

Prof. Joachim M. Buhmann

Examination Questions

January 26th, 2016

First and Last name: _____

ETH number: _____

Signature: _____

General Remarks

- The examination lasts for two hours. There are five sections.
- Please write your answers directly on the exam sheets. At the end of the exam you will find supplementary sheets, feel free to separate them from the exam. If you submit the supplementary sheets, put your name and ETH number on top of each.
- Please answer the questions in English. Do not use a pencil or red color pen.
- You may provide at most one valid answer per question. Invalid solutions must be clearly canceled out.

We wish you great success for this examination!

	Topic	Max. Points	Points	Signature
1	ML & Bayesian inference	20		
2	Kernels	20		
3	Neural Networks	20		
4	Gaussian processes	20		
5	Unsupervised learning	20		
Total		100		

Grade:

This page has been intentionally left blank.

Question 1: Maximum Likelihood and Bayesian Learning (20 pts.)

1. Given are n independently distributed binary vectors, each of dimension d .
 They are distributed according to:

$$X_{ij} \sim \text{ber}(1 - \theta), \quad \lambda_j = 1, \\ X_{ij} \sim \text{ber}(\theta), \quad \lambda_j = 0,$$

for $1 \leq j \leq d$ and $1 \leq i \leq n$ where λ_j selects the Bernoulli source.

- (a) What is the likelihood function $p(\mathcal{X}|\theta, \lambda)$?

$$\begin{aligned} p(\mathcal{X}|\theta, \lambda) &= \prod_{i,j} p(X_{ij}|\theta, \lambda) \\ &= \prod_{i,j} (\theta)^{2\lambda_j X_{ij} + 1 - \lambda_j - X_{ij}} (1-\theta)^{X_{ij} - 2\lambda_j X_{ij} + \lambda_j} \\ &= (1-\theta)^{\sum_{i,j} (2\lambda_j X_{ij} + 1 - \lambda_j - X_{ij})} \theta^{\sum_{i,j} (X_{ij} - 2\lambda_j X_{ij} + \lambda_j)} \end{aligned}$$

3 pts.

- (b) Derive the maximum likelihood estimate for θ given λ .

$$\begin{aligned} \frac{\partial \log p(\mathcal{X}|\theta, \lambda)}{\partial \theta} \\ = \frac{\partial}{\partial \theta} \left[\sum_{i,j} \left(\frac{(2\lambda_j X_{ij} + 1 - \lambda_j - X_{ij}) \log(1-\theta) + (X_{ij} - 2\lambda_j X_{ij} + \lambda_j) \log \theta}{\theta} \right) \right] \\ = - \frac{\sum_{i,j} (2\lambda_j X_{ij} + 1 - \lambda_j - X_{ij})}{1-\theta} + \frac{\sum_{i,j} (X_{ij} - 2\lambda_j X_{ij} + \lambda_j)}{\theta} \\ = 0 \\ \Rightarrow \theta = \frac{\sum_{i,j} (X_{ij} - 2\lambda_j X_{ij} + \lambda_j)}{n} = \frac{\sum_{j \in d} (n_j - 2n_i \lambda_j + n \lambda_j)}{n} \end{aligned}$$

4 pts.

$$\text{let } \sum_{i \in h} X_{ij} = h_j$$

- (c) Derive the posterior distribution $p(\theta|\mathcal{X}, \lambda, \alpha)$ given that the parameter θ is beta distributed: $p(\theta|\alpha) = B(\alpha, \alpha)^{-1}\theta^{(\alpha-1)}(1-\theta)^{(\alpha-1)}$, where $B(\alpha, \alpha)$ normalizes the distribution. Please determine the normalization of the posterior explicitly.

$$\begin{aligned}
 p(\theta|\mathcal{X}, \lambda, \alpha) &\propto \frac{p(x|\theta, \lambda) p(\theta|\alpha)}{p(x|\lambda, \alpha)} \\
 &\propto \frac{\beta(\lambda)}{(1-\theta)^{\lambda} \theta^{\lambda}} \frac{\theta^{\lambda-1} (1-\theta)^{\lambda-1}}{p(x|\lambda, \alpha)} \\
 &= \frac{\beta(\lambda)^{-1} (1-\theta)^{\lambda} \theta^{\lambda-1} [\sum_j (2\lambda x_{ij} + 1 - \lambda) - \lambda] + \lambda - 1}{p(x|\lambda, \alpha)}
 \end{aligned}$$

3 pts.

- (d) Derive the evidence $p(\mathcal{X}|\lambda, \alpha)$. Please determine the normalization of the evidence explicitly.

$$p(x|\lambda, \alpha) = \frac{\int_{\theta=0}^1 \beta(\lambda, \alpha)^{-1} \theta^{\lambda} (1-\theta)^{\alpha} d\theta}{\beta(\lambda, \alpha)^{-1} (1-\theta)^{\lambda} \theta^{\lambda} [\sum_j (2\lambda x_{ij} + 1 - \lambda) - \lambda] + \lambda - 1}$$

3 pts.

- (e) Under what condition is the maximum likelihood estimator for θ equal to the maximum a posterior estimator?

$$\text{MLE} = \frac{\sum_j (n_j - 2\lambda_j \lambda_j + n \lambda_j) + \lambda - 1}{n + 2\lambda - 2} = \frac{\sum_j (n_j - 2\lambda_j \lambda_j + n \lambda_j)}{n}$$

$$\Rightarrow \cancel{\lambda} \quad \lambda = 1 \quad \text{prior is uniform}$$

or when $\lambda \rightarrow \infty$ 1 pts.

2. Given is an input vector $X = (X_1, \dots, X_d)^T$ and model parameters $\beta \in \mathbb{R}^d$. A linear regression model predicts the response Y :

$$Y = X^T \beta + \epsilon, \quad \epsilon \sim N(\epsilon | 0, \sigma^2),$$

- (a) What is the likelihood function for n i.i.d. observations of Y ?

$$\text{Likelihood Function: } L = \frac{1}{(2\pi)^{\frac{n}{2}} \sigma^n} e^{-\frac{\sum_{i=1}^n (y_i - x_i^T \beta)^2}{2\sigma^2}}$$

2 pts.

