
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

$$f(ax + (1 - a)y) \leq af(x) + (1 - a)f(y) \quad (1)$$

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 $\epsilon \rightarrow 0$ 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]

Rearranging the Equation 1

$$f(y + a(x - y)) \leq f(y) + a(f(x) - f(y)).$$

Or

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

Letting $a = \epsilon/(x - y)$ **then**

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

Taking the limit $\epsilon \rightarrow 0$ **then using**

$$\lim_{\epsilon \rightarrow 0} \frac{f(y + \epsilon) - f(y)}{\epsilon} = f'(y)$$

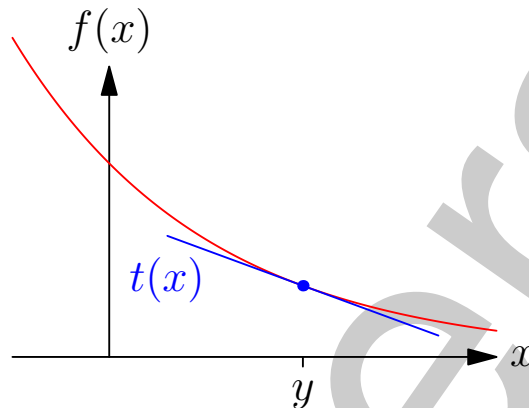
So that

$$(x - y)f'(y) \leq f(x) - f(y)$$

or

$$f(x) \geq f(y) + (x - y)f'(y) = t(x).$$

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



(c) Starting from the inequality for a convex function

$$f(x) \geq f(y) + (x - y)f'(y) \quad (2)$$

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 ϵ^2 show that for a convex function $f''(x) \geq 0$. [4 marks]

We use the expansions

$$\begin{aligned} f(x + \epsilon) &= f(x) + \epsilon f'(x) + \frac{\epsilon^2}{2} f''(x) + O(\epsilon^3) \\ f'(x + \epsilon) &= f'(x) + \epsilon f''(x) + O(\epsilon^2) \end{aligned}$$

Substituting into the Equation (2)

$$f(x) \geq f(x) + \epsilon f'(x) + \frac{\epsilon^2}{2} f''(x) + O(\epsilon^3) - \epsilon (f'(x) + \epsilon f''(x) + O(\epsilon^2))$$

Or subtraction $f(x)$ on both sides

$$0 \geq -\frac{\epsilon^2}{2} f''(x) + O(\epsilon^3)$$

Since this has to be true for all $\epsilon > 0$ this requires $f''(x) \geq 0$.

(d) Prove that x^4 is convex.

[1 marks]

$$\frac{d^2 x^4}{dx^2} = 12x^2 \geq 0$$

End of question 1

2

(a) Show by writing out in component for that $\text{tr} \mathbf{AB} = \text{tr} \mathbf{BA}$ where $\text{tr} \mathbf{M} = \sum_i M_{ii}$ (i.e. the trace of a matrix is equal to the sum of terms down the diagonal). [2 marks]

$$\text{tr } \mathbf{AB} = \sum_{i,j} A_{ij} B_{ji}$$

$$\text{tr } \mathbf{BA} = \sum_{k,l} B_{kl} A_{lk}$$

These are identical (which is clearer if we let $l = i$ and $k = j$).

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

$$\text{tr } \mathbf{M} = \text{tr } (\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T) = \text{tr } (\mathbf{V}^T\mathbf{V}\mathbf{\Lambda}) = \text{tr } \mathbf{\Lambda} = \sum_i \lambda_i$$

where we have used (1) $\text{tr } \mathbf{AB} = \text{tr } \mathbf{BA}$ with $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}$ and $\mathbf{B} = \mathbf{V}^T$ and (2) $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

