

Note:

- During the attendance check a sticker containing a unique code will be put on this exam.
- This code contains a unique number that associates this exam with your registration number.
- This number is printed both next to the code and to the signature field in the attendance check list.

Mining Massive Datasets

Exam: IN2323 / Endterm

Date: Friday 9th August, 2019

Examiner: Prof. Dr. Stephan Günnemann

Time: 13:30 – 15:00

	P 1	P 2	P 3	P 4	P 5	P 6	P 7
I							

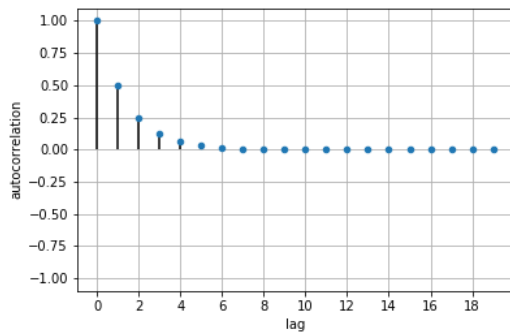
Working instructions

- This exam consists of **8 pages** with a total of **7 problems**.
Please make sure that you received a complete copy of the exam.
- You can earn 38 points.
- **Detaching pages from the exam is prohibited!**
- Allowed resources:
 - A4 sheet of handwritten notes (two sides)
 - **no other materials (e.g. books, cell phones, calculators) are allowed!**
- Only write on the sheets given to you by supervisors. If you need more paper, ask the supervisors.
- Last two pages can be used as scratch paper.
- All sheets (including scratch paper) have to be returned at the end.
- **Only use a black or a blue pen (no pencils, red or green pens)!**
- Write your answers only in the provided solution boxes or the scratch paper.
- **For problems that say "Justify your answer" or "Show your work" you only get points if you provide a valid explanation.** Otherwise it's sufficient to only provide the correct answer.
- Exam duration - 90 minutes.

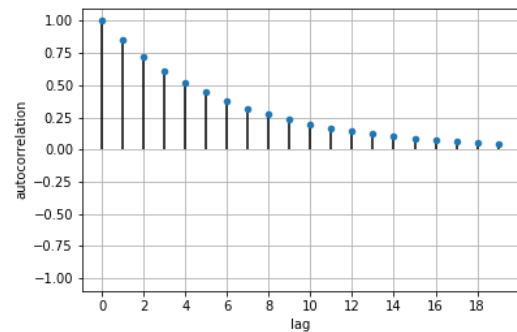
Left room from _____ to _____ / Early submission at _____

Problem 1 AR models: correlation function (4 points)

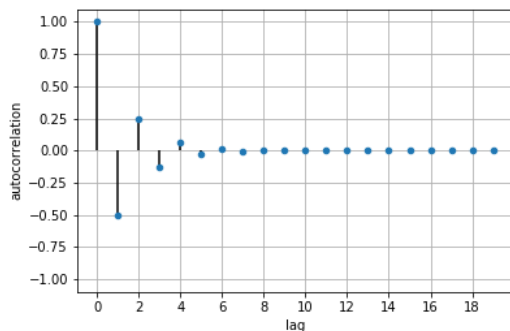
For each of the AR(1) and AR(2) models below find the corresponding autocorrelation function plot from Figure 1.1. Each plot was generated using one of the AR models from the list so that there is a one-to-one correspondence between the plots (1)-(4) and the AR processes (a)-(d). Everywhere $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ with a positive σ .



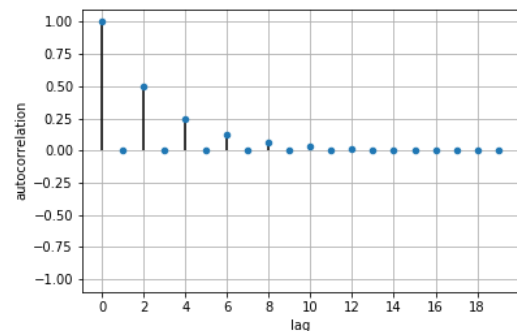
(1)



(2)



(3)



(4)

Figure 1.1: AR models: correlation function

a) $\mathcal{X}_t = -0.5\mathcal{X}_{t-1} + \epsilon_t$

☐ (1) ☐ (2) ☐ (3) ☐ (4)

b) $\mathcal{X}_t = 0.85\mathcal{X}_{t-1} + \epsilon_t$

☐ (1) ☐ (2) ☐ (3) ☐ (4)

c) $\mathcal{X}_t = 0.5\mathcal{X}_{t-2} + \epsilon_t$

☐ (1) ☐ (2) ☐ (3) ☐ (4)

d) $\mathcal{X}_t = 0.5\mathcal{X}_{t-1} + \epsilon_t$

☐ (1) ☐ (2) ☐ (3) ☐ (4)

Problem 2 Hidden Markov Models (6 points)

Consider the following Hidden Markov Model where Z_t are latent variables and X_t are continuous observed variables. We parametrize the prior and transition probabilities $P(Z_1 = i) = \pi_i$ and $P(Z_{t+1} = j | Z_t = i) = A_{ij}$ by:

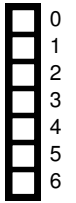
$$\pi = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}, \quad A = \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix}$$

We parametrize the emission probabilities with $X_t | Z_t = 0 \sim \text{Uniform}([8, 12])$ and $X_t | Z_t = 1 \sim \text{Uniform}([11, 13])$. We assume we observed $X = [8.5, 11, 13, 9.5]$.

Perform the Forward algorithm, compute α_t and compute the most probable state Z_t given X_1, \dots, X_t at every step t .