- (b) Recall that ridge regression minimizes the following function:

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^d \beta_j^2$$

Given that $\sigma = 1$, how should one choose the prior on β such that the maximum a posteriori estimator is equal to $\hat{\beta}_{\text{ridge}}$? Please explain your answer.

$$\text{Ansatz: } \beta \sim N(0, \Lambda^{-1})$$

$$\Lambda = \frac{\lambda}{\sigma^2} \text{Id}$$

in this case

$$\max P(\beta | n, X, Y) = \max \lambda \cdot e^{-\frac{\sum_{i=1}^n (y_i - x_i^T \beta)^2}{2\sigma^2} + \frac{\beta^T \Lambda \beta}{2}}$$

$$\text{when } \Lambda = \frac{\lambda}{\sigma^2} \text{Id}$$

4 pts.

Question 2: Maximum Likelihood and Bayesian Learning (20 pts.)

1. Given are n random binary vectors, each of dimension d , $\mathcal{X} = \{X_1, \dots, X_n\}, X_i \in \{0, 1\}^d$. The components of the binary vectors are generated from one of the two sources:

$$\begin{aligned} X_{ij} &\sim \text{ber}(1 - \theta), & \lambda_j = 1, \\ X_{ij} &\sim \text{ber}(\theta), & \lambda_j = 0, \end{aligned}$$

where λ_j selects the Bernoulli source and $1 \leq j \leq d$.

- (a) What is the likelihood function $p(\mathcal{X} | \theta, \lambda)$?

$$p(\mathcal{X} | \theta, \lambda) = \prod_{i \leq n} \prod_{j \leq d} \text{ber}(X_{ij} | 1 - \theta)^{\lambda_j} \text{ber}(X_{ij} | \theta)^{1 - \lambda_j}$$

2 pts.

- (b) Derive the maximum likelihood estimate for θ given λ . Use the following abbreviation $n_j = \sum_{i \leq n} X_{ij}$.

$$\begin{aligned} \frac{d}{d\theta} (-\log p(\mathcal{X} | \theta, \lambda)) &= \frac{d}{d\theta} \left(-\sum_{j \leq d} \lambda_j [n_j \log(1 - \theta) + (n - n_j) \log \theta] \right. \\ &\quad \left. + (1 - \lambda_j) [n_j \log \theta + (n - n_j) \log(1 - \theta)] \right) \\ &= -\sum_{j \leq d} \lambda_j \left[\frac{n_j}{1 - \theta} + \frac{n - n_j}{\theta} \right] + (1 - \lambda_j) \left[\frac{n_j}{\theta} + \frac{n - n_j}{1 - \theta} \right] \\ &\propto -\sum_{j \leq d} \lambda_j [\theta n_j + (n - n_j)(1 - \theta)] + (1 - \lambda_j) [n_j(1 - \theta) + (n - n_j)\theta] \\ &= -\sum_{j \leq d} [\theta(2n_j - n) + n - n_j] + (1 - \lambda_j) [\theta(n - 2n_j) + n_j] = 0 \end{aligned}$$

$$\hat{\theta} = \frac{\sum_{j \leq d} \lambda_j (n - n_j) + (1 - \lambda_j) n_j}{\sum_{j \leq d} \lambda_j (2n_j - n) + (1 - \lambda_j) (n - 2n_j)}$$

4 pts.

- (c) Derive the posterior distribution $p(\theta | \mathcal{X}, \lambda, \alpha)$ in closed form given that the prior distribution on θ is beta distributed: $p(\theta | \alpha) = B(\alpha, \alpha)^{-1} \theta^{(\alpha-1)} (1 - \theta)^{(\alpha-1)}$, where $B(\alpha, \alpha)$ normalizes the distribution.

$$\begin{aligned}
p(\theta \mid \mathcal{X} \lambda, \alpha) &= a \prod_{j \leq d} (1 - \theta)^{(n_j \lambda_j + (n - n_j)(1 - \lambda_j) + \alpha - 1)} \theta^{((n - n_j)\lambda_j + n_j(1 - \lambda_j) + \alpha - 1)} B(\alpha, \alpha)^{-1} \\
&= \prod_{j \leq d} \text{beta}(\theta \mid n_j \lambda_j(n - n_j)(1 - \lambda_j) + \alpha, (n - n_j)\lambda_j + n_j(1 - \lambda_j) + \alpha)
\end{aligned}$$

3 pts.

- (d) Derive the evidence $p(\mathcal{X} \mid \lambda, \alpha)$.

$$p(\mathcal{X} \mid \lambda, \alpha) = \prod_{j \leq d} \frac{B(n_j \lambda_j + (n - n_j)(1 - \lambda_j) + \alpha, (n - n_j)\lambda_j + n_j(1 - \lambda_j) + \alpha)}{B(\alpha, \alpha)}$$

3 pts.

- (e) Under what conditions is the maximum likelihood estimator for θ equal to the maximum a posterior estimator?

When the prior is uniform or when $n \rightarrow \infty$.

2 pts.

2. Given is an input vector $X = (X_1, \dots, X_d)^T$ and model parameters $\beta \in \mathbb{R}^d$. A linear regression model predicts the response Y :

$$Y = X^T \beta + \epsilon, \quad \epsilon \sim N(\epsilon \mid 0, \sigma^2),$$

- (a) Given n i.i.d. observations, what is the likelihood function for β ?

$$p(Y \mid \beta, X) = \mathcal{N}(Y \mid X^T \beta, \mathbb{I}\sigma^2)$$

2 pts.

- (b) Recall that ridge regression maximizes the following function:

$$\hat{\beta}_{\text{ridge}} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2$$

Given that $\sigma = 1$, how should one choose the prior on β such that the maximum a posteriori estimator is equal to $\hat{\beta}_{\text{ridge}}$? Please explain your answer.

