambda I + A^\mathsf{T}A)^{-1} A^\mathsf{T}y$$

$A^\mathsf{T}A$ is positive definite in that its diagonal entries are all positive. Hence, $\nabla^2 g = 2(\lambda I + A^\mathsf{T}A)$ is positive definite as well, which indicates $x_*$ is unique.

# 3  Problem 3

The solution for this problem has been finished in the Colab, and I will upload this notebook *HW1_CS5525_YangyangLi.ipynb* as well. Feel free to check that.

Have a good day!