

# ML1 Main Exam 2022

52041MAL6Y Machine Learning 1 22/23 (Period 1) · 13 exercises · 40.0 points

## 1 MC: Evil Likelihood Model

1.0 point · 1 question

Suppose you wish to devise an "evil likelihood model", which will make your model likelihood  $p(\mathbf{x}|\mathbf{w})$  as low as possible. In this case, we wish to find the minimum likelihood solution for the model parameters  $\mathbf{w}$ . Which statements are true?

1.0 point · Multiple choice · 4 alternatives

- ☒ The goal of such model is to maximize the function:  $\frac{1}{p(\mathbf{x}|\mathbf{w})}$ .
- ☒ The goal of such model is to maximize the function:  $-\log p(\mathbf{x}|\mathbf{w})$ .
- ☐ In the case of linear regression, the minimum likelihood estimate is equal to the maximum likelihood solution, but with opposite signs.
- ☐ In the case of a linear regression, we can find the unique solution to this problem.

Feedback

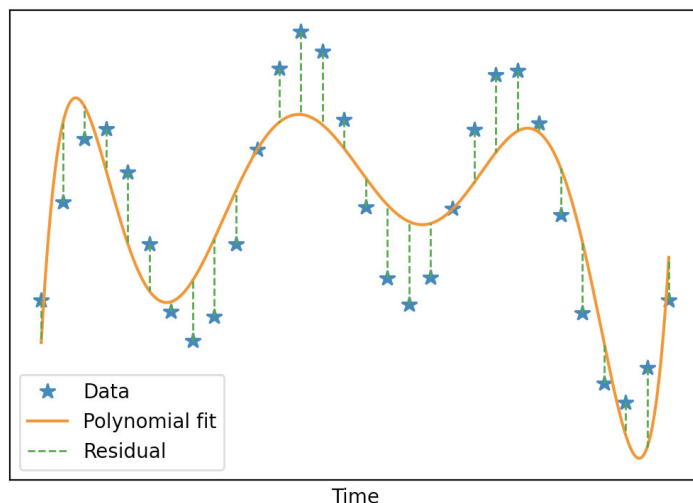
Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

## 2 MC: FLAC compression

1.0 point · 1 question



Audio is stored as sequence of measurements of the waveform at discrete time intervals. FLAC is a way of compressing audio by least-squares fitting a polynomial and storing the residual. Weights are stored at full precision, while the residual is compressed further. Which statements are true?

1.0 point · Multiple choice · 4 alternatives

- ☐ Adding L2 regularisation can improve compression performance.
- ☒ The weights for the polynomial fit have a closed-form solution.
- ☐ The audio data is i.i.d.
- ☒ Increasing the order  $M$  of the polynomial causes the residual to shrink

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

### 3 MC: I.i.d., Conditional Independence and GPs

1.0 point · 1 question

Let  $\mathcal{D} = \{(x_i, t_i)\}_{i=1}^N$  with samples given by  $t_i = f(x_i)$  with  $f \sim GP(m(\cdot), k(\cdot, \cdot))$  a random function according to a Gaussian process with mean function  $m$  and kernel  $k$ . Which statements are true?

1.0 point · Multiple choice · 6 alternatives

- ☐  $t_i$  is not a random variable.
- ☐  $t_i \sim p(t)$  is i.i.d. relative to some  $p(t)$ .
- ☒  $t_i \sim p(t|x_i)$  is i.i.d. relative to some  $p(t|x)$

Feedback

Yes, namely  $p(t|x) = N(t|m(x), k(x, x))$ .

- ☐  $\text{Cov}[t_i, t_j] = 0$  for any  $i \neq j$
- ☐  $\text{Cov}[t_i, t_j] = k(t_i, t_j)$
- ☒  $\text{Cov}[t_i, t_j] = k(x_i, x_j)$

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

## 4 MC: Valid Kernels

1.0 point · 1 question

Which of the following kernels are valid? Let  $x, x' \in \mathbb{R}^d$  be two vectors of the same dimensionality  $d$ .