The MAP estimator is given by:

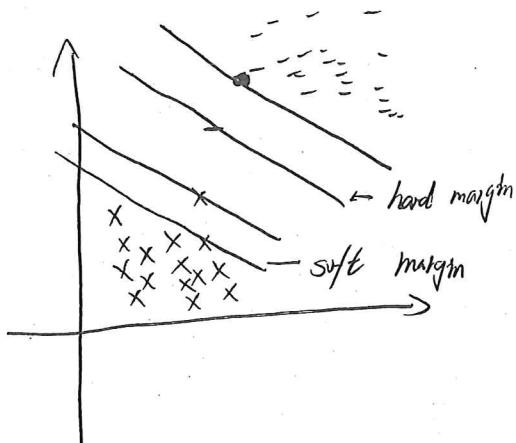
$$\begin{aligned}\hat{\beta} &= \arg \max_{\beta} \log p(\mathbf{y} | \mathbf{X}, \beta, \sigma^2) + \log p(\beta | \Lambda) \\ &\stackrel{(a)}{=} \arg \max_{\beta} -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2 - \frac{1}{2} \beta^T \Lambda \beta,\end{aligned}$$

where in (a) we use a Gaussian prior on β . If we choose $\Lambda = \mathbb{I}\lambda$ we obtain the ridge regression.

4 pts.

Question 3: Kernels (20 pts.)

- What is the difference between hard-margin and soft-margin SVM for classification? Please illustrate your argument with a figure. What is the meaning of the support vectors in each case?



in the hard margin
it is the
vectors that
lie on the boundary
that compose the
support vectors
in the latter one,
the boundary and the slack
vector lie on the variable

3 pts.

- Suppose your training dataset is perfectly linearly separable. Is there any reason (motivation) why to use soft-margin SVM for classification in this case?

avoid overfitting

1 pts.

3. What are the three main (distinct) advantages of using the Kernel trick?

- (1) do not need to explicitly express the feature transformation, more expressive
- (2) computational efficiency, since the kernel function
- (3) ~~can't~~ non-linear transformation, high-dimension representation

3 pts.

4. The Sigmoid Kernel is defined as: $K(x, y) = \tanh(\alpha x^T y + r)$ where α is a scaling parameter and r is a shifting parameter.

(a) Prove that the Sigmoid Kernel is not always a valid kernel. Which property of a valid kernel is violated by the Sigmoid Kernel? Please give a short description.

semi-positive

$$\left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)$$

Let $\alpha = 1$

$$\begin{bmatrix} \tanh(\frac{1}{k}) & \tanh(1) \\ \tanh(1) & \tanh(k^2) \end{bmatrix}$$
$$x_1 = \frac{1}{k}$$
$$x_2 = k$$

1 pts.

- (b) Suppose that x and y are feature vectors in $\{-1, 0, 1\}^N$. Is there any circumstance in which the Sigmoid Kernel is a valid kernel? If yes, please provide such an example.

$$\boxed{\alpha > 0} \quad \begin{array}{l} \text{是这样} \\ y \rightarrow 0 \end{array}$$

4 pts.

5. Support Vector Machines are inherently binary classifiers. How would you adapt the original SVM technique in order to perform multi-class classification? In particular:

- (a) Describe the main idea of your technique.

$$\min \frac{1}{2} w^T w$$

s.t. $(w_1^T x_i + w_2^T x_i) - \max_{2 \neq j} (w_j^T x_i + w_2) \geq 1$

$$\text{or} \quad \min \frac{1}{2} w^T w$$

s.t. $w^T \psi(z_i, y_i) - \max_{j \neq i} (w^T \psi(z_j, y_j) + \delta(z_j)) \geq -\rho_i$

1 pts. $\psi(z, y)$

(b) Provide a formal definition of the optimization problem during training.

3 pts.

(c) Provide the formula that you would use to do classification in test phase.

$$Z = \max(w_2^T x + w_{2,0})$$

2 pts.

(d) Mention one advantage and one disadvantage of your solution.

disadvantage: no probabilistic explanation did not discriminate the difference between of the class differences, may not be feasible

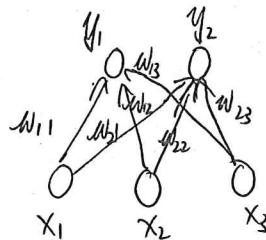
advantage: simple data model, easy to train.

In linear case
separable

2 pts.

Question 4: Neural Networks (20 pts.)

1. Sketch a neural-network-diagram of a linear discriminant model $y_k = \sum_{j=0}^2 w_{kj}x_j$ with two outputs $k = 1, 2$ and label all inputs, outputs and weights.



1 pts.

2. Why are activation functions in multilayer networks in general non-linear?

in order to ~~generalize~~ generalize of linear model
allow more non-linear for both regression
and classification

no ~~~symmetry~~ break the symmetry 1 pts.

3. Explain why weights are generally initialized to small values close to zero?

after normalization

- (1) Expectation is 0 since the backpropagation is sort of gradient descent, when the initial weight is large, the error will be large, so the step will be huge.

(2) since the expectation for weight is around 0

Vanish gradient

1 pts.

- (3) ~~Logistic function is close to flat for large positive or negative inputs.~~ 13

- (4) the same values of 0 will cause - output and does not break the symmetry of the net and

4. What is the computational advantage of activation functions like the sigmoid or the hyperbolic tangent?

1

$$f(x) = f(x)(1-f(x))$$

① Since the derivative of them are simple

and can be reused.

② they are differentiable.

③ yield a zero gradient

1 pts.

5. What difficulty is caused by step activation functions for the learning algorithm backpropagation?

non-differentiable

Vanishing gradient problem especially when the neural networks are deep, $f(x) \rightarrow 0$ when $x \rightarrow \infty$

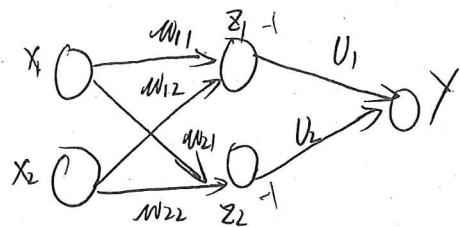
1 pts.

6. Consider a two layer feed-forward neural network with a 2 dimensional input (x_1, x_2) . The hidden layer (z_1, z_2) is also two dimensional. We will use the sigmoid $\sigma(x)$ as activation function in the hidden layer. In addition, both layer have biases, defined as b_1 and b_2 . The output y is one dimensional. Thus we can define the network by the following equations:

$$z_j = \sigma\left(\sum_{k=1}^2 w_{jk}x_k + b_1\right)$$

$$y = \sum_{k=1}^2 U_k z_k + b_2$$

(a) Sketch the neural-network-diagram of the above model.



1 pts.

