

## Past Papers unofficial solutions

While studying and unable to have the solutions for these papers, I decided to create this document which will (hopefully) contain all (unofficial) solutions from every question from every past paper in the MLPR course. There are only a few days till the exam, so I won't be writing answers in details; rather a sentence or a phrase to which you can use to expand your solution for that question.

These answers have been written by me or with some assistance from someone who knew the topic or material more than I do. Some questions contain multiple possible answers. It's important to note that **these answers are not official and not fact checked with a course organizer**, so mistakes are expected to exist. The answers will be updated and checked daily by myself and help from my peers for some answers, until the day of the examination. Take them with a grant of salt.

I spoke with Dr. Arno, and he likes the idea and even encouraged me (and everyone) to share past paper solutions with each other. Discussion with your peers is a great way of learning. Your peer may know something that you don't and vice versa.

If you have any suggestion or any edits I should make, please email me and I will immediately act on it.

Written by: George Karabassis.

Special thanks from every peer who has helped by coming up with answers from the past papers.

### 2021 Dec

1. A) *(Pick one solution)*
    - i. Gaussian Elimination -> epoch form -> # of Non 0 rows  $\Sigma(2)$ , so not full rank  $\therefore \Sigma(2)$  causes the problem.
    - ii. Cholesky Decomposition – `numpy.linalg.cholesky( $\Sigma$ )` ->  $\Sigma(2)$  throws an error.
  2. A) *(Any two solutions)*
    - i. **(possibly wrong)** L2 Regularization to the diagonal of the matrix – makes it positive definite.
    - ii. **(possibly wrong)** Use Partial Eigenvalue Decomposition on  $\Sigma(2) = Q\Lambda Q^T$ . replace negative eigenvalues with 0 -> reconstruct  $\Sigma(2)$ .
    - iii. Modify training data.
    - iv. Remove 3<sup>rd</sup> feature.
- B)

- i. **(mod required)** Function based on description (large values close to edges): sigmoid with range (3,5):  $a_{t+1} = 3 + \sigma(\eta \cdot \nabla_a c(a_t)) \cdot 2$ , with cost function  $c(a) = \frac{1}{2m} \sum_{i=1}^m (y_i - f(x_i; a))^2$  where the partial derivative is:  $\frac{\partial c}{\partial a} = -\frac{1}{m} \sum_{i=1}^m (y_i - f(x_i; a)) \frac{\partial f(x_i; a)}{\partial a}$ . Use a small learning rate,  $\eta = 0.001$ .
- ii. ☹
- iii. Apply Gaussian Naives optimization and then RMSE (check Assignment, Q5).

3. A)

- i. D features, N data points  $\rightarrow O(D \cdot N)$
- ii. Full Covariance  $\rightarrow O(N \cdot D^2)$

B)

- i. Pro: Faster computationally.  
Pro: Suitable for large databases, thanks to the use of mini batches.  
Con: Harder to implement.  
Con: Higher Computational Complexity (greater Big "Oh").

4. A)

- i. W: H x D (H hidden neurons, D dimensional layer x) order matters!
- ii. b: H x 1

B) Matrix A: H x K

Matrix B: K x D

Bias c : H

Extra parameters per country: H x K + K x D + H.

C)  $X @ (W + A @ B) . T + c [None, :]$

D) Less risk on overfitting, less likely to memorize training data, less memory consumption, more generative. (pick any two of these reasons)

Handwritten mathematical derivations for backpropagation through matrix multiplication. The forward pass is  $h = g((W+AB)x+c)$ . The backward pass starts with  $D = g'(D) \odot h$ . Then  $F = (W+AB)x$ . Since  $D = E + c \Rightarrow \bar{D} = \bar{E}$  and  $\bar{D} = \bar{c}$ , we have  $\bar{F} = \bar{W} + \bar{G}$ . Since  $F = W + AB$ , we have  $\bar{F} = \bar{W} + \bar{G}$ . Using  $C = AB \Rightarrow \bar{A} = \bar{C}\bar{B}$ ,  $\bar{B} = \bar{A}\bar{C}$ . Then  $\bar{E} = \bar{F}x \Rightarrow \bar{F} = \bar{E}x^T$ . Since  $F = W + G \Rightarrow \bar{F} = \bar{W}$ . Since  $G = AB \Rightarrow \bar{A} = \bar{C}\bar{B}$ ,  $\bar{B} = \bar{A}\bar{C}$ . Prepropagation:  $\bar{B} = \bar{A}^T g'((W+AB)x+c) \odot h \cdot x^T$ .

E)

F) Initialize AB to be a zero matrix by setting the matrix B to zero. Replace bias c to b. For the first epoch, the model will be back to its original form. The gradients from B will slowly be generated after training.

G) Early Stopping (to prevent overfitting).

5. A) ☹

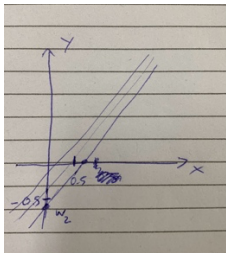
B) **(possible idea)** take  $k(x_1, x_2) = E[f_1 * f_2] - E[f_1] * E[f_2]$ , where  $f = a+bx+cx^2$ . No, we can't have a gaussian process  $k(x_1, x_2) < 0$ .