1.0 point · Multiple choice · 6 alternatives

- ☐  $k(x, x') = \min(x, x')$ , for  $x, x' \in \mathbb{R}$
- ☒  $k(x, x') = 1$
- ☒  $k(x, x') = (1 + x \cdot x')^2$
- ☒  $k(x, x') = 1 + x \cdot x'$
- ☒  $k(x, x') = \exp(x + x')$ , for  $x, x' \in \mathbb{R}$
- ☒  $k(x, x') = x \cdot x'$

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

## 5 MC: Overfitting and model complexity

1.0 point · 1 question

Which of the following statements are true? Check all that apply

1.0 point · Multiple choice · 4 alternatives

- ☐ Higher complexity models are more prone to overfitting and typically have lower variance
- ☒ Only adding more data for training a learner with high bias may not reduce the test error.
- ☒ Overfitting may arise when relevant features are missing in the data
- ☐ Increasing the depth of a neural network will always reduce the test error.

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

## 6 MC: Two clusters

1.0 point · 1 question

We have the following dataset, where the samples belong to two different classes  $C1, C2$ . Which of the following statements are true? (Note:  $I_2$  denotes the  $2 \times 2$  identity matrix.)

1.0 point · Multiple choice · 6 alternatives

- ☐ K-means with 2 means may fit the data well.
- ☒ A Gaussian Mixture Model with 2 Gaussian components may fit the data well.
- ☐ The conditional distribution for  $C1$  can be accurately modeled with a Gaussian  $\mathcal{N}(\mu_1, s \cdot I_2)$  for some mean vector  $\mu_1 \in \mathbb{R}^2$  and some scalar  $s > 0$ .
- ☒ The conditional distribution for  $C2$  can be accurately modeled with a Gaussian  $\mathcal{N}(\mu_2, s \cdot I_2)$  for some mean vector  $\mu_2 \in \mathbb{R}^2$  and some scalar  $s > 0$ .
- ☒ The conditional distribution for  $C1$  can be accurately modeled with a Gaussian  $\mathcal{N}(\mu_1, \Sigma_1)$  for some mean vector  $\mu_1 \in \mathbb{R}^2$  and some covariance matrix  $\Sigma_1$ .
- ☒ The conditional distribution for  $C2$  can be accurately modeled with a Gaussian  $\mathcal{N}(\mu_2, \Sigma_2)$  for some mean vector  $\mu_2 \in \mathbb{R}^2$  and some covariance matrix  $\Sigma_2$ .

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

## 7 MC: Probabilistic models

1.0 point · 1 question

In classification we consider three models: discriminant functions, probabilistic generative models and probabilistic discriminative models. The following are statements about probabilistic generative models and probabilistic discriminative models. Which are true?

1.0 point · Multiple choice · 4 alternatives

- ☒ Logistic regression is a probabilistic discriminative model.
- ☒ In probabilistic discriminative models, the posterior probabilities  $p(C|\mathbf{x})$  are modeled directly.
- ☐ In probabilistic discriminative models, the prior probability of class  $p(C)$  is modeled.
- ☒ In generative models, the class conditional probability  $p(\mathbf{x}|C)$  are modeled.

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly



**8 MC: SVM**

1.0 point · 1 question

Consider the following SVM optimization problem. Which statements are true?

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\lambda} \sum_{n=1}^N \xi, \quad \text{s.t.} \quad \begin{cases} \forall_n : t_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi \\ \forall_n : \xi \geq 0 \end{cases}$$

1.0 point · Multiple choice · 4 alternatives

- ☐ For large  $\lambda$  we expect a more complex decision boundary than for small  $\lambda$ .
- ☒ For large  $\lambda$  we expect a less complex decision boundary than for small  $\lambda$ .
- ☒ For large  $\lambda$  we expect more support vectors than for small  $\lambda$ .
- ☐ For large  $\lambda$  we expect less support vectors than for small  $\lambda$ .

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

## 9 MC: Neural Networks

1.0 point · 1 question

Consider a neural network with  $L = 5$  layers. Let us denote  $w_{ij}^{(l)}$  the weights in each layer, the number of features (the width) in each hidden layer with  $M$ , the used hidden activation functions with  $h$ , and the error of the model with respect to the  $n^{th}$  data point with  $E_n$ . Which statements are true?

1.0 point · Multiple choice · 5 alternatives

- ☐ Updating the weights in layer  $l = 2$  requires to perform the forward pass only up to layer 2.

Feedback

False, you need to compute the full forward pass.

- ☐ Updating the weights in layer  $l = 3$  requires computation of the backward pass down to all layers.

Feedback

No, we don't need to backprop to layer 1 and 2 if we want to update the weights of layer 3

- ☐ In order for back-propagation to work the network is not allowed to contain skip connections.