(b) Initialize all w_{ij} and U_i to -1 , and set all biases b_i to 1 . The loss function is $L(t, y) = \frac{1}{2}(t - y)^2$ and you are given one data point $(1, \frac{1}{2}, \frac{1}{2})$ with label $t = 1$. With backpropagation, calculate the gradient of the loss with respect to the weight w_{11} , and use it, together with a learning rate $\eta = 4$, to update w_{11} .

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_{ij}}$$

$$= (t - y) \frac{\partial y}{\partial z_1} \frac{\partial z_1}{\partial w_{11}}$$

$$= (t - y) \sigma'(\cancel{\frac{1}{\sum} w_{ij} x_j + b_i})$$

$$= (t - y) \cdot \sigma'(\cancel{\frac{1}{\sum} w_{ij} x_j + b_i}) x_1 = (1 - 0) \times (-1) \times \cancel{\sigma'(\cancel{\frac{1}{\sum} w_{ij} x_j + b_i})}$$

$$= -1 \times \cancel{\frac{1}{2} \times \cancel{\frac{1}{2}}} \times \cancel{\frac{1}{2}} + \cancel{\frac{1}{He^1}} + \cancel{\frac{1}{He^1}}$$

$$= -\cancel{\frac{1}{8}} - \cancel{\frac{1}{(1+e^{-1})^2}} \frac{1}{8}$$

thus

$$w_{11}' = w_{11} + \eta \left(-\frac{\partial L}{\partial w_{11}} \right)$$

4 pts.

$$= -1 + 4 \times \cancel{\frac{1}{8} \times \cancel{\frac{1}{(1+e^{-1})^2}}}$$

$$= \cancel{\frac{1}{2}} - \cancel{\frac{4e^{-1}}{(1+e^{-1})^2}} = \frac{1}{2}$$

$$y = -\frac{1}{2} - \frac{1}{2} + 1 = 0$$

7. Multilayer Perceptrons

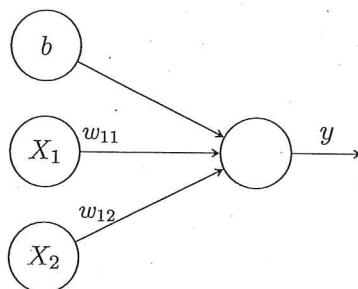
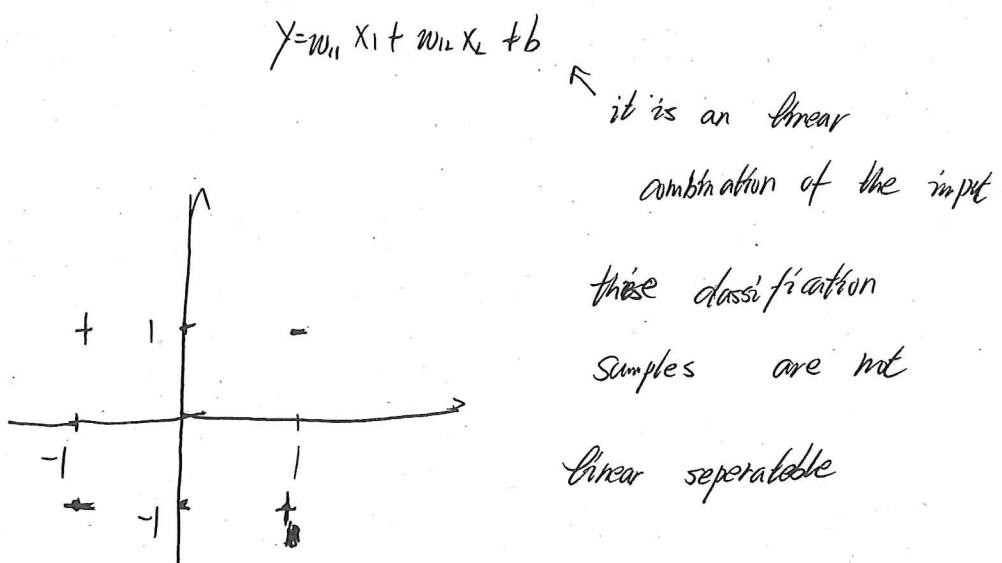


Figure 1: Single layer network.

- (a) Formally prove, that the single layer perceptron shown in Figure ?? can not correctly classify the following training data set \mathcal{Z} :

$$\mathcal{Z} = \{((1, 1), A), (-1, -1), A), ((1, -1), B), ((-1, 1), B)\}$$

Each element $((x_1, x_2), t)$ denotes a pattern and its associated label. A pattern is classified as class A if the network output y is less than zero and as class B else. The third unit b denotes the bias. Note that the activation function of the output is the identity.



From add (1, 1) $b < 0$

from add (1, -1) $b > 0$
contradict

suppose we have

$$\left\{ \begin{array}{l} w_{11} + w_{12} + b < 0 \quad \text{--- ①} \\ -w_{11} + (-w_{12}) + b < 0 \quad \text{--- ②} \\ w_{11} - w_{12} + b > 0 \quad \text{--- ③} \\ -w_{11} + w_{12} + b > 0 \quad \text{--- ④} \end{array} \right.$$

4 pts.

