PROBLEM SHEET 2 FOR ADVANCED MACHINE LEARNING (COMP6208)

This problem sheet asks you to prove some well known results. Although the algebra is easy the proofs are not entirely straightforward. There are marks assigned to the readability of the solution and also how well laid out and explained the steps you make are. (A good proof needs to be easy to follow: you need not comment on trivial algebra, but there should not be steps that are difficult to follow).

This looks very mathematical, but it helps to develop the tools and language that is used to describe machine learning.

**1**

(a) Starting from the definition of a convex function where, for $a \in [0,1]$,

$$f(ax + (1-a)y) \leq af(x) + (1-a)f(y) \tag{1}$$

Let $a = \epsilon/(x-y)$ and rearrange the inequality to give

$$(x-y)\left(\frac{f(y+\epsilon) - f(y)}{\epsilon}\right)$$

on the left-hand side. Taking the limit $\epsilon \to 0$ show that the function $f(x)$ lies above the tangent line $t(x) = f(y) + (x-y)f'(y)$ going through the point $y$.

[4 marks]

Let $a = \epsilon/(x-y)$, for the left-hand side:

$f(ax + (1-a)y) = f\left(\frac{x\epsilon}{x-y} + (1-\frac{\epsilon}{x-y})y\right) = f(y+\epsilon)$

for the right-hand side:

$af(x) + (1-a)f(y) = f(y) + a(f(x) - f(y)) = f(y) + \frac{\epsilon(f(x)-f(y))}{x-y}$

so the inequality becomes:

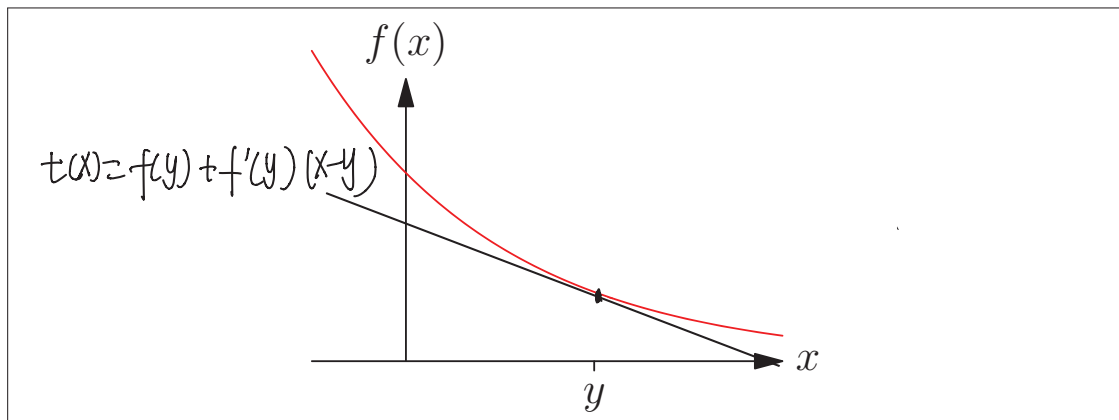$f(y+\epsilon) \leq f(y) + \frac{\epsilon(f(x)-f(y))}{x-y}$

finally:

$(x-y)\left(\frac{f(y+\epsilon)-f(y)}{\epsilon}\right) \leq f(x) - f(y)$

when $\epsilon \to 0$, inequality becomes: $(x-y)f'(y) \leq f(x) - f(y)$

So that: $f(y) + (x-y)f'(y) \leq f(x) \Rightarrow t(x) \leq f(x)$

so $f(x)$ lies above the tangent line $t(x)$ going through the point $y$.

$\overline{4}$

(b) Sketch the tangent line, $t(x)$, at the point $y$ in the graph shown below. [1 mark]

$f(x)$

$t(x) = f(y) + f'(y)(x-y)$

$x$

$y$

1

(c) Starting from the inequality for a convex function, $f$,

$$f(x) \geq f(y) + (x - y)f'(y) \tag{2}$$

consider the case $y = x + \epsilon$, then by Taylor expanding $f(x + \epsilon)$ and $f'(x + \epsilon)$ around $x$ and keeping all terms up to order $\epsilon^2$, show that $f''(x) \geq 0$. [4 marks]

by second-order Taylor expanding:
$$f(x) = f(y) + f'(y)(x-y) + f''(z)\frac{(x-y)^2}{2} \quad z \in [x,y]$$
So the inequality becomes:
$$f(y) + f'(y)(x-y) + f''(z)\frac{(x-y)^2}{2} \geq f(y) + (x-y)f'(y)$$
So $f''(z)\frac{\epsilon^2}{2} \geq 0$, as $\frac{\epsilon^2}{2} \geq 0$
So $f''(z) \geq 0$ for $z \in [x,y]$.
So that $f''(x) \geq 0$.