Feedback

False, the NN should be feed forward, but is allowed to have skip connections

- ☒ Optimizing a neural network with stochastic gradient descent means updating the weights via  $w_{ij}^{(l)} = w_{ij}^{(l)} - \eta \frac{\partial E_n}{\partial w_{ij}^{(l)}}$ , with  $\eta$  a hyperparameter.

Feedback

Correct

- ☐ Using  $h(a) = a^2$ , the neural network can in theory represent any function  $\phi(x)$  up to arbitrary precision by scaling up  $M$ .

Feedback

Feedback when the question is answered correctly

Feedback when the question is answered partially correctly

Feedback when the question is answered incorrectly

## 10 Gamma Distribution

5.5 points · 3 questions

In linear regression we assume that  $y = \phi(\mathbf{x})^T \mathbf{x} + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$  for all our datapoints. We saw that this formulation could equivalently be expressed as  $y$  being random according to a (predictive) distribution

$$p(y \mid \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(\phi(\mathbf{x})^T \mathbf{x}, \beta^{-1}).$$

In such a model, the mean parameter of the Normal distribution is thus modeled by a linear model  $\mu(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$ .

In general, we can expect the observations to follow different type of distributions, depending on the type of noise or the type of values target value can take on. E.g., many stochastic processes can not have negative outcomes (e.g. rainfall, waiting times, loan defaults), in which case regressing with a normal distribution would be undesirable!

For example, when modeling waiting times  $y$  in stores, given some input features  $\mathbf{x}$ , we know we are predicting a quantity that is always positive as one cannot have negative waiting times! An appropriate distribution for such random variables is the gamma distribution. We say that a random variable  $y > 0$  is gamma-distributed with shape  $\alpha$  and rate  $\beta$ , denoted as  $y \sim \text{Gamma}(\alpha, \beta)$ , if

$$p(y \mid \alpha, \beta) = \frac{y^{\alpha-1} e^{-\beta y} \beta^\alpha}{\Gamma(\alpha)},$$

where  $\Gamma$  denotes the gamma function (whose explicit form we do not need). Let us have a look at optimizing the parameters of the gamma distribution.

Text

a Suppose we observe a dataset  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$  and want to directly model a distribution for  $y_n$ , regardless of  $\mathbf{x}_n$ . I.e. we aim to find a single gamma distribution that models all  $y_n$ . Assume  $\alpha$  to be given. Give the Maximum Likelihood (ML) estimate of  $\beta$  in terms of  $\alpha$  and the data.

*Note that, given our modeling assumptions, the solution will not depend on  $\mathbf{x}_n$ .*

2.5 points · Open · 9/10 Page

**+0.5 points**

Correctly specify the likelihood

**+0.5 points**

Correctly specify the log-likelihood (give points to previous item if log-likelihood is directly given, which is fine)

**+0.5 points**

Provide objective (either as argmax or set derivative to zero)

**+0.5 points**

Correctly compute derivative

**+0.5 points**

Correctly solve for  $\beta$

b Suppose we have the following prior distribution of the  $\beta$  parameter:  $\beta \sim \text{Gamma}(a, b)$

. Assume hyperparameters  $\alpha, a, b$  all to be known. Give the Maximum A-Posteriori (MAP) estimate for  $\beta$  in terms of the data and the known parameters.

3.0 points · Open · 9/10 Page

**+1 point**

Correctly specify objective (either explicitly via argmax with posterior, or as likelihood x prior, or in terms of derivative of likelihood x prior = 0)

**+0.5 points**

Correctly compute log of the prior

**+0.75 points**

Correctly compute derivative of the log-prior

**+0.75 points**

Correctly solve for beta

c **[BONUS]** Suppose now we want our  $\beta$  parameter to depend on some input feature vector  $\phi_n := \phi(\mathbf{x}_n)$ . Specifically, we want to model the following predictive distribution:

$$p(y_n | \mathbf{x}_n, \mathbf{w}, \alpha) = \text{Gamma}(y_n | \alpha, \cosh(\phi_n^T \mathbf{w})).$$

For example, a chain of stores wants to model the customer waiting times as a function of location, time of day and other parameters. Then,  $y_n$  correspond to waiting time of the  $n$ -th observed customer, while  $\phi_n$  would be a vector with information about store location, time of day, etc.

As machine learners, we will make use of our good friend *gradient descent/ascent*. Provide the gradient-based update for model parameters  $\mathbf{w}$  with the aim of maximizing the log-likelihood.

