

Name :  
 Study program :  
 ID. NR. :

1. For each of the following sub-questions, you are asked to provide a *short but essential* answer. You should not need more than three sentences per answer.

- a. Consider a binary classification problem with two classes  $\{y_1, y_2\}$  and input vector  $x$ . We are given a data set to train the parameters  $\theta$  for a likelihood model of the form

$$p(y_k = 1|x, \theta) = \frac{1}{1 + e^{-\theta_k^T x}}$$

There are two fundamentally different ways to train  $\theta$ , namely through a generative model or by discriminative training.

- (1) Explain shortly how we train  $\theta$  through a generative model. No need to work out all equations for Gaussian models, but explain the strategy in probabilistic modeling terms.
  - (2) Explain shortly how we train  $\theta$  through a discriminative approach.
- b. Explain shortly how Bayes rule relates to machine learning. In your answer, you may assume a model  $\mathcal{M}$  with prior distribution  $p(\mathcal{M})$  and an observed data set  $D$ .
- c. What is the difference between supervised and unsupervised learning? Express the goals of these two learning methods in terms of a probability distribution. (I'm looking here for a statement such as: "Given ..., the goals of supervised/unsupervised learning is to estimate  $p(\cdot|\cdot)$ ".)
- d. In a particular model with hidden variables, the log-likelihood can be worked out to the following expression:

$$L(\theta; D) = \sum_n \log \left( \sum_k \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right)$$

Do you prefer a gradient descent or EM algorithm to estimate maximum likelihood values for the parameters? Explain your answer. (No need to work out the equations. )

- e. The maximum likelihood estimate (MLE) of the class-conditional mean in a classification problem can be expressed as

$$\hat{\mu}_k = \frac{\sum_n y_n^k x_n}{\sum_n y_n^k}$$

and the M-step update for the cluster mean in a clustering problem is given by

$$\hat{\mu}_k = \frac{\sum_n \gamma_n^k x_n}{\sum_n \gamma_n^k}$$

Explain the relation between  $y_n^k$  and  $\gamma_n^k$ . Is  $y_n^k$  a binary variable? And what about  $\gamma_n^k$ ?

2. The lifetime  $x > 0$  of a light bulb is postulated to be exponentially distributed with unknown mean  $\mu > 0$ , i.e.

$$p(x|\mu) = \frac{1}{\mu} e^{-x/\mu}$$

In order to estimate  $\mu$ , the lifetimes  $\mathbf{X} = \{x_1, \dots, x_N\}$  of  $N$  independent bulbs are observed.

- a. Work out the log-likelihood  $\log p(\mathbf{X}|\mu)$ .
- b. What is the maximum likelihood estimate for  $\mu$  based on observations  $\mathbf{X}$ ?

In a separate experiment  $M$  independent bulbs were tested, but the individual lifetimes were not recorded. We will use the symbols  $\mathbf{Z} = \{z_1, \dots, z_M\}$  for the *unobserved* lifetimes of bulbs  $N + 1, \dots, N + M$  and  $\{\mathbf{X}, \mathbf{Z}\} = \{x_1, \dots, x_N, z_1, \dots, z_M\}$  for the complete data set. Instead of lifetimes, we only recorded if a bulb had failed at time  $t$ . We record  $y_m = 1$  if bulb  $N + m$  still burned at time  $t$  and  $y_m = 0$  if the bulb had already failed at time  $t$ . We will now derive an EM algorithm to estimate  $\mu$ , based all  $N + M$  observations.

- c. Complete the following two formula's for the EM algorithm:

$$\begin{aligned} \mathbf{E}\text{-step} : & \quad \text{evaluate } p(\cdot | \cdot, \mu^{\text{old}}) \\ \mathbf{M}\text{-step} : & \quad \mu^{\text{new}} = \arg \max \sum p(\cdot | \cdot, \cdot) \log p(\cdot, \cdot | \cdot) \end{aligned}$$

- d. Proof that the expected complete-data log-likelihood  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mu)]$  equals

$$-(N + M) \log \mu - \frac{1}{\mu} \left( N\bar{x} + \sum_{m=1}^M \mathbb{E}[z_m] \right)$$

where  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$ .