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

$$\text{tr } \mathbf{X}^T\mathbf{X} = \text{tr } \begin{pmatrix} \mathbf{x}_1^T\mathbf{x}_1 & \mathbf{x}_1^T\mathbf{x}_2 & \dots & \mathbf{x}_1^T\mathbf{x}_n \\ \mathbf{x}_2^T\mathbf{x}_1 & \mathbf{x}_2^T\mathbf{x}_2 & \dots & \mathbf{x}_2^T\mathbf{x}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_n^T\mathbf{x}_1 & \mathbf{x}_n^T\mathbf{x}_2 & \dots & \mathbf{x}_n^T\mathbf{x}_n \end{pmatrix} = \sum_{i=1}^n \mathbf{x}_i^T\mathbf{x}_i = \sum_{i=1}^n \|\mathbf{x}_i\|^2$$

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

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

using the previous result show that $\|\mathbf{X}\|_F^2 = \text{tr } \mathbf{X}^T\mathbf{X}$ [2 marks]

$$\|\mathbf{X}\|_F^2 = \sum_{i,j} X_{ij}^2 = \sum_j \left(\sum_i X_{ij}^2 \right) = \sum_j \|\mathbf{x}_j\|^2 = \text{tr } \mathbf{X}^T\mathbf{X}$$

where we have used the fact that $\sum_i X_{ij}^2$ is a sum over all elements in the j^{th} column but the j^{th} column is equal to the vector \mathbf{x}_j .

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

Now

$$\mathbf{X}^T\mathbf{X} = (\mathbf{V}\mathbf{S}^T\mathbf{U}^T)(\mathbf{U}\mathbf{S}\mathbf{V}^T) = \mathbf{V}\mathbf{S}^T\mathbf{S}\mathbf{V}^T$$

And so using

$$\text{tr } \mathbf{X}^T \mathbf{X} = \text{tr } \mathbf{S}^T \mathbf{S}$$

but $\mathbf{S}^T \mathbf{S}$ is a diagonal matrix with elements s_i^2 so $\text{tr } \mathbf{X}^T \mathbf{X} = \sum_i s_i^2$. But we have shown that $\|\mathbf{X}\|_F^2 = \text{tr } \mathbf{X}^T \mathbf{X}$, so $\|\mathbf{X}\|_F^2 = \sum_i s_i^2$.

End of question 2

3 The p -norm of a matrix \mathbf{M} is defined to satisfy

$$\|\mathbf{M}\|_p = \max_{\mathbf{x} \neq 0} \frac{\|\mathbf{M}\mathbf{x}\|_p}{\|\mathbf{x}\|_p} \quad (3)$$

$$= \max_{\mathbf{x}: \|\mathbf{x}\|_p=1} \|\mathbf{M}\mathbf{x}\|_p \quad (4)$$

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

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

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

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

We note that

$$\|\mathbf{U}\mathbf{v}\|_2^2 = \mathbf{v}^T \mathbf{U}^T \mathbf{U} \mathbf{v} = \mathbf{v}^T \mathbf{v} = \|\mathbf{v}\|_2^2$$

since norms are positive $\|\mathbf{U}\mathbf{v}\|_2 = \|\mathbf{v}\|_2$. We note that $\|\mathbf{U}\mathbf{A}\mathbf{x}\|_2 = \|\mathbf{A}\mathbf{x}\|_2$, thus the unit vector \mathbf{x} that maximises $\|\mathbf{A}\mathbf{x}\|_2$ will also maximise $\|\mathbf{U}\mathbf{A}\mathbf{x}\|_2$. By Equation (4) this implies $\|\mathbf{U}\mathbf{A}\|_2 = \|\mathbf{A}\|_2$.

