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 five sentences per answer.

a.  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 $D$.

$$\underbrace{p(\mathcal{M}|D)}_{\text{posterior}} = \frac{p(D|\mathcal{M})}{p(D)} \underbrace{p(\mathcal{M})}_{\text{prior}}$$

Bayes rule relates what we know about a model before (prior) and after (posterior) having seen the data. The difference between the prior and posterior distributions for the model can be interpreted as a 'machine learning' effect. (Alternative answers are also possible).

b.  Explain the relation between Bayesian estimation, Maximum a Posteriori (MAP) estimation and Maximum Likelihood (ML) estimation. You may assume a context of a given model structure with unknown parameters $\theta$ and an observed data set $D$.

$$\hat{\theta}_{\text{bayes}} = \int_{\theta} \theta p(\theta|D)\, d\theta \quad \text{(Bayes est.)}$$

$$\hat{\theta}_{\text{map}} = \arg\max_{\theta} p(\theta|D) = \arg\max_{\theta} p(D|\theta)p(\theta) \quad \text{(MAP)}$$

$$\hat{\theta}_{\text{ml}} = \arg\max_{\theta} p(D|\theta) \quad \text{(ML)}$$

Bayes estimation picks the mean from the posterior $p(\theta|D)$. MAP picks the mode from $p(\theta|D)$. ML is MAP with uniform prior. (Alternative answers are also possible).

The following two sub-questions relate to a (Factor Analysis) model $x_n = \Lambda z_n + v_n$ for an observed data set $D = \{x_1, \ldots, x_N\}$. The modeling assumptions include $z_n \sim \mathcal{N}(0, I)$, $v_n \sim \mathcal{N}(0, \Psi)$ and $\varepsilon[z_n v_n^T] = 0$.

c.  Show that the covariance matrix of the observed data $x_n$ is equal to $\Lambda\Lambda^T + \Psi$.

$$\epsilon[x] = \epsilon[\Lambda z + v] = \Lambda\epsilon[z] + \epsilon[v] = 0$$
$$\text{var}[x] = \epsilon[(x - \epsilon[x])(x - \epsilon[x])^T] = \epsilon[(\Lambda z + v)(\Lambda z + v)^T]$$
$$= \Lambda\epsilon[zz^T]\Lambda^T + \epsilon[vv^T] = \Lambda\Lambda^T + \Psi$$

d.  Why is this model not interesting for unconstrained $\Psi$? How does probabilistic PCA handle this problem?

(a) Because setting $\Lambda = 0$ would result in a 'regular' gaussian model with covariance matrix $\Psi$; i.o.w. it's no more interesting than any other gaussian model.
(b) If $\Psi$ is diagonal, then all correlations between the components in $x$ *must* be absorbed ('explained') by the rank-$K$ matrix $\Lambda\Lambda^T$. In pPCA, this is achieved by the constraint $\Psi = \sigma^2 I$.

e.  Which of the following statements are justified? You can pick more than one and read the sign
    '∼' as: 'is similar to'. (Just pick the correct statements; no explanation needed).
    1: discriminative classification $\sim$ density estimation
    2: generative classification $\sim$ density estimation
    3: hidden Markov model $\sim$ factor analysis through time
    4: Kalman filtering $\sim$ unsupervised regression through time
    5: clustering $\sim$ supervised classification

    > 2 and 4 are correct.

**2.** (EM for 2-component Gaussian mixture). Consider an observed IID data set $D = \{x_1, \ldots, x_N\}$
    and a proposed model,

$$p(x_n) = \sum_{k=0}^{1} p(x_n, z_n = k|\pi)$$
$$= p(z_n = 1|\pi)p(x_n|z_n = 1) + p(z_n = 0|\pi)p(x_n|z_n = 0)$$
$$= \pi\mathcal{N}_1(x_n) + (1 - \pi)\mathcal{N}_0(x_n)$$

where we used shorthand notation $\mathcal{N}_k(x_n) \equiv (2\pi\sigma_k^2)^{-1/2} \exp\left(-(x_n - \mu_k)^2/(2\sigma_k^2)\right)$ for the Gaussian distribution. We assume that the parameters $\theta = (\mu_0, \sigma_0^2, \mu_1, \sigma_1^2)$ are known, but the *mixing proportion* parameter $\pi$ is unknown. The random variable $z_n \in \{0, 1\}$ is an *unobserved* 'cluster label'. In this question we will derive an EM-algorithm for maximum likelihood estimation of $\pi$. Let's assume that a estimate $\hat{\pi} = \pi^{(j)}$ is available from the previous iteration. We will now focus on the $(j + 1)$-th iteration in the EM algorithm.