*Hint: make use of the property that  $\frac{d}{dx} \cosh(x) = \sinh(x)$ .*

2.0 points · Bonus · Open · 4/5 Page

### +0.75 points

Correctly specify a gradient descent step. This can be:

- (i) SGD on a single datapoint likelihood,
- (ii) mini-batch SGD on K datapoints, or
- (iii) GD on the dataset likelihood.

(No points awarded for just writing down the generic GD update rule)

### +0.5 points

Correctly compute the likelihood needed for the gradient descent step.

(This can be any of SGD, mini-batch SGD, and GD, as long as it matches the gradient step)

### +0.75 points

Correctly compute the derivative of the likelihood needed for the gradient descent step.

### -0.3 points

For doing gradient descent on the log-likelihood, instead of gradient ascent

## 11 Pet Detective

6.0 points · 5 questions

Consider yourself to be a pet detective highly specialized in determining the species of odd looking pets of the "cat" and "dog" variety. In your work, clients come to you with pets of which they are uncertain about their species. Your approach is to collect appearance characteristics (furriness, color, weight, etc.) which you collect in a numeric vector  $\mathbf{x} \in \mathbb{R}^d$ . Your approach is based on probability theory, and you consider both  $\mathbf{x}$  and  $c \in \{\text{cat}, \text{dog}\}$  random variables according to some joint distribution  $\mathbf{x}, c \sim p(\mathbf{x}, c)$ .

Text

a You figured that you can best determine pet species  $c$  based on the probability for that class given the appearance vector  $\mathbf{x}$ . You know that any posterior for binary random variables can be written in the form  $p(c = \text{dog} \mid \mathbf{x}) = f(a(\mathbf{x}))$ , with  $f(a) = \frac{1}{1 + \exp(-a)}$ . What is the name for this function  $f$ ?

1.0 point · Open · 1/10 Page

**+1 point**

Logistic sigmoid (sigmoid is also fine)

b Give (or derive) the expression for  $a(\mathbf{x})$  in terms of the joint distribution  $p(\mathbf{x}, c)$ .

1.0 point · Open · 7/20 Page

**+1 point**

$a(\mathbf{x}) = \log( p(\mathbf{x}, c=\text{dog})/p(\mathbf{x}, c=\text{cat}) )$

**+0.75 points**

$a(\mathbf{x}) = -\log \left( \frac{p(\mathbf{x})}{p(\mathbf{x}, c=\text{dog})} - 1 \right)$

**-0.25 points**

Small mistake

**+0.25 points**

$p(\mathbf{x}, c) = p(c|\mathbf{x}) p(\mathbf{x})$

c What are the function values  $a(\mathbf{x})$  called, and why?

1.0 point · Open · 3/10 Page

**+0.5 points**

Logits, or log odds

**+0.5 points**

Because they give the log of the odds (ratio) for the dog class over the cat class

**+0.25 points**

(Incorrect, but partial credit) Activation functions

**+0.25 points**

Explanation for activation functions

d When modeling joint  $p(\mathbf{x}, c)$  with a Gaussian Mixture Model under some assumptions, the function  $a(\mathbf{x})$  takes on the shape of a linear function  $a(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ . The posterior then thus becomes a generalized linear model! Under what assumptions does this happen?

1.0 point · Open · 1/4 Page

**+1 point**

When  $\Sigma_1 = \Sigma_2$  with  $\Sigma_1$  and  $\Sigma_2$  the covariances of the two class conditional distributions.

e Suppose you are not interested in the joint  $p(\mathbf{x}, c)$ , but want to directly model the posterior with a generalized linear model using a database of solved cases  $\mathcal{D} = \{(\mathbf{x}_n, c_n)\}_{n=1}^N$ . You intend to find the best model parameters  $\mathbf{w}$  and  $b$  by defining an appropriate loss function and optimize via gradient descent. At the same time you would like to learn which of the measurements in  $\mathbf{x}$  are most important, and adapt the loss such that this becomes possible. We want to do feature selection as reducing the number of features could save you time in future investigations! Which loss should you minimize? (Answer in words or equations are both fine)

Are there any hyperparameters to tune?

2.0 points · Open · 1/2 Page

**+0.67 points**

Given the probabilistic model, the default loss is the binary cross-entropy loss

**+0.67 points**

This can be augmented with an L1 loss, which can be used to sparsify the weights and thus do feature selection.

**+0.67 points**

The parameter in front of the L1 loss, higher means more sparsification.

**+0.5 points**

alternative answer: For mentioning PCA, though the question specified how the loss function could be adapted, not an entirely different method.