- e. You can find the (local) optimum of  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mu)]$  by setting its derivative w.r.t.  $\mu$  to zero. Now differentiate  $\mathbb{E}[\log p(\mathbf{X}, \mathbf{Z} | \mu)]$  (see answer 2d) w.r.t.  $\mu$  and set to zero to obtain the re-estimation formula (**M-step**) for  $\mu$ .

We do not derive the  $\mathbb{E}[z_m]$  for the **E-step**. Use the following equation instead

$$\mathbb{E}[z_m] = \begin{cases} t + \mu & \text{if } y_m = 1 \\ \mu - \frac{t \exp(-\frac{t}{\mu})}{1 - \exp(-\frac{t}{\mu})} & \text{if } y_m = 0 \end{cases}$$

- f. In total we found that  $r$  out of  $M$  bulbs had failed at time  $t$ . Derive an expression for  $\sum_{m=1}^M \mathbb{E}[z_m]$  in terms of  $r$  and  $M$ .
- g. Put the results of the last two exercises together and derive the re-estimation formula (**M-step**) for  $\mu$  (in terms of a previous estimate of  $\mu^{\text{old}}$ ).

3. Let  $B$  be a positive real valued random variable with probability density

$$p_B(b) = e^{-b}, \quad \text{for all } b > 0.$$

Also  $A$  is a real valued random variable with conditional density

$$p_{A|B}(a|b) = \sqrt{\frac{b}{\pi}} e^{-a^2 b}, \quad \text{for all } a \in (-\infty, \infty) \text{ and } b \in (0, \infty).$$

- a. Give an (integral) expression for  $p_A(a)$ .  
Do not try to evaluate the integral.
- b. Approximate  $p_A(a)$  using the Laplace approximation.  
Give the detailed derivation, not just the answer.

4. The *Bayesian Information Criterion* is results in

$$\underbrace{\log \frac{p(\mathcal{M}_1 | x^N)}{p(\mathcal{M}_2 | x^N)}}_{(*1)} \approx \underbrace{\log \frac{p(\mathcal{M}_1)}{p(\mathcal{M}_2)}}_{(*2)} + \underbrace{\log \frac{p(x^N | \mathcal{M}_1, \hat{\theta}_1)}{p(x^N | \mathcal{M}_2, \hat{\theta}_2)}}_{(*3)} + \underbrace{\frac{1}{2} (k_1 - k_2) \log N}_{(*4)}.$$

Here  $x^N$  is a binary data sequence of length  $N$ ,  $k_1$  and  $k_2$  are the number of free parameters in respectively model  $\mathcal{M}_1$  and  $\mathcal{M}_2$ , and  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are the estimated (ML) parameter vectors.

- a. Explain the four terms marked by (\*1), (\*2), (\*3), and (\*4).
- b. The binary data  $x^N = x_1, x_2, \dots, x_N$  is generated by a Bernoulli process, i.e.

$$p(x^N | \mathcal{M}, \theta) = (1 - \theta)^{n(0|x^N)} \theta^{n(1|x^N)}.$$

The parameter prior  $p(\theta | \mathcal{M})$  is given by the Beta distribution:

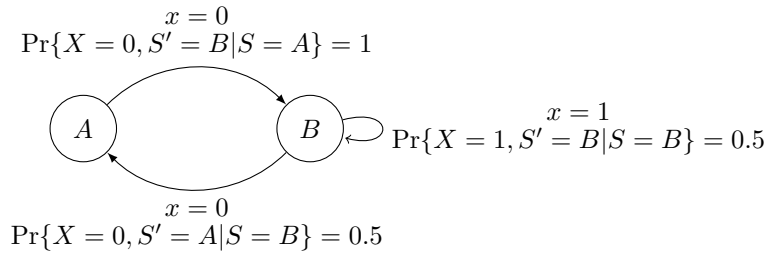
$$p(\theta | \mathcal{M}) = \frac{1}{\pi} \frac{1}{\sqrt{\theta(1-\theta)}}.$$

Let  $N = 10$  and  $x^{10} = 1001101101$ . Determine  $p(x^N | \mathcal{M})$ . Give the complete derivation starting with the information given above.