a.  Describe shortly the E- and M-steps in the $(j + 1)$-th iteration of the EM-algorithm. In particular, complete the following equation set (fill in the stars) for the $(j + 1)$-th iteration and shortly describe the meaning of the equations: (Note: the expression $\langle f(x) \rangle_{p(x)}$ stands for the expectation of $f(x)$ w.r.t. probability distribution $p(x)$.)

$$q_n^{(j+1)} = p(\star|\star) \quad \text{(E-step)}$$
$$\pi^{(j+1)} = \arg\max_{\pi}\langle\star\rangle_{\star} \quad \text{(M-step)}$$

> $$q_n^{(j+1)} = p(z_n|x_n, \pi^{(j)}) \quad \text{(E-step)}$$
> $$\pi^{(j+1)} = \arg\max_{\pi}\langle\sum_n p(x_n, z_n|\pi)\rangle_{q_n^{(j+1)}} \quad \text{(M-step)}$$
>
> **E-step**: $q_n^{(t+1)}$ is the posterior probability of $z_n$, given observation $x_n$ and an estimate $\pi^{(j)}$ from the previous iteration. $q_n$ represents our knowledge about $z_n$.
> **M-step**: Maximizes the expected complete-data log-likelihood. Through Jensen's inequality it can be proved that this procedure increases the (observed data) log-likelihood $p(D|\pi)$.

b.  Work out $p(x_n, z_n = 1|\pi)$ (hint: use product rule). Work out $p(x_n, z_n = 0|\pi)$. And now work out the joint distribution $p(x_n, z_n|\pi)$ to a Bernoulli distribution (as in eq.A1, see formula cheat sheet). In this question, you need to work out the probabilities in terms of $z_n$, $\mathcal{N}_0(x_n)$, $\mathcal{N}_1(x_n)$ and $\pi$.

> $p(x_n, z_n = 1) = p(x_n|z_n = 1)p(z_n = 1) = \pi\mathcal{N}_1(x_n)$
> $p(x_n, z_n = 0) = p(x_n|z_n = 0)p(z_n = 0) = (1 - \pi)\mathcal{N}_0(x_n)$
> $p(x_n, z_n|\pi) = [\pi\mathcal{N}_1(x_n)]^{z_n}[(1 - \pi)\mathcal{N}_0(x_n)]^{1-z_n}$

c.  Show that the complete-data log-likelihood $\ell_c(\pi) = \sum_n \log p(x_n, z_n|\pi)$ can be worked out to

$$\ell_c(\pi) = \sum_n z_n \log \frac{\pi\mathcal{N}_1(x_n)}{(1 - \pi)\mathcal{N}_0(x_n)} + \sum_n \log(1 - \pi)\mathcal{N}_0(x_n) \tag{1}$$

$$\ell_c(\pi) = \sum_n \log p(x_n, z_n | \pi)$$

$$= \sum_n \log \left( [\pi \mathcal{N}_1(x_n)]^{z_n} [(1 - \pi) \mathcal{N}_0(x_n)]^{1 - z_n} \right)$$

$$= \sum_n z_n \log \pi \mathcal{N}_1(x_n) + \sum_n (1 - z_n) \log(1 - \pi) \mathcal{N}_0(x_n)$$

$$= \sum_n z_n \log \frac{\pi \mathcal{N}_1(x_n)}{(1 - \pi) \mathcal{N}_0(x_n)} + \sum_n \log(1 - \pi) \mathcal{N}_0(x_n)$$

To finalize the E-step, we now take the expectation of the complete-data log-likelihood with respect to the posterior distribution $p(z_n | x_n, \pi^{(j)})$. It follows from Eq.1 that we need to compute the expected value of $z_n$. We'll compute the expected value of $z_n$ in two stages:

d.   First show that the expectation $\sum_{\{z_n\}} z_n \cdot p(z_n | x_n, \pi^{(j)})$ can be worked out as follows:

$$\sum_{\{z_n\}} z_n p(z_n | x_n, \pi^{(j)}) = p(z_n = 1 | x_n, \pi^{(j)})$$

$$\sum_{\{z_n\}} z_n p(z_n | x_n, \pi) = 0 \cdot p(z_n = 0 | x_n, \pi) + 1 \cdot p(z_n = 1 | x_n, \pi)$$

$$= p(z_n = 1 | x_n, \pi)$$

e.   And now use Bayes rule to work out an expression for $p(z_n = 1 | x_n, \pi^{(j)})$ in terms of $\pi^{(j)}$, $\mathcal{N}_0(x_n)$ and $\mathcal{N}_1(x_n)$.