4

(d) Prove that $x^4$ is convex. [1 mark]

as $f''(x) = 12x^2 \geq 0$ for $\forall x \in C$. $x^4$ is convex.

1

End of question 1

| (a) $\frac{}{4}$ | (b) $\frac{}{1}$ | (c) $\frac{}{4}$ | (d) $\frac{}{1}$ | Total $\frac{}{10}$ |

**2**

(a) Show by writing out in component for that $\operatorname{tr}\mathbf{AB} = \operatorname{tr}\mathbf{BA}$ where $\operatorname{tr}\mathbf{M} = \sum_i M_{ii}$ (i.e. the trace of a matrix is equal to the sum of terms down the diagonal).

[2 marks]

Let $A$ is $n \times m$ matrix, $B$ is $m \times n$ matrix. $A = a_{ij}$, $B = b_{ij}$.

$(AB)_{ii} = \sum_{k=1}^{m} a_{ik} b_{ki}$     $\operatorname{tr}(AB) = \sum_{j=1}^{n} \sum_{k=1}^{m} a_{jk} b_{kj}$

$(BA)_{ii} = \sum_{k=1}^{n} b_{ik} a_{ki}$     $\operatorname{tr}(BA) = \sum_{k=1}^{m} \sum_{j=1}^{n} b_{kj} a_{jk}$

rearrange $\operatorname{tr}(BA) = \sum_{j=1}^{n} \sum_{k=1}^{m} a_{jk} b_{kj} = \operatorname{tr}(AB)$, so $\operatorname{tr}(AB) = \operatorname{tr}(BA)$.

$\boxed{2}$

(b) Using the fact that we can write a symmetric matrix $\mathbf{M}$ as $\mathbf{M} = \mathbf{V \Lambda V}^\mathsf{T}$ where $\mathbf{V}$ is an orthogonal matrix and $\mathbf{\Lambda} = \operatorname{diag}(\lambda_1, \lambda_2, \ldots)$ (i.e. a diagonal matrix with $\Lambda_{ii} = \lambda_i$). Show that $\operatorname{tr}\mathbf{M} = \sum_i \lambda_i$. [2 marks]

$\operatorname{tr}(M) = \operatorname{tr}(V \Lambda V^\mathsf{T}) = \operatorname{tr}(\Lambda V V^\mathsf{T})$

$= \operatorname{tr}(\Lambda I) \Leftarrow V V^\mathsf{T} = I$

$= \operatorname{tr}(\Lambda)$

$= \sum_i \lambda_i$

$\boxed{2}$

(c) Consider the matrix $\mathbf{X} = (\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_n)$ where the $i^{th}$ column of $\mathbf{X}$ is the vector $\boldsymbol{x}_i$. Compute $\operatorname{tr}\mathbf{X}^\mathsf{T}\mathbf{X}$ [2 marks]

$X^\mathsf{T} X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} x_1, x_2 \cdots x_n \end{bmatrix} = \begin{bmatrix} x_1^2 & x_1 x_2 & \cdots & x_1 x_n \\ x_2 x_1 & x_2^2 & & \vdots \\ \vdots & & \ddots & \\ x_n x_1 & \cdots & \cdots & x_n^2 \end{bmatrix}$

$\operatorname{tr}(X^\mathsf{T} X) = \sum_i x_i^2$

$\boxed{2}$

(d) The Frobenius norm, $\|\mathbf{X}\|_F$ for a matrix $\mathbf{X}$ is given by

$$\|\mathbf{X}\|_F = \sqrt{\sum_{i,j} X_{ij}^2},$$

where $X_{ij}$ is the $(i,j)$ component of $\mathbf{X}$. Using the previous result, show that $\|\mathbf{X}\|_F^2 = \operatorname{tr}\mathbf{X}^\mathsf{T}\mathbf{X}$ [2 marks]

Following by $c$, $\operatorname{tr}(X^TX) = \sum_i X_i^2$, $X_i$ is the $i^{th}$ column of $X$.