**+0.33 points**

alternative answer: For giving L2 instead of L1 loss

**+0.67 points**

For mentioning some relevant hyperparameter



## 12 Modeling Muons

8.0 points · 5 questions

Muons are elementary particles, similar to electrons. On earth, muons constantly enter the atmosphere from space and decay into electrons and neutrinos. This time it takes for a muon to decay will be denoted with  $x$ , and it follows an exponential distribution

$$\text{Exp}(x|\tau) = \tau e^{-\tau x}.$$

I.e., the decay time  $x$  of a muon is random and follows the above distribution specified by half-life time  $\tau$ . You are a physics student who wants to measure the half-life  $\tau$  of a muon. In order to do this, you've bought a secondhand detector that measures decay time  $x$  of incoming muons and you've let it running all night to get many measurements.

But oh no! The detector is broken. Sometimes it works fine, but other times, it returns random noise. The figure below shows the situation. On the left is your measurement. On the right is what you think might have happened.

Since you studied ML1, you decide to model this as a mixture distribution. You assume that the faulty measurements are normally distributed according to

$$\mathcal{N}(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}.$$

Let  $\mathcal{D} = \{x_i\}_{i=1}^N$  denote the collected dataset.

Text

a Write down the log-likelihood of the data. You can give your answer in terms of  $\mathcal{N}(x|\mu, \sigma)$  and  $\text{Exp}(x|\tau)$  to save yourself some writing. Also, use  $\pi_1$  to denote the probability for receiving a correct measurement, and  $\pi_2$  for the probability of a faulty measurement.

1.5 points · Open · 2/5 Page

**+1 point**

For the correct answer

$$\log p(\mathcal{D}|\pi_1, \pi_2, \tau, \mu, \sigma) = \sum_{i=1}^N \log \left( \pi_1 \text{Exp}(x_i|\tau) + \pi_2 \mathcal{N}(x_i|\mu, \sigma) \right)$$

**+0.5 points**

For correct use of i.i.d. assumption, or mention of it

**-0.25 points**

For mixing up  $\pi_1$  and  $\pi_2$

**-0.5 points**

For forgetting about the log, or writing it in incorrect places

b Write down an expression for the posterior probability that a data point  $x_n$  was generated by a faulty measurement as well as the posterior probability that it was created by a good measurement. Again, you can write it in terms of  $\mathcal{N}(x|\mu, \sigma)$  and  $\text{Exp}(x|\tau)$ .

1.5 points · Open · 2/5 Page

**+0.75 points**

p(faulty|x\_n) correct

$$p(\text{faulty}|x_n) = \frac{\pi_2 \mathcal{N}(x_n|\mu, \sigma)}{\pi_1 \text{Exp}(x_n|\tau) + \pi_2 \mathcal{N}(x_n|\mu, \sigma)}$$

**+0.75 points**

p(correct|x\_n) correct

$$p(\text{correct}|x_n) = \frac{\pi_1 \text{Exp}(x_n|\tau)}{\pi_1 \text{Exp}(x_n|\tau) + \pi_2 \mathcal{N}(x_n|\mu, \sigma)}$$

**+0.5 points**

Correct application of Bayes rule (written for either scenario)

$$p(\text{faulty}|x_n) = \frac{p(x_n|\text{faulty}) p(\text{faulty})}{p(x_n)}$$

c Based on the obtained probabilistic model, you decide to throw away a data point if it is more likely to be noise than to be true data and compute the conditions for when you throw away a measurement  $x_n$ . Write the expression in the form of a quadratic inequality i.e.  $ax_n^2 + bx_n \geq c$ .  
*Hint: note that a decision boundary based on positive quantities does not change when applying a log on both sides. i.e.,  $p_1 \geq p_2 \Leftrightarrow \log p_1 \geq \log p_2$ .*

2.0 points · Open · 4/5 Page

### +0.5 points

For correct specification of the decision rule in terms of the posteriors

$$p(\text{faulty}|x_n) \geq p(\text{good}|x_n)$$

or, equivalently:

$$p(\text{faulty}|x_n) \geq 0.5$$

### +1.5 points

For the computation:

$$\begin{aligned} p(\text{faulty}|x_n) &\geq p(\text{good}|x_n) \\ \pi_2 \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x_n-\mu}{\sigma}\right)^2} &\geq \pi_1 \tau e^{-\tau x_n} \\ \log\left(\frac{\pi_2}{\sqrt{2\pi}\sigma}\right) - \frac{1}{2\sigma^2}(x_n - \mu)^2 &\geq \log(\pi_1 \tau) - \tau x_n \\ \tau x_n - \frac{1}{2\sigma^2}(x_n - \mu)^2 &\geq \log\left(\frac{\sqrt{2\pi}\pi_1\sigma\tau}{\pi_2}\right) \\ -\frac{1}{2\sigma^2}x_n^2 + \left(\tau + \frac{\mu}{\sigma^2}\right)x_n &\geq \left(\frac{\sqrt{2\pi}\pi_1\sigma\tau}{\pi_2}\right) + \frac{\mu^2}{2\sigma^2} \end{aligned}$$