from

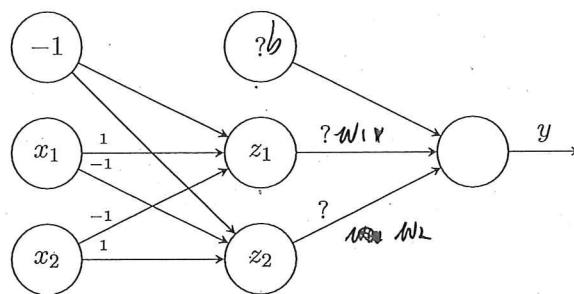
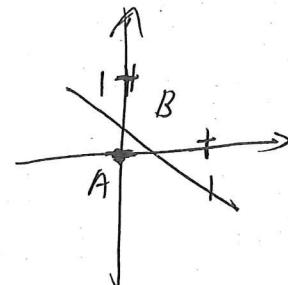
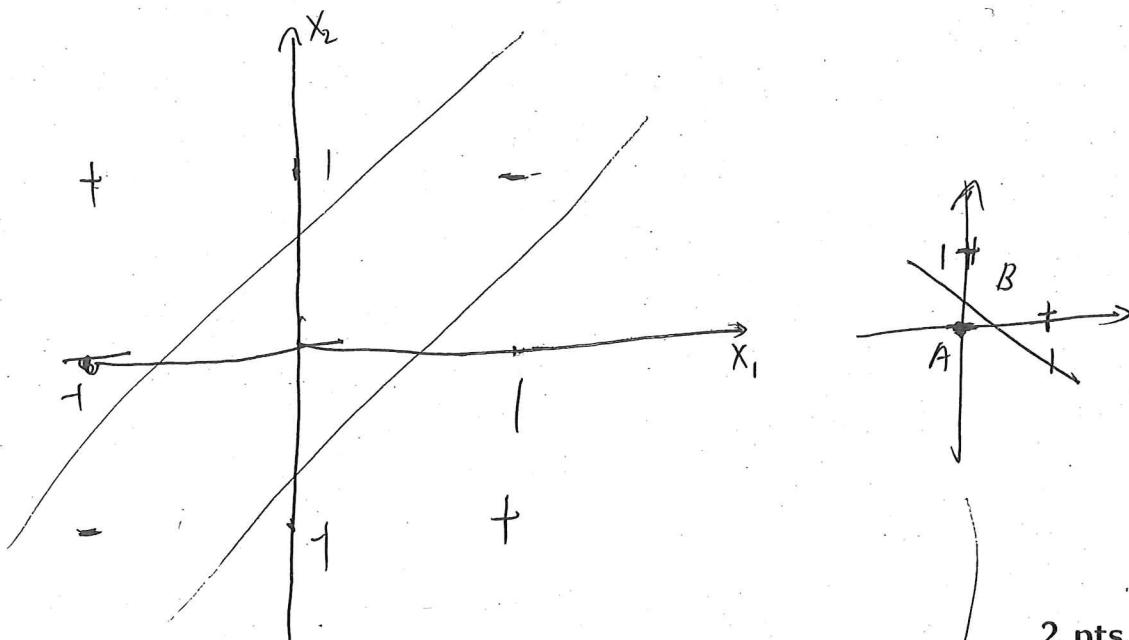


Figure 2: Two layer network.

- (b) Now consider the two-layer network in Figure ???. Assume that the activation function in the hidden layer is a step functions, which is 0 if their input is less than zero and 1 otherwise. Show that it can correctly classify the training data set \mathcal{Z} by selecting an appropriate set of second layer weights and bias.

$$\begin{aligned}
 & \text{input}(z_1) = x_1 - x_2 - 1 \\
 & \text{input}(z_2) = -x_1 + x_2 - 1 \\
 & \text{for } z \text{ with } A, \text{ output}(z) = 0 \\
 & \quad \dots B, \quad \begin{cases} (1, -1), \\ (1, 1), \end{cases} \quad 0 \\
 & \text{for } z \text{ with } A, \text{ output}(z) = 0 \\
 & \quad B, \quad (1, -1) = 0 \\
 & \quad (1, 1) \quad | \\
 & \text{in this class case} \quad \begin{cases} (-1, 1) \rightarrow (0, 0) \\ (1, 1) \rightarrow (0, 0) \end{cases} \\
 & \quad (1, -1) \rightarrow (1, 0) \\
 & \quad (-1, 1) \rightarrow (0, 1) \\
 & Y = w_1 z_1 + w_2 z_2 + b \\
 & b < 0 \\
 & \left\{ \begin{array}{l} w_1 + b > 0 \\ w_2 + b > 0 \end{array} \right. \quad \text{which can be met} \\
 & \quad \text{by selection} \quad 2 \text{ pts.}
 \end{aligned}$$

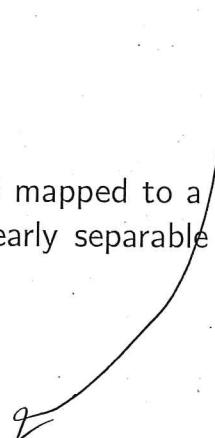
- (c) Plot the decision regions in the (x_1, x_2) plane, i.e. the areas classified as A and B , respectively.



2 pts.

- (d) Plot how the original input space (x_1, x_2) is mapped to a new feature space (z_1, z_2) . Is the training data set linearly separable in the new feature space?

Yes



2 pts.

Question 5: Gaussian process regression (20 pts.)

Consider the Gaussian process regression model, $y = f(\mathbf{x}) + \epsilon$, where $\mathbf{x} \in \mathbb{R}^d$ is the input vector, y is the response, ϵ is the Gaussian noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$, $f(\mathbf{x}) \sim \mathcal{GP}(\mathbf{0}, k(\mathbf{x}, \mathbf{x}'))$ with mean $\mathbf{0}$ and covariance function $k(\mathbf{x}, \mathbf{x}')$.

1. There are n training vectors associated with the responses: $X \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^n$. A single test vector is \mathbf{x}_* , assume the predictive target of it is y_* . For notation simplicity, let

$$K(X, X) := \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix}, K(\mathbf{x}_*, X) := [k(\mathbf{x}_*, x_1), \dots, k(\mathbf{x}_*, x_n)]$$

- (a) Derive the joint distribution of $\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix}$:

$$\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix} \mid \begin{bmatrix} \cancel{K(X)} & C \\ C^T & K_{**} \end{bmatrix} \right)$$

$$\text{with } Q = K(X, X) + \sigma^2 I$$

$$C = k(\mathbf{x}_*, \mathbf{x}_*) + \sigma^2$$

$$k = k(\mathbf{x}_*, \mathbf{x})^T$$

$$K^T = K(\mathbf{x}_*, \mathbf{x})$$

2 pts.

- (b) According to the joint distribution, determine the predictive distribution $y_* | X, \mathbf{y}, \mathbf{x}_* \sim \mathcal{N}(\bar{y}_*, V[y_*])$. Hint: using the conditioning rule: Let \mathbf{a} and \mathbf{b} be jointly Gaussian random vectors

$$\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix} \right)$$

the conditional distribution of \mathbf{b} given \mathbf{a} is: $\mathbf{b}|\mathbf{a} \sim \mathcal{N}(\mu_b + C^T A^{-1}(\mathbf{a} - \mu_a), B - C^T A^{-1}C)$.

$$\tilde{y}_* = \frac{k(x_*, x)}{k(x, x) [k(x, x) + \sigma^2]^{-1}} y$$

$$V(y_*) = k(x_*, x_*) - k(x_*, x) [k(x, x) + \sigma^2]^{-1} k(x_*, x)^T$$

2 pts.