So $X_i = [X_{i1}, X_{i2}, X_{i3} \cdots X_{im}]$ If $X$ is $n \times m$ matrix

hence, $\operatorname{tr}(X^TX) = \sum_i \sum_j X_{ij}^2$

$\|X\|_F^2 = \sum_i \sum_j X_{ij}^2 = \operatorname{tr}(X^TX)$

$\boxed{2}$

(e) By using the SVD $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^\mathsf{T}$ where $\mathbf{S} = \operatorname{diag}(s_1, s_2, \ldots, s_n)$ (i.e. a diagonal matrix where $S_{ii} = s_i$—the $i^{th}$ singular value) and using the previous results, show that $\|\mathbf{X}\|_F^2 = \sum_i s_i^2$. [2 marks]

$\|X\|_F^2 = \operatorname{tr}(X^TX) = \operatorname{tr}(XX^T)$
$= \operatorname{tr}(USV^TVS^TU^T) \Leftarrow V^TV = I$
$= \operatorname{tr}(USS^TU^T) \Leftarrow S = S^T$
$= \operatorname{tr}(US^2U^T)$
$= \operatorname{tr}((S^2U^T)U) \Leftarrow \operatorname{tr}(AB) = \operatorname{tr}(BA)$.
$= \operatorname{tr}(S^2U^TU) \Leftarrow U^TU = I$
$= \operatorname{tr}(S^2)$
$= \sum_i S_i^2$

$\boxed{2}$

End of question 2

| (a) $\frac{\phantom{x}}{2}$ | (b) $\frac{\phantom{x}}{2}$ | (c) $\frac{\phantom{x}}{2}$ | (d) $\frac{\phantom{x}}{2}$ | (e) $\frac{\phantom{x}}{2}$ | Total $\frac{\phantom{x}}{10}$ |
|---|---|---|---|---|---|

**3** The $p$-norm of a matrix **M**, for $p \geq 1$ is defined to satisfy

$$\|\mathbf{M}\|_p = \max_{\boldsymbol{x} \neq \mathbf{0}} \frac{\|\mathbf{M}\boldsymbol{x}\|_p}{\|\boldsymbol{x}\|_p} \tag{3}$$

$$= \max_{\boldsymbol{x}:\|\boldsymbol{x}\|_p=1} \|\mathbf{M}\boldsymbol{x}\|_p \tag{4}$$

where $\|\boldsymbol{x}\|_p$ is the $p$ norm of a vector defined by

$$\|\boldsymbol{x}\|_p = \left( \sum_i |x_i|^p \right)^{1/p}.$$

Note that with this definition $\|\mathbf{M}\boldsymbol{x}\|_p \leq \|\mathbf{M}\|_p \|\boldsymbol{x}\|_p$ (where the inequality is tight, i.e. there exists a vector where the inequality becomes an equality).

(a) If **U** is an orthogonal matrix show that for any vector $v$ that $\|\mathbf{U}v\|_2 = \|v\|_2$. Use this to show $\|\mathbf{U}\mathbf{A}\|_2 = \|\mathbf{A}\|_2$. **[2 marks]**

$$\|Uv\|_2^2 = (UX)^T(UX) \qquad \|VA\|_2 = \sqrt{(UA)^T(UA)}$$
$$= X^TU^TUX \qquad\qquad = \sqrt{A^TU^TUA}$$
$$= \|X\|_2^2 \qquad\qquad\qquad = \sqrt{A^TA} = \|A\|_2$$
$$\text{So } \|Uv\|_2 = \|X\|_2$$

$\boxed{2}$

(b) If $V$ is an orthogonal matrix show that $\|\mathbf{A}\mathbf{V}^{\mathsf{T}}\|_2 = \|\mathbf{A}\|_2$. **[2 marks]**

$$\|AV^T\|_2 = \max_{\|X\|_p=1} \|AV^TX\| = \max_{\|V^TX\|_p=1} \|AY\| = \max_{\|Y\|_p=1} \|AY\|_2$$
$$= \|A\|_2$$

$\boxed{2}$

(c) Use the SVD $\mathbf{M} = \mathbf{U}\mathbf{S}\mathbf{V}^{\mathsf{T}}$ and the results of part (a) and part (b) to show that $\|\mathbf{M}\|_2 = \|\mathbf{S}\|_2$. **[1 mark]**

$$\|M\|_2 = \|USV^T\|_2 = \|SV^T\|_2 = \|S\|_2$$

$\boxed{1}$

(d) Compute $\|\mathbf{S}x\|_2^2$ where $x = (x_1, x_2, ..., x_n)$ and $\mathbf{S} = \operatorname{diag}(s_1, s_2, ..., s_n)$ is the diagonal matrix of singular values, $s_i$.