### -0.5 points

If a mistake is made along the way, but the general steps are ok

d The mixture model can be optimized via the Expectation Maximization (EM) algorithm. Derive the M-step equation for the muon half-life  $\tau$ . Write it in terms of the responsibilities for the faulty class (found in sub question b), which you may denote with the symbol  $\gamma$ .

2.0 points · Open · 4/5 Page

**+0.5 points**

Give objective (maximize log likelihood w.r.t. tau)

$$\max_{\tau} \log p(\mathcal{D}) = \max_{\tau} \sum_{i=1}^N \log p(x_i)$$

**+0.5 points**

Correct computation of the derivative

$$\frac{\partial}{\partial \tau} \log p(\mathcal{D}) = \sum_{i=1}^N \frac{1}{p(x_i)} \pi_1 \tau \exp(-\tau x_i) \left( \frac{1}{\tau} - x_i \right)$$

**+0.5 points**

Correct replacement of the posterior probability from part (b) with gamma, or at least identifying the posterior

$$\gamma_i = (\pi_1 \tau \exp(-\tau x_i)) / p(x_i)$$

**+0.5 points**

Correct final solution in terms of the gamma

$$\tau = \left( \sum_{i=1}^N \gamma_i \right) / \left( \sum_{i=1}^N \gamma_i x_i \right)$$

e Let's assume you successfully derived and applied the EM algorithm to obtain optimal values for  $\pi_1, \pi_2, \tau, \mu$ , and  $\sigma$ . Give an expression for the percentage of data you expect to throw away in terms of the relevant model parameters.

1.0 point · Open · 1/4 Page

**+1 point**

Correct answer:  $\pi_2 * 100$

**-0.5 points**

For mixing up  $\pi_1$  with  $\pi_2$

### 13 Polar Coordinate SVM

11.5 points · 8 questions

**[Problem setting]** We have a dataset  $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$  where  $\mathbf{x}_n \in \mathbb{R}^2$  and  $t_n \in \{-1, +1\}$ . The data is centered around the origin and we expect the data to be almost perfectly separable by a decision boundary with a shape similar (up to scaling) to a curve  $M$ . See figure below in which the blue points correspond to datapoints for which  $t_n = -1$  and the red points correspond to datapoints for which  $t_n = +1$ .

Text

**[Derivation of the maximum margin objective]** The boundary  $M$  is parameterized by a continuous function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , where  $f(\mathbf{x}) \in \mathbb{R}$  indicates the distance from the origin to the boundary in the direction of the vector  $\mathbf{x} \in \mathbb{R}^2$  (see Figure above). The boundary  $M$  is then given by the following set of points

$$M = \left\{ f \left( \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \right) \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \in \mathbb{R}^2 : \theta \in [0, 2\pi) \right\}.$$

Text

We want to find a scale  $\sigma \geq 0$  such that  $\sigma M$  separates our dataset. For simplicity, we measure the distance of a point from the decision boundary in the radial direction, i.e. the signed distance  $d(\mathbf{x}, \sigma)$  of a point  $\mathbf{x}$  from the boundary  $\sigma M$  is given by

$$d(\mathbf{x}, \sigma) = \|\mathbf{x}\|_2 - \sigma f(x).$$

Text

We can use  $d(\mathbf{x}, \sigma)$  to define a decision boundary and assign label  $t_n = +1$  if  $d(\mathbf{x}_n, \sigma) \geq 0$  and  $t_n = -1$  otherwise. Two observations are important. Firstly, we have for all correct classifications

$$t_n d(\mathbf{x}_n, \sigma) \geq 0.$$

Text

Secondly, the decision boundary does not change when scaling points by a factor  $\alpha$ , i.e.,  $d(\alpha \mathbf{x}, \sigma) = 0 \Leftrightarrow \alpha d(\mathbf{x}, \sigma) = 0$ . This leads to arbitrariness when we want to define a margin. We can get rid of this arbitrary scaling by introducing a variable  $\alpha$  which we use to scale the signed distances such that the margin (smallest scaled distance) has value 1. We then have  $t_n \alpha d(\mathbf{x}_n, \sigma) \geq 1$  for all  $n$ , which fully written out gives

$$\forall n, \quad \alpha t_n (\|\mathbf{x}_n\|_2 - \sigma f(\mathbf{x}_n)) \geq 1.$$