2. Consider the special situation: Given one scalar data point (x, y) , and the squared exponential kernel with length-scale l : $k(x, x') = \exp(-\frac{(x - x')^2}{2l^2})$, for the test point x_* .

- (a) Write down the specific predictive mean and variance.

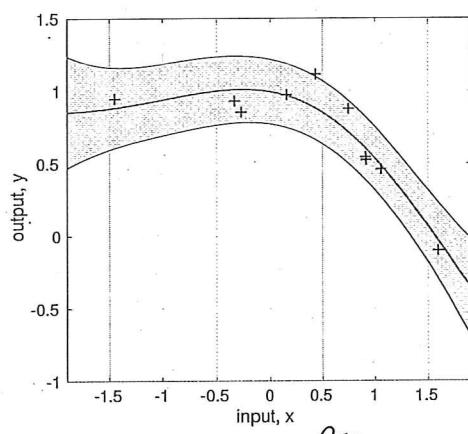
$$\tilde{y}_* = \exp\left(-\frac{(x_* - x)^2}{2l^2}\right) [1 + \sigma^2]^{-1} y$$

$$= \frac{y}{1 + \sigma^2} \cdot e^{-\frac{(x_* - x)^2}{2l^2}}$$

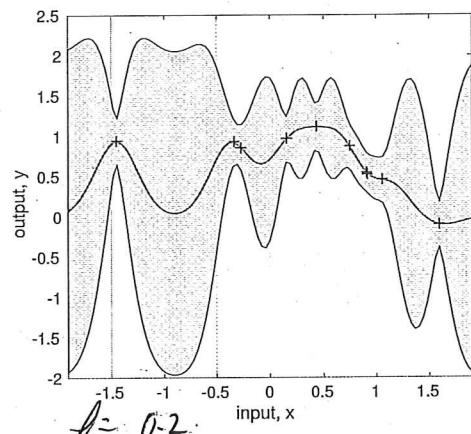
$$V(y_*) = 1 - \frac{\exp\left(-\frac{(x_* - x)^2}{2l^2}\right)}{1 + \sigma^2}$$

2 pts.

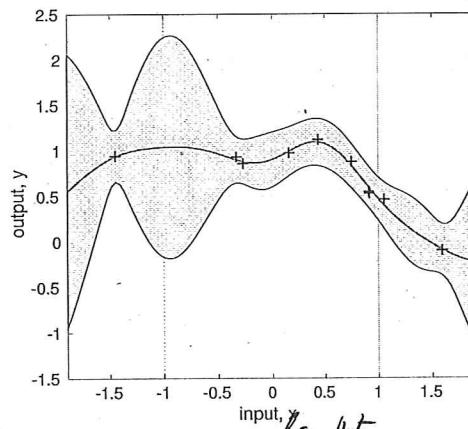
- (b) Fixing the noise level $\sigma = 0.1$, in Figure ?? you can see four different plots showing predictions made using length-scales of $l = \{0.2, 0.5, 1.2, 2\}$, with ten scalar training data points (crosses in the Figure). Specify which l was used for each of these four predictions. Please indicate your choice by writing the value of l below the corresponding plot.



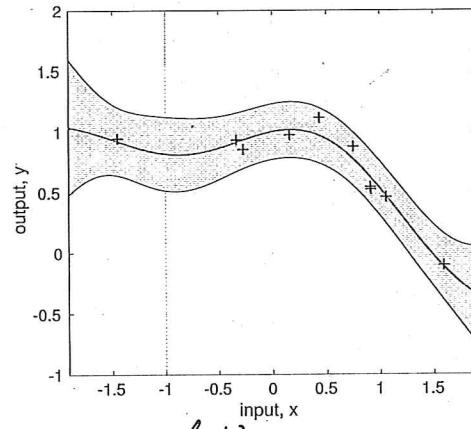
$$l=2$$



$$l=0.2$$



$$l=0.5$$



$$l=1.2$$

Figure 3: Predictions using four different length-scales.

$$\frac{y}{1+0.1} e^{-\frac{(x-x)^2}{2l^2}}$$

$$1-\frac{exp}{1.01} - \frac{exp}{1.01}$$

3 pts.

(c) Briefly describe what role the length-scale l plays.

the scale of the dissimilarity, how to discriminate

Sensitivity (regularization)

$\ell \uparrow$, more smooth

$\ell \downarrow$, more adapt to the data change

3 pts.

3. Let θ represent all the hyperparameters of the GP (including hyperparameters in the kernel, e.g. l , and noise level σ). Usually the optimal hyperparameters are found by maximizing the log marginal likelihood ($\log \mathbb{P}(y|X, \theta)$) with respect to θ .

- (a) Prove that the log marginal likelihood is

$$\log \mathbb{P}(y|X, \theta) = -\frac{1}{2} y^T [K(X, X) + \sigma^2 I]^{-1} y - \frac{1}{2} \log |K(X, X) + \sigma^2 I| - \frac{n}{2} \log 2\pi$$

where I is the identity matrix, $|\cdot|$ is the determinant of a matrix.

$$\log P(y|x, \theta) = \frac{1}{\sqrt{(2\pi)^n |V(x)|}} \left[-\frac{1}{2} (y - \bar{y})^T V^{-1}(x) (y - \bar{y}) \right]$$

$$\log P(y|x, \theta) = -\frac{n}{2} \log 2\pi - \frac{1}{2} \log |K(x, x) + \sigma^2 I| - \frac{1}{2} y^T [K(x, x) + \sigma^2 I]^{-1} y$$

4 pts.

- (b) Generally, can you find an analytic expression for the optimal hyperparameter θ^* that maximize the log marginal likelihood? If not, suggest an approach to compute θ^* .

On hard to find the close-form solution
since the inversion of the matrix can be too huge

use Newton update

$$\theta' = \theta_0 - \frac{\cancel{\log p(y|x_0)}}{\cancel{\partial \log p(y|x_0)} / \cancel{\partial \theta}}$$

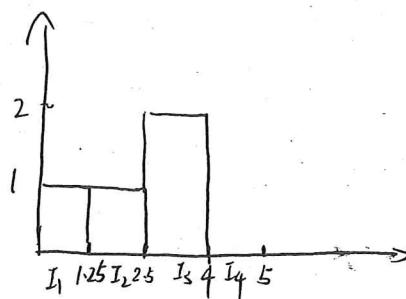
4 pts.