[1 mark]

$$\|SX\|_2^2 = \|X\|_2^2 = \sum_i (X_i)^2$$

1

(e) Write down the Lagrangian, $L$, to maximise $\|\mathbf{S}x\|_2^2$ subject to $\|x\|_2^2 = 1$. Compute the extremum conditions given by $\partial L/\partial x_i = 0$. Let $(s_\alpha | \alpha = 1, 2, ...)$ be the set of unique singular values and $I_\alpha$ the set of indices such that $s_i = s_\alpha$ if $I \in I_\alpha$. Using the extremum condition and the constraint, write down the set of extremum values for $\|\mathbf{S}x\|$ and hence show that $\|\mathbf{M}\|_2 = s_{max}$ where $s_{max}$ is the maximum singular value and $\mathbf{M} = \mathbf{USV}^\mathsf{T}$.

[4 marks]

To maximum $\|SX\|_2^2$ subject to $\|X\|_2^2 = 1$.

to Introduce a multiplier $\lambda$,

$$L(X, \lambda) = X^T S^T S^T X - \lambda (X^T X - 1)$$

$$\frac{dL(X,\lambda)}{dX} \Rightarrow 2\delta X^T (S^T S X - \lambda X) = 0$$

$$\frac{d(X,\lambda)}{d\lambda} \Rightarrow \delta \lambda (X^T X - 1) = 0$$

therefore, the stationary points satisfy:

$$S^T S X = \lambda X \quad \text{with} \quad \|X\|_2 = 1.$$

So $X$ needs to be an eigen vector with $S^T S$

So $X$ is the $S_{max}$ when $\|SX\|_2^2$ is maximum.

$$\|M\|_2 = \|USV^T\|_2 = \|S\|_2 = S_{max}.$$

4

End of question 3

| (a) | (b) | (c) | (d) | (e) | Total |
|-----|-----|-----|-----|-----|-------|
| — 2 | — 2 | — 1 | — 1 | — 4 | — 10 |

**4**

(a) We consider the mapping $y = \mathbf{M}x$ where $\mathbf{M}$ is an $n \times n$ matrix. Suppose there is some noise in $x$ so that $x' = x + \epsilon$ and under the mapping $y' = \mathbf{M}x'$. Compute an upper bound for $\|y' - y\|_2$ in terms of $\|\epsilon\|$ and $s_{max}$, where $s_{max}$ follows the same definition as in Q3(e). [2 marks]

$$\| y'-y \|_2 = \| Mx' - Mx \|_2 = \| M\epsilon \|_2$$

follow the Q3 (e) the upper bound is $s_{max} \| \epsilon \|_2$

$\boxed{2}$

(b) For a matrix $\mathbf{M} = \mathbf{USV}^\mathsf{T}$ show that

$$\|\mathbf{M}x\|_2 = \|\mathbf{S}a\|_2 \|x\|_2$$

where $a = \mathbf{V}^\mathsf{T}x / \|x\|_2$ so that $\|a\|_2 = 1$. Show that we can lower bound $\|\mathbf{S}a\|_2^2$ by $s_{min}^2$ and hence prove

$$\|\mathbf{M}x\|_2 \ge s_{min} \|x\|_2.$$

where $s_{min}$ is the minimum non-zero singular value analogous to the definition of $s_{max}$. [3 marks]

$$\| Mx \|_2 = \| USV^\mathsf{T}x \|_2 = \| SU^\mathsf{T}x \|_2 \quad \text{since } U \text{ is orthogonal}$$