Finally, it will be convenient in our derivations later on to apply the change of variable  $\beta = \alpha \sigma$  such that

$$\forall n, \quad t_n (\alpha \|\mathbf{x}_n\|_2 - \beta f(\mathbf{x}_n)) \geq 1.$$

Text

**[The objective]** We want to maximize the original distances of closest points on the margin to the decision boundary, given by  $d(\mathbf{x}_n, \sigma)$ . Given the constraint  $t_n \alpha d(\mathbf{x}_n, \sigma) \geq 1$ , maximizing  $d(\mathbf{x}_n, \sigma)$  implies that we want to minimize  $\alpha$  if we want to keep the margin at 1. Hence, we will consider the following equivalent problem:

$$\min_{\alpha} \quad \frac{1}{2} \alpha^2, \quad \text{subject to} \quad \begin{cases} t_n (\alpha \|\mathbf{x}_n\|_2 - \beta f(\mathbf{x}_n)) & \geq 1, \\ \alpha \beta & \geq 0. \end{cases}$$

Note that we replaced the original condition that  $\sigma = \beta/\alpha \geq 0$  with  $\alpha \beta \geq 0$ , which enforces  $\alpha$  and  $\beta$  to agree on their signs and, therefore,  $\sigma$  to be non-negative.

Text

a Suppose you have solved the optimization and found the optimal values of  $\alpha$  and  $\beta$ . What is the size of the margin?

1.5 points · Open · New page · 4/5 Page

**+0.5 points**

Solution determined by point on the margin, which satisfies  $t_n \alpha (\|\mathbf{x}\| - \sigma f(\mathbf{x})) = 1$

**+0.5 points**

Substitute  $t_n = +1$

**+0.5 points**

Obtain  $d(x_n, \sigma) = \frac{1}{\alpha}$

**+1.5 points**

Argue that the scaled (by  $\alpha$ ) margin will be 1 after optimization, hence the margin will be  $\frac{1}{\alpha}$

b Unfortunately, the points are not exactly separable. To allow for some error in the classification, introduce the slack variables  $\{\xi_n\}_{n=1}^N$  and a penalty  $C$  for the misclassified points. State the final optimization problem (and explicitly enumerate all the constraints):

1.0 point · Open · 21/50 Page

**+1 point**

For correct modification of the original problem

$$\min \frac{1}{2} \alpha^2 + C \sum_{n=1}^N \xi_n$$

,  
subject to

$$t_n (\alpha \|x_n\|_2 - \beta f(x_n)) \geq 1 - \xi_n$$

$$\alpha \beta \geq 0$$

$$\xi_n \geq 0$$

.

with  $n = 1, \dots, N$ .

**-0.25 points**

If the  $n = 1, \dots, N$  is not specified

**-0.25 points**

For minor mistakes (missing the slack variables in the optimization problem).

c Write down the primal Lagrangian. Use the following Lagrange multipliers for each constraint listed above:  $\{\lambda_n\}_{n=1}^N$  for the constraints on the margins;  $\{\mu_n\}_{n=1}^N$  for those on  $\xi_n$ ; and  $\gamma$  for the one on  $\alpha\beta$ .

*Indicate which variables are the primal variables and which ones are the dual variables.*

2.0 points · Open · 9/20 Page

**+0.75 points**

For the right terms in the Lagrangian:

$$L = \frac{1}{2}\alpha^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \lambda_n (t_n(\alpha||x_n|| - \beta f(x_n)) - 1 + \xi_n) - \sum_{n=1}^N \mu_n \xi_n - \gamma \alpha \beta$$

**+0.75 points**

If the sign of the Lagrange multipliers are correct

**+0.5 points**

For mentioning which variables are primals, and which duals

**-0.5 points**

If  $\gamma\alpha\beta$  is inside a sum over  $n$ , if notation is ambiguous be lenient



d Write down all of the Karush-Kuhn-Tucker (KKT) conditions. How many KKT conditions do we have in total? (Here we do not consider stationarity as part of the KKT conditions, see next question 13e)