- c. Why do we think that the probability estimate  $p(x^N | \mathcal{M})$  is a useful and good estimate for the actual, but unknown, probability  $p(x^N | \mathcal{M}, \theta)$ ? And how close will  $p(x^N | \mathcal{M})$  be to the probability  $p(x^N | \mathcal{M}, \theta)$ , for any  $x^N$  and any  $\theta$ , if  $\mathcal{M}$  is a binary memoryless model?
5. Consider the following binary finite state model (Markov source). This model produces outputs  $X_t$  where the probability of the next output symbol depends on the current state of the source. We list all non-zero probabilities.

$$\begin{aligned} \Pr\{X_t = 0, S_{t+1} = B | S_t = A\} &= 1, \\ \Pr\{X_t = 0, S_{t+1} = A | S_t = B\} &= 0.5, \\ \Pr\{X_t = 1, S_{t+1} = B | S_t = B\} &= 0.5, \end{aligned}$$

The following figure depicts this model.



- a. Compute the stationary probabilities  $q(A)$  and  $q(B)$  where

$$q(s) = \lim_{t \rightarrow \infty} \Pr\{S_t = s\} \quad \text{for } s \in \{A, B\}.$$

- b. Compute the following probabilities assuming that the model is stationary (i.e.  $\Pr\{S_1 = A\} = q(A)$  and  $\Pr\{S_1 = B\} = q(B)$ ).

$$\begin{aligned} \Pr\{X_1 = 1\} \\ \Pr\{X_2 = 1 | X_1 = 0\} \end{aligned}$$

- c. Let  $\mathcal{M}_0$  be the i.i.d. model with

$$\underline{\theta}_0 = (\Pr\{X_1 = 1\}).$$

Also  $\mathcal{M}_1$  is the first order model with

$$\underline{\theta}_1 = (\theta_{10}, \theta_{11}) = (\Pr\{X_2 = 1 | X_1 = 0\}, \Pr\{X_2 = 1 | X_1 = 1\}).$$

And  $\mathcal{M}_2$  is the second order model with

$$\begin{aligned} \underline{\theta}_2 &= (\theta_{200}, \theta_{201}, \theta_{210}, \theta_{211}) \\ &= (\Pr\{X_3 = 1 | X_1 = 0, X_2 = 0\}, \Pr\{X_3 = 1 | X_1 = 0, X_2 = 1\}, \\ &\quad \Pr\{X_3 = 1 | X_1 = 1, X_2 = 0\}, \Pr\{X_3 = 1 | X_1 = 1, X_2 = 1\}). \end{aligned}$$

The Markov model produces a ‘typical’ sequence so

$$-\log_2 \Pr\{X^n = x^n | \mathcal{M}_i, \theta_i\} \approx H_i(X^n),$$

where  $H_i(X^n)$  is the entropy rate of the  $i^{\text{th}}$  model.

Given is that

$$H_0(X^n) = 0.9183 \cdot n$$

$$H_1(X^n) = 0.8742 \cdot n$$

$$H_2(X^n) = 0.7925 \cdot n$$

**Determine for what range of  $n$  you should use  $\mathcal{M}_0$ . And when  $\mathcal{M}_1$  and when  $\mathcal{M}_2$ ?**

Use the idea of stochastic complexity and **motivate your answer**.

---

Points that can be scored per question:

Question 1: each sub-question a through e: 2 points. Total 10 points.

Question 2: a) 2 points; b) 1 point; c) 2 points; d) 2 points; e) 1 point; f) 1 point; g) 1 point.  
Total 10 points.

Question 3: a) 1 point; b) 5 points. Total 6 points.

Question 4: a) 4 points; b) 2 points; c) 2 points. Total 8 points.

Question 5: a) 1 point; b) 2 points; c) 3 points. Total 6 points.

Max score that can be obtained: 40 points.

The final grade is obtained by dividing the score by 4 and rounding to the nearest integer.