$$\| SU^\mathsf{T}x \|_2 = \| S(U^\mathsf{T}x) \|_2$$

$$= \sum_i s_i^2 \| U^\mathsf{T}x \|_2 \ge s_{min} \| V^\mathsf{T}x \|_2$$

$$= s_{min} \| x \|_2.$$

$\boxed{3}$

(c) Using the previous results, obtain an upper bound for the relative error

$$\frac{\|y' - y\|_2}{\|y\|_2}$$

in terms of $s_{max}$, $s_{min}$, $\|\epsilon\|_2$ and $\|x\|_2$. [1 mark]

$$\frac{s_{max} \| \epsilon \|_2}{s_{min} \| x \|_2}$$

$\boxed{1}$

(d) The condition number for an invertible square matrix $\mathbf{M}$ is given by $\kappa_2(\mathbf{M}) = \|\mathbf{M}\|_2 \|\mathbf{M}^{-1}\|_2$ (there are different condition numbers for different norms.) Write down the condition number of $\mathbf{M}$ in terms of $s_{max}$ and $s_{min}$. [1 mark]
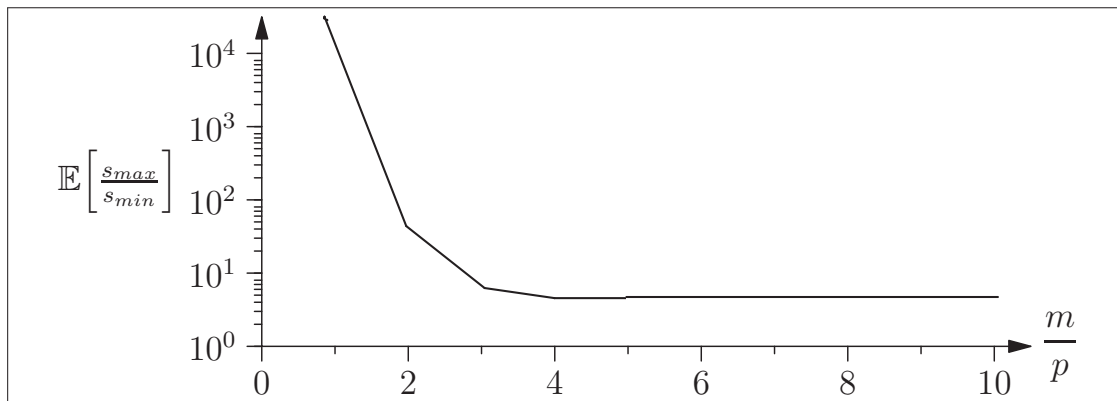
$$k_2(mn) = \|M\|_2 \, \|M^{-1}\|_2 = \frac{Smax}{Smin}$$

$\boxed{\overline{1}}$

(e) In linear regression we make predictions $\hat{y} = \boldsymbol{x}^\mathsf{T}\boldsymbol{w}$ given an input $\boldsymbol{x}$ where $\boldsymbol{w} = \mathbf{X}^+\boldsymbol{y}$, where $\mathbf{X}^+ = (\mathbf{X}^\mathsf{T}\mathbf{X})^{-1}\mathbf{X}^\mathsf{T}$ is the pseudo inverse of the design matrix $\mathbf{X}$ and $\boldsymbol{y}$ is a vector of training examples. There are bounds on the accuracy of linear regression depending on $\mathbb{E}\left[s_{max}/s_{min}\right]$ where $s_{max}$ and $s_{min}$ are respectively the maximum and minimum no-zero singular values of the design matrix. Consider randomly drawn feature vectors

$$\boldsymbol{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

Use python to generate the $m \times p$ dimensional design matrix $\mathbf{X}$ with rows $\boldsymbol{x}_i^\mathsf{T}$. By computing the singular values for $\mathbf{X}$ with $m = i \times p$ where $i \in \{1, 2, ..., 10\}$, find $s_{max}/s_{min}$. Repeat this 10 times for each $i$ to obtain an estimate of $\mathbb{E}\left[s_{max}/s_{min}\right]$. Plot a graph of your estimate for $\mathbb{E}\left[s_{max}/s_{min}\right]$ (on a log-axis) versus $m/p$ for $p = 10, 50$ and $100$. [3 marks]



$\boxed{\overline{3}}$

End of question 4

| (a) $\frac{}{2}$ | (b) $\frac{}{3}$ | (c) $\frac{}{1}$ | (d) $\frac{}{1}$ | (e) $\frac{}{3}$ | Total $\frac{}{10}$ |