### *2021 Dec Comments*

Harder than expected. The exam was open book, so questions are expected to be harder and require higher levels of technical knowledge to answer them. Don't let that bring you down. 😊

2019 Dec

- i. A)  
 $\mathbb{E}_{p(\mathbf{x}, y)}[L(y, f(\mathbf{x}))]$  Where:  $L$  = loss,  $y$  = training result,  $f(\mathbf{x})$  = model function,  $p$  = true data distribution.  
**OR**  $\int L(y, f(\mathbf{x}))p(\mathbf{x}, y)d\mathbf{x}dy$  Where:  $L$  = loss,  $y$  = training result,  $f(\mathbf{x})$  = model function,  $p$  = true data distribution.  
**OR** The expected loss that the model (function) will achieve/yield/score on future/unseen test datapoints.
  - ii.  $\frac{1}{M} \sum_{m=1}^M L(y^{(m)}, f(\mathbf{x}^{(m)}))$
  - iii.  $\sqrt{\frac{1}{M} \sum_{m=1}^M (Lm - L)^2}$
  - B)
    - i. —
    - ii. —
  - C)
    - i. Generalization error doesn't predict risk to overfitting. Model could be prone to overfitting.
    - ii. Provide validation data sets to test the model and check if it overfits.
2. A)  $W: K \times D$ ,  $c: D \times 1$ ,  $v: 1 \times K$ ,  $b: 1 \times 1$ .  
B) Gradients will be zero and the model will not learn.  
C) Prevents from the gradients (derivatives) from growing too large, causing optimization issues.



3. A) i)  
3 lines close to the means on each axis, each line close to each other (small variance), straight lines (linear function). Not necessarily parallel, small angle should be considered, given the means.  
ii) Increase the variance, so there's less emphasis on the mean.
- B) i) We have prior knowledge, based on the parameters (mean and variance). Very efficient with less data, just like in our example.  
ii) Minimizing the expected square error, will essentially get closer to the mean of the samples. Similarly, the predictive distribution, will choose a predictive value closer to the mean (mode) of the gaussian. Therefore, both ideas estimate the  $y$  in a similar way.
- C) i)  $\alpha_1, \alpha_2, l_1, l_2$ .  $\alpha_1 > 0, \alpha_2 > 0, l_1 \neq 0, l_2 \neq 0$ .  
ii)  $\alpha_1, \alpha_2$ , give more emphasis on one kernel over another. The  $l$  (lengthscale) changes

the variation of the kernel.

Assume ~~two~~ hyperparameters exist:

$$k_{\text{comb}}(x^{(i)}, x^{(j)}) = a_1 \exp\left(-\frac{1}{2\sigma_1^2} \|x\|^2\right) + a_2 \exp\left(-\frac{1}{2\sigma_2^2} \|x\|^2\right)$$

(used  $X = (x^{(i)} - x^{(j)})^T (x^{(i)} - x^{(j)})$  to save time)

Take logs:

$$\log(k_{\text{comb}}(x^{(i)}, x^{(j)})) = \log(a_1) + \log\left(\exp\left(-\frac{1}{2\sigma_1^2} \|x\|^2\right)\right) + \log(a_2) + \log\left(\exp\left(-\frac{1}{2\sigma_2^2} \|x\|^2\right)\right)$$

$$= \log(a_1 a_2) + \left(\frac{1}{2\sigma_1^2} + \frac{1}{2\sigma_2^2}\right) \|x\|^2$$

Take logs on single:

$$\log(k_{\text{single}}(x^{(i)}, x^{(j)})) = \log(a_1) + \log\left(\exp\left(-\frac{1}{2\sigma_1^2} \|x\|^2\right)\right)$$

Equate comb & single:  $\log(a_1 a_2) = \log(a_1) \Rightarrow a_1 a_2 = a_1$

Answer: No.  $\frac{1}{2\sigma_1^2} = \frac{1}{2\sigma_1^2} + \frac{1}{2\sigma_2^2}$  only if  $\sigma_2 = \infty$  or  $\sigma_2 = 1$  ✓

iii)

4. A)  $B(\text{dash}) = A^T \times C(\text{dash}) \rightarrow 1^{\text{st}}$  Standard Result)  
 B)  $O(N \times K \times D^2)$   
 C) PCA, pro: doesn't require data labels, con: error when re-assembling data (some data loss).
5. A)  $p(y = 1 | x, D) \approx \sum_{s=1}^S s(w^{(s)T} x)$   
 B)  $p(y = 1 | x, D) = \int \sigma(w^T x) p(w | D) dw$   
 C) SVI, Laplace, Sampling, Monte Carlo, Sequential Monte Carlo (any two).

## 2019 Dec Comments

1)a) has been fact checked from the last exam preparation lecture.

## 2018 Dec

1. A) i) Converting categorical data to numerical values. One-hot encoding converts categorical data to binary values, where the selected category will be 1 and the rest 0.  
 ii) Assumes features are independent to each other, while in one-hot encoding, all features are dependent to each other.  
 iii) Neural Network, Logistic Regression with arg max (?).  
 iv) Apply logarithmic transformation on each salary feature. Simplifies the data.

## FAQs

### **Q. How do I come up with the solutions?**

**A.**

1. Attempt the whole past paper on my own under exam-like conditions.
2. Check my answers through the notes from the website.
3. Do a **lot** of research to confirm my answers are correct or if they are wrong, spend a lot of time to understand why and what could've done better.
4. Ask my peers. Most likely they will know the answers partially or give me suggestions or ideas.
5. Write the mathematical answers in a concise form and include them here.
6. Pass the answers from paper to Word.

### **Q. How long does it take you to come up with the answers?**

**A.** A lot! No coincidence why I study for about 16 hours a day. Completing a past paper takes me a good 2 hour period. 1 hour checking my answers and fixing my mistakes, 1 extra hour researching, 2 hours extra to cover a topic I feel very weak on, 1 hour to pass my answers from the paper to the Word document.

### **Q. Some mathematical format is wrong. Why is that?**

**A.** No experience with writing Mathematical Equations in LaTeX! That's the main reason why some answers are written by hand. Saves time. 😊