Hint 1: You can check your computations by comparing with $\alpha_4 = \begin{bmatrix} 3/3200 \\ 0 \end{bmatrix}$

Hint 2: Maintaining non-decimal fractions makes calculations easier (e.g. use $8/10$ instead of 0.8).



Problem 3 Recurrent neural networks (8 points)

Given a sequence of integers $\{x_1, \dots, x_n\}$, $x_t \in \{0, 1\}$, the goal is to output 1 at step t if the sum of all the inputs up to step t is even. That is, $h_t = 1$ if $\sum_{i=0}^t x_i$ is even, otherwise 0. For example, the outputs for a sequence 01101 are 10110. To model this you use an RNN with the following update equations:

$$\begin{aligned}\tilde{h}_t &= a \cdot h_{t-1} + b \cdot x_t \\ z_t &= f(c \cdot h_{t-1} + d \cdot x_t) \\ h_t &= f((1 - z_t) \cdot h_{t-1} + z_t \cdot \tilde{h}_t)\end{aligned}$$

where $h_0 = 1$ and $a, b, c, d \in \mathbb{R}$ and $f(x) = 1$ if $x > 0.5$, otherwise 0.

- 0 ☐ a) Fill out the table with all the possible input values x_t and input hidden state values h_{t-1} , and calculate the corresponding output h_t according to the given requirements.

1 ☐

2 ☐

h_{t-1}	x_t	h_t

- 0 ☐ b) Find the values for a, b, c and d such that this model achieves 100% accuracy. Show your work.

1 ☐

2 ☐

3 ☐

4 ☐

5 ☐

6 ☐

Hint: think about what can happen with a new input and what should gate z_t do in those cases.

Problem 4 Graph laws (4 points)

You are given an Erdős-Rényi graph $G(n, p)$, where n is the number of nodes and p is the probability of an edge.

a) What is the expected number of cliques of size m ?

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2

b) We add a new node and connect it to every existing node with probability q . What is the expected number of newly created triangles?

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2

Problem 5 Deep Generative Models (6 points)

The loss used in generative adversarial networks (GANs) can be written in the following form:

$$\min_{\theta} \max_{\phi} \mathcal{L}(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{p^*(\mathbf{x})}[\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{p(\mathbf{z})}[\log(1 - D_{\phi}(f_{\theta}(\mathbf{z})))]$$

where $p^*(\mathbf{x})$ is the true data distribution, $p(\mathbf{z})$ is the distribution of the noise, f_{θ} is the generator, and D_{ϕ} is the discriminator.

- 0 ☐ a) For a given generator (fixed parameters θ) assume there exists a discriminator $D_{\phi^*}(\mathbf{x})$ with parameters ϕ^*
 1 ☐ such that for all \mathbf{x} :

$$D_{\phi^*}(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + p_{\theta}(\mathbf{x})}$$

where $p_{\theta}(\mathbf{x})$ is the distribution learned by the generator. Show that D_{ϕ^*} is **optimal**, i.e. $\phi^* = \arg \max_{\phi} \mathcal{L}(\theta, \phi)$.
 Hint: $\max_y [a \log(y) + b \log(1 - y)] = \frac{a}{a+b}$ for any $a, b \in \mathbb{R}_0^+$, $a + b > 0$.

- 0 ☐ b) Show that $\mathcal{L}(\theta, \phi^*) = -\log(4) + 2 \cdot \text{JSD}(p^*(\mathbf{x}) || p_{\theta}(\mathbf{x}))$ where
 1 ☐
 2 ☐
 3 ☐

$$\text{JSD}(p^*(\mathbf{x}) || p_{\theta}(\mathbf{x})) = \frac{1}{2} \left[\text{KL}(p^*(\mathbf{x}) || m(\mathbf{x})) + \text{KL}(p_{\theta}(\mathbf{x}) || m(\mathbf{x})) \right]$$

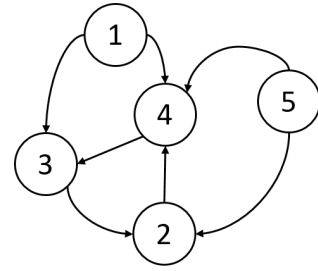
is the Jensen–Shannon divergence, KL is the Kullback–Leibler divergence, and $m(\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}) + p^*(\mathbf{x})}{2}$.

Problem 6 Ranking (6 points)

Given the following graph G let $e = (s, t)$ be a new edge, i.e. an edge which is **not** yet in the graph G . We denote with G_{new} the resulting graph when adding e to G . Your task is to identify the edge e based on the following observation:

The PageRank vector of G_{new} based on a teleport set $S = \{1, 5\}$ and $\beta = 0.85$ (i.e. the probability to teleport is 0.15) is given by

$$\pi = [0.0750, 0.1935, 0.2668, 0.2764, 0.1884]$$



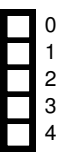
What is the edge e (specify the source node s and target node t) that has been added? Justify your answer.

Problem 7 Spectral clustering (4 points)

You are given a graph G with $N = 1000$ nodes that is generated from a Stochastic Block Model with the following parameters

$$\pi = \begin{bmatrix} 0.5 & 0.5 \end{bmatrix} \quad \eta = \begin{bmatrix} 0.05 & 0.5 \\ 0.5 & 0.05 \end{bmatrix}$$

Assume that G is connected. Do you expect spectral clustering to recover the true communities \mathbf{z} when applied to the graph G (yes or no)? Justify your answer.



Additional space for solutions—clearly mark the (sub)problem your answers are related to and strike out invalid solutions.