2.0 points · Open · 3/5 Page

**+0.67 points**

Depending on how solution is given, points for 1/3 of the constraints correctly handled, or 1/3th of type of KKT is correct (primal feasibility, dual feasibility, complementary slackness)

**+0.67 points**

Depending on how solution is given, points for 1/3 of the constraints correctly handled, or 1/3th of type of KKT is correct (primal feasibility, dual feasibility, complementary slackness)

**+0.67 points**

Depending on how solution is given, points for 1/3 of the constraints correctly handled, or 1/3th of type of KKT is correct (primal feasibility, dual feasibility, complementary slackness)

**-0.67 points**

If not all constraints are enumerated (with  $n$ )

or if  $\gamma$  is also enumerated with  $n$ ...

**+2 points**

KKT conditions:

$$\lambda_n \geq 0$$

$$t_n(\alpha \|x_n\| - \beta f(x_n)) - 1 + \xi_n \geq 0$$

$$\lambda_n t_n(\alpha \|x_n\| - \beta f(x_n)) - 1 + \xi_n \geq 0$$

$$\xi_n \geq 0$$

$$\mu_n \geq 0$$

$$\xi_n \mu_n \geq 0$$

$$\alpha \beta \geq 0$$

$$\gamma \geq 0$$

$$\gamma \alpha \beta \geq 0$$

The total amount of constraints are:  $6N+3$ .

**-0.5 points**

For forgetting to mention the amount of constraints.

e Optimize the primal Lagrangian with respect to the primal variables. I.e, derive the stationarity conditions. *(answer box continues on the next page)*

2.0 points · Open · 1 3/10 Page

### +1 point

Correctly specify the objective (take the derivative set to zero).

Set the following derivatives equal to zero:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 0$$

### +1 point

Correctly perform derivation and rearrange:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = 0 \Leftrightarrow \alpha = \sum_{n=1}^N \lambda_n t_n ||x_n|| + \gamma \beta$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = 0 \Leftrightarrow 0 = \sum_{n=1}^N \lambda_n t_n f(x_n) - \gamma \alpha$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \Leftrightarrow C = \mu_n + \lambda_n$$

f The identities derived from the stationarity conditions can be used to derive expressions for  $\gamma$ . In our derivations we would have to consider three cases:

1. The case  $\beta = 0$  and any  $\gamma \geq 0$
2. The case  $\beta > 0$  and  $\gamma > 0$
3. The case  $\beta > 0$  and  $\gamma = 0$

Using the KKT conditions we can interpret these results. Which of these three cases gives us a sensible classifier, and why can we discard the other two as degenerate solutions?

2.0 points · Open · 1/2 Page

**+0.67 points**

Correct argumentation for why 1. is not sensible:  $\beta = 0$  implies that  $\sigma = \beta/\alpha = 0$  and thus every point gets classified as +1.

**+0.67 points**

Correct argumentation for why 2. is not sensible:  $\beta, \gamma > 0$  implies (due to complementary slackness) that  $\alpha = 0$ . This implies that  $\sigma = \beta/\alpha = \infty$  and thus all points lie within the boundary and are classified as -1.

**+0.67 points**

Correct conclusion is 3.

g Given the stationarity conditions, it is possible to show that one derives the dual Lagrangian as follows

$$\hat{L}(\{\lambda_n\}) = -\frac{1}{2} \left( \sum_{n=1}^N \lambda_n t_n \|x_n\|_2 \right)^2 + \sum_{n=1}^N \lambda_n \quad \text{with constraints } \forall n : \begin{cases} 0 \leq \lambda_n \leq C \\ \sum_{n=1}^N \lambda_n t_n f(\mathbf{x}_n) = 0 \end{cases}$$

Give the expression for the kernel function  $k(\mathbf{x}_n, \mathbf{x}_m)$  that is effectively used in the above problem.

1.0 point · Open · 1/5 Page

**+1 point**

For correct answer  $k(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m\| \|\mathbf{x}_n\|$

h **[BONUS]** We note that our kernel does not depend on the shape our boundary  $M$ , i.e., it does not depend on the function  $f$ . Why is this to be expected? How does the boundary  $M$  still influence the solution?

1.0 point · Bonus · Open · 1/2 Page

**+0.5 points**

The kernel defines a notion of similarity between points. This notion should not depend on the specific problem specification.

**+0.5 points**

$M$  still influences the constraints of the dual problem, and thus influences the solution.