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

We consider a unit-vector \mathbf{x} that maximises $\|\mathbf{A}\mathbf{x}\|_2$ then by definition (4) $\|\mathbf{A}\mathbf{x}\|_2 = \|\mathbf{A}\|_2$. But we note that $\mathbf{y} = \mathbf{V}\mathbf{x}$ will be a unit vector that maximises $\|\mathbf{A}\mathbf{V}^T\mathbf{y}\|_2$ since

$$\|\mathbf{A}\mathbf{V}^T\mathbf{y}\|_2 = \|\mathbf{A}\mathbf{V}^T\mathbf{V}\mathbf{x}\|_2 = \|\mathbf{A}\mathbf{x}\|_2$$

using $\mathbf{V}^T \mathbf{V} = \mathbf{I}$. But as \mathbf{x} is a unit vector that maximises $\|\mathbf{A}\mathbf{x}\|_2$ then \mathbf{y} will be a unit vector that maximise $\|\mathbf{A}\mathbf{V}^T\mathbf{y}\|_2$ and $\|\mathbf{A}\mathbf{V}^T\|_2 = \|\mathbf{A}\|_2$.

- (c) Using the SVD $\mathbf{M} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ and the results of part (a) and part (b) show that $\|\mathbf{M}\|_2 = \|\mathbf{S}\|_2$. [1 marks]
-

We note that $\|\mathbf{M}\|_2 = \|\mathbf{U}\mathbf{S}\mathbf{V}^T\|_2$. But by part (a) $\|\mathbf{U}\mathbf{S}\mathbf{V}^T\|_2 = \|\mathbf{S}\mathbf{V}^T\|_2$ and by part (b) $\|\mathbf{S}\mathbf{V}^T\|_2 = \|\mathbf{S}\|_2$.

- (d) Compute $\|\mathbf{S}\mathbf{x}\|_2^2$ where $\mathbf{S} = \text{diag}(s_1, s_2, \dots, s_n)$ is the diagonal matrix of singular values [1 marks]

Since $\mathbf{S} = \text{diag}(s_1, s_2, \dots, s_n)$ then

$$\|\mathbf{S}\mathbf{x}\|_2^2 = \sum_{i=1}^n s_i^2 x_i^2$$

- (e) Write down the Lagrangian to maximise $\|\mathbf{S}\mathbf{x}\|_2^2$ subject to $\|\mathbf{x}\|_2^2 = 1$. Compute the extremum conditions given by $\partial L / \partial x_i = 0$. Let $(s_\alpha | \alpha = 1, 2, \dots)$ be the set of unique singular values and I_α the set of indices, i such that $s_i = s_\alpha$. Using the extremum condition and the constraint write down the set of extremum values for $\|\mathbf{S}\mathbf{x}\|$ and hence show that $\|\mathbf{M}\|_2 = s_{\max}$ where s_{\max} is the maximum singular value. [4 marks]

The Lagrangian is equal to

$$L = \sum_{i=1}^n s_i^2 x_i^2 - \lambda \left(\sum_{i=1}^n x_i^2 - 1 \right)$$

The extremum conditions are

$$\frac{\partial L}{\partial x_i} = x_i (s_i^2 - \lambda) = 0$$

(and the constraint). The extremum conditions are thus either $x_i = 0$ or $s_i^2 = \lambda$. Thus the value of Lagrange multipliers are $\lambda = s_\alpha^2$ where $x_i = 0$ if $i \notin I_\alpha$. Thus

$$\|\mathbf{S}\mathbf{x}\|_2^2 = s_\alpha^2 \sum_{i \in I_\alpha} x_i^2 = s_\alpha^2$$

where we have used $\sum_i x_i^2 = 1$. The maximum condition is just given by the case s_{\max} (the largest maximum eigenvalue). Thus using Equation (4) we have $\|\mathbf{S}\|_2 = s_{\max}$. But since $\|\mathbf{M}\|_2 = \|\mathbf{S}\|_2$ it follows that $\|\mathbf{M}\|_2 = s_{\max}$.