Question 6: Unsupervised learning (20 pts.)

1. Non-parametric density estimation.

Consider the following datasets of 1-dimensional data points: $D^+ = \{d_1, d_2, d_3, d_4\}$ where $d_1 = 1, d_2 = 1.5, d_3 = 3$ and $d_4 = 3.1$ and $D^- = \{d_5, d_6, d_7\}$ where $d_5 = 1.3, d_6 = 3.9, d_7 = 4$.

- (a) Show the histogram of the dataset D^+ with intervals $I_1 = [0, 1.25)$, $I_2 = [1.25, 2.5)$ and $I_3 = [2.5, 4)$.



1 pts.

- (b) Kernel density estimators and nearest neighbor estimates use the following formula to approximate the probability density for any point x :

$$p(x) = \frac{K(x)}{nV}$$

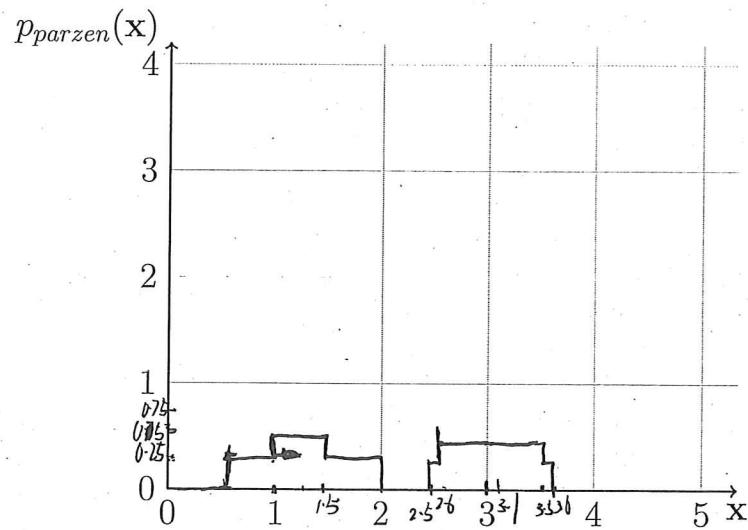
where $K(x)$ is the number of samples falling within a window around x , n is the total number of training samples and V is the area of the window around x .

Let $k(\frac{x-x_i}{h})$ be a Parzen window of the unit cube function scaled by the bandwidth h . The unit cube function is defined as follows:

$$k(u) = \begin{cases} 1 & \text{if } |u_i| \leq 1/2 \text{ for all } i \in 1, \dots, d \\ 0 & \text{otherwise} \end{cases}$$

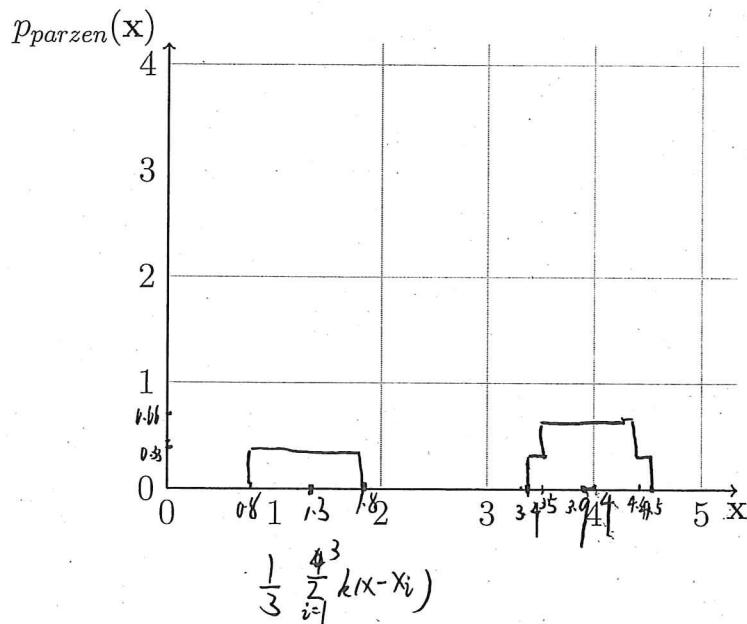
- i. Plot (draw) the resulting probability density function $p_{parzen}(x)$ for the 1-dimensional dataset D^+ and for $h = 1$. Use the provided plotting space.

$$\frac{1}{n} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right) = \frac{1}{4} \sum_{i=1}^4 k(x-x_i)$$



2 pts.

- ii. Plot (draw) the resulting probability density function $p_{parzen}(x)$ for the 1-dimensional dataset D^- and for $h = 1$. Use the provided plotting space.



2 pts.

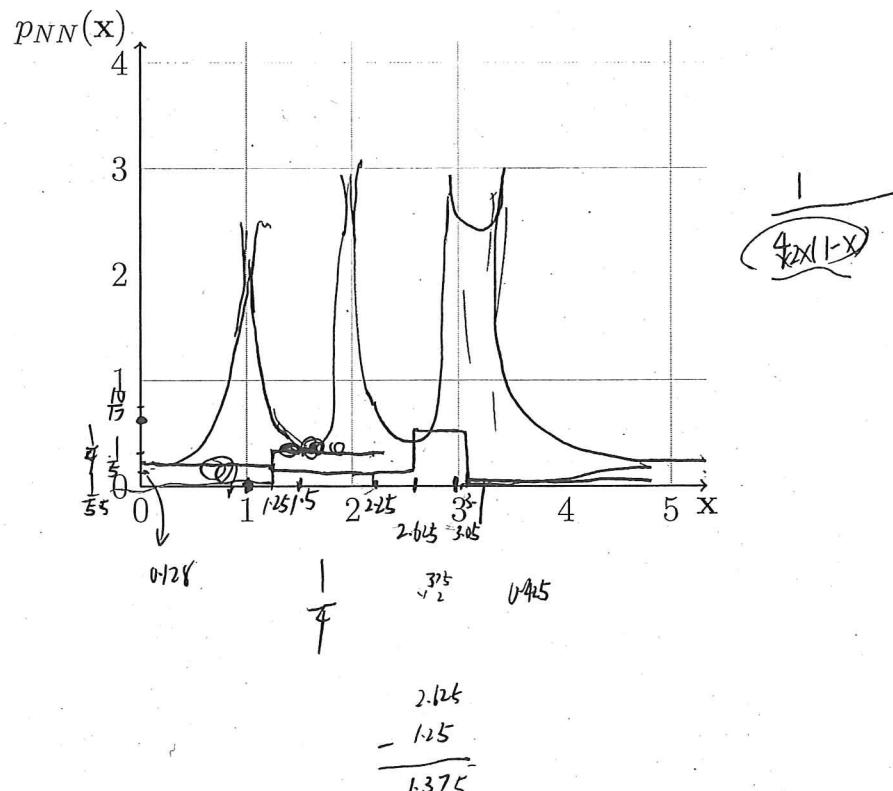
- iii. We are now given the point $x = 1.2$, which we want to classify in either one of the following two classes $C = \{c^+, c^-\}$. The training data is $D = D^+ \cup D^-$ where D^+ and D^- are defined as before and contain the samples of the class c^+ and c^- . In order to do classification based on the previously introduced density we need to estimate the following table of probabilities. Fill in the table using