End of question 3

4

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

$$\|\mathbf{y}' - \mathbf{y}\|_2 = \|\mathbf{M}\boldsymbol{\epsilon}\|_2 \leq \|\mathbf{M}\|_2 \|\boldsymbol{\epsilon}\|_2 = s_{\max} \|\boldsymbol{\epsilon}\|_2$$

- (b) For a matrix $\mathbf{M} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ show that

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

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

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

[3 marks]

$$\|\mathbf{X}\mathbf{x}\|_2 = \|\mathbf{U}\mathbf{S}\mathbf{V}^T \mathbf{x}\|_2 = \|\mathbf{S}\mathbf{V}^T \mathbf{x}\|_2$$

but using $\mathbf{a} = \mathbf{V}^T \mathbf{x} / \|\mathbf{x}\|_2$

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

where $\|\mathbf{a}\|_2 = 1$. But

$$\|\mathbf{S}\mathbf{a}\|_2^2 = \sum_i a_i^2 s_i^2 \geq s_{\min}^2 \sum_i a_i^2 = s_{\min}^2 \|\mathbf{a}\|_2^2 = s_{\min}^2.$$

Thus $\|\mathbf{S}\mathbf{a}\|_2 \geq s_{\min}$ and

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

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

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

in terms of s_{\max} , s_{\min} , $\|\boldsymbol{\epsilon}\|_2$ and $\|\mathbf{x}\|_2$.

[1 marks]

$$\frac{\|\mathbf{y}' - \mathbf{y}\|_2}{\|\mathbf{y}\|_2} \leq \frac{s_{\max}}{s_{\min}} \frac{\|\boldsymbol{\epsilon}\|_2}{\|\mathbf{x}\|_2}$$

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

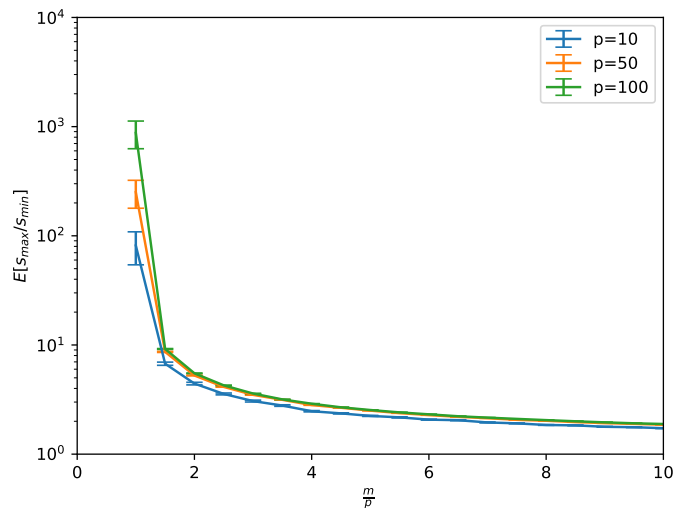
We observe that $\|\mathbf{A}\|_2 = s_{\max}$ and $\|\mathbf{A}^{-1}\|_2 = s_{\min}^{-1}$ so that

$$\kappa_2(\mathbf{M}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 = \frac{s_{\max}}{s_{\min}}$$

(e) In linear regression we make predictions $\hat{\mathbf{y}} = \mathbf{x}^T \mathbf{w}$ given an input \mathbf{x} where $\mathbf{w} = \mathbf{X}^+ \mathbf{y}$ where $\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is the pseudo inverse of the design matrix \mathbf{X} and \mathbf{y} is a vector of training examples. There are bounds on the accuracy of linear regression depending on $\mathbb{E}[s_{\max}/s_{\min}]$ where s_{\max} and s_{\min} are the maximum and minimum non-zero singular value of the design matrix. Consider randomly drawn feature vectors

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

Using python generate the $m \times p$ dimensional design matrix \mathbf{X} with rows x_i^\top . By computing the singular values for \mathbf{X} for $m = i \times p$ where $i = 1, 2, \dots, 10$ find s_{\max}/s_{\min} . Repeat this 10 times to obtain an estimate of $\mathbb{E}[s_{\max}/s_{\min}]$. Plot a graph of you estimate for $\mathbb{E}[s_{\max}/s_{\min}]$ (on a log-axis) versus m/p for $p = 10, 50$ and 100. [3 marks]



End of question 4

END OF PAPER