your results from (i) and (ii) and applying Bayes rule:

For $x = 1.2$ and using a Parzen window estimate from (b)				
	$p(c x)$	$p(x c)$	$p(c)$	$p(x)$
$c = c^+$	$\frac{2}{3}$	$\frac{1}{2}$	$\frac{4}{7}$	$\frac{3}{7}$
$c = c^-$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{3}{7}$	$\frac{3}{7}$

$$\sum p(x|c) p(c) = \frac{1}{8} \times \frac{4}{7} + \frac{1}{3} \times \frac{3}{7} = \frac{3}{7} \text{ pts.}$$

- iv. Now plot the related $p_{NN}(x)$ for D^+ using a nearest neighbor density estimator. In nearest neighbor density estimation the size of the window in $p(x) = \frac{K(x)}{nV}$ is chosen such that exactly one sample falls into the hyper-sphere around x (i.e. $K(x) = 1$). You can resolve cases where it is not possible to set $K(x) = 1$ (i.e. tie-breaks) arbitrarily. Use the provided space below.



4 pts.

v. Compute the values of $\int_{-\infty}^{\infty} p_{parzen}(x)dx$ and $\int_{-\infty}^{\infty} p_{NN}(x)dx$ obtained in (i) and (iv) and compare them.

$$\int_{-\infty}^{\infty} p_{parzen}(x) = \frac{1}{3} + \frac{1}{3}x(0.2) + \frac{2}{3}x(0.9)$$

$$= 1$$

$$\int_{-\infty}^{\infty} p_{NN}(x) dx$$

$$= \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

2 pts.

$$\int_0^1 \frac{1}{8(1-x)} dx + \int_1^{1.5} \frac{1}{8(x-1)} dx + \int_{1.5}^{3} \frac{1}{8(3-x)} dx$$

$$+ \int_{3}^{3.1} \frac{1}{8(x-3)} dx + \int_{3.1}^{3.05} \frac{1}{8(3.1-x)} dx$$

$$+ \cancel{\int_{3.05}^{3.1} \frac{1}{8(3.1-x)} dx} + \int_{3.1}^{\infty} \frac{1}{8(x-3.1)} dx$$

$\leftarrow D$

2. *Expectation-Maximization algorithm.*

Assume the following mixture of d -dimensional densities, consisting of K components:

$$p(\mathbf{x}|\mu_1, \dots, \mu_K) = \frac{1}{K} \sum_{c=1}^K p(\mathbf{x}|\mu_c),$$

where

$$p(\mathbf{x}|\mu_c) = (2\pi)^{-d/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^T(\mathbf{x} - \mu_c)\right)$$

Here, the parameters μ_c , $c = 1, \dots, K$ are d -dimensional vectors.

Let $\{\mathbf{x}_i\}_{i=1}^N$ be a dataset generated from the mixture above. For the EM-algorithm we introduce latent variables $\gamma_{i,c}$, $i = 1, \dots, N$; $c = 1, \dots, K$.

$\gamma_{i,c}$ indicates the probability of the point \mathbf{x}_i belonging to class c .

- (a) What is the difference between this model and the traditional Gaussian Mixture Model?

the prior are the uniform distribution
and same σ^2 for each gaussian.

2 pts.

- (b) Derive the update step for $\gamma_{i,c}$ given the model parameters μ_c .

Hint: recall that $\gamma_{i,c}$ is defined to be $\gamma_{i,c} = p(c|\mathbf{x}_i)$

$$\begin{aligned} \gamma_{ic} &= \frac{p(c|\mathbf{x}_i, \theta) p(\theta)}{p(\mathbf{x}_i)} \\ &= \frac{\cancel{p(c)}}{\cancel{p(\mathbf{x}_i)}} \cdot \frac{\exp(-\frac{1}{2}(\mathbf{x}_i - \mu_c)^T(\mathbf{x}_i - \mu_c))}{\cancel{\sum_{j=1}^K} \exp(-\frac{1}{2}(\mathbf{x}_i - \mu_j)^T(\mathbf{x}_i - \mu_j))} \end{aligned}$$

3 pts.

- (c) Derive the update step for the model parameters $\mu_c, c = 1, \dots, K$ given the latent variables $\gamma_{i,c}$.

$$\frac{\sum_i^N \gamma_{i,c} x_i}{\sum_i^N \gamma_{i,c}} = \mu_c'$$

$$Q(\theta; \theta') = \sum_{i=1}^N \sum_{c=1}^K \gamma_{i,c} \log \frac{1}{K} \cdot (2\pi)^{-\frac{d}{2}} \exp \left(-\frac{1}{2} (x_i - \mu_c)^T (x_i - \mu_c) \right) \quad 3 \text{ pts.}$$

$$= \sum_{i=1}^N \sum_{c=1}^K \gamma_{i,c} \left(-\log K - \frac{d}{2} \log 2\pi - \frac{1}{2} (x_i - \mu_c)^T (x_i - \mu_c) \right) \quad \cancel{\text{---}}$$

$$\frac{\partial Q}{\partial \mu_c} = \sum_{x \in X} \gamma_{i,c} \Sigma_c^{-1} (x - \mu_c) = 0$$

Supplementary Sheet

Supplementary Sheet

Supplementary Sheet