$$p(z_n = 1 | x_n, \pi^{(j)}) = \frac{p(x_n | z_n = 1) p(z_n = 1 | \pi^{(j)})}{\sum_k p(x_n | z_n = k) p(z_n = k | \pi^{(j)})}$$

$$= \frac{\pi^{(j)} \mathcal{N}_1(x_n)}{\pi^{(j)} \mathcal{N}_1(x_n) + (1 - \pi^{(j)}) \mathcal{N}_0(x_n)}$$

If we use shorthand notation $\zeta_n = p(z_n = 1 | x_n, \pi^{(j)})$, then the expected complete-data log-likelihood can be written as

$$\langle \ell_c(\pi) \rangle = \sum_n \zeta_n \log \frac{\pi \mathcal{N}_1(x_n)}{(1 - \pi) \mathcal{N}_0(x_n)} + \sum_n \log(1 - \pi) \mathcal{N}_0(x_n)$$

f.   Set $\partial \langle \ell_c(\pi) \rangle / \partial \pi = 0$ and obtain an expression for $\pi^{(j+1)}$ in terms of $\pi^{(j)}$, $\mathcal{N}_0(x_n)$ and $\mathcal{N}_1(x_n)$ (i.e. write down the $(j + 1)$-th iteration of the M-step).

$$\frac{\partial \langle \ell_c(\pi) \rangle}{\partial \pi} = \sum_n \frac{\zeta_n}{\pi} + \sum_n \frac{\zeta_n}{1 - \pi} - \sum_n \frac{1}{1 - \pi}$$

$$= \frac{1}{\pi(1 - \pi)} \sum_n (\zeta_n - n\pi)$$

Set derivative to zero and it follows that

$$\pi^{(j+1)} = \frac{1}{N} \sum_n \zeta_n$$

$$= \frac{\pi^{(j)} \mathcal{N}_1(x_n)}{\pi^{(j)} \mathcal{N}_1(x_n) + (1 - \pi^{(j)}) \mathcal{N}_0(x_n)}$$

**3.** You observe some data $x^n$. You ask two experts to explain the data.

Expert $A$ uses a data compression system that needs 1537 bits to describe the parameters of the model and 438 bits to describe the data given the model.

Expert $B$ gives you a system that needs 1325 bits for the parameters and 650 bits for the data, given the model.

a. Which expert's result do you prefer?
Explain (briefly) why you select that experts result.

> The total description length of $A$'s result is $1537 + 438 = 1975$ bits. For expert $B$ the total description length is $1325 + 650 = 1975$ bits. So both experts achieve the same complexity. In accordance with Occam's razor I prefer expert $B$'s explanation because his/her model is less complex.

b. You ask two additional experts.
Expert $C$ gives you a model with a parameter description length of 1471 bits and a data description that needs 450 bits.
Expert $D$ gives you a model with a parameter description length of 1464 bits and a data description that needs 543 bits.
Of the four experts $A$ to $D$, which result do you prefer, and why?

> The total complexity for expert $C$ is $1471 + 450 = 1921$ bits and for expert $D$ it is $1464 + 543 = 2007$ bits. Expert $D$ explanation is more complex than any of the three others so I reject it in accordance with the MDL principle. For the same reason I prefer expert $C$'s explanation, because it has the smallest overall complexity although the model complexity is larger than for expert $B$.

**4.** Let $X$ be a real valued random variable with probability density

$$p_X(x) = \frac{e^{-x^2/2}}{\sqrt{2\pi}}, \quad \text{for all } x.$$

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

$$p_{Y|X}(y|x) = \frac{e^{-(y-x)^2/2}}{\sqrt{2\pi}}, \quad \text{for all } x \text{ and } y.$$

a. Give an (integral) expression for $p_Y(y)$.
Do not try to evaluate the integral.

> $$p_Y(y) = \int_{-\infty}^{\infty} p_X(x) p_{Y|X}(y|x)\, dx = \int_{-\infty}^{\infty} \frac{e^{-\frac{1}{2}(x^2 + (y-x)^2)}}{2\pi}\, dx$$

b.　Approximate $p_Y(y)$ using the Laplace approximation.
Give the detailed derivation, not just the answer.
Hint: You may use the following results.

Let

$$g(x) = \frac{e^{-x^2/2}}{\sqrt{2\pi}}, \quad \text{and}$$

$$h(x) = \frac{e^{-(y-x)^2/2}}{\sqrt{2\pi}}, \quad \text{for some real value } y.$$

Then

$$\frac{\partial}{\partial x}g(x) = -xg(x)$$

$$\frac{\partial^2}{\partial x^2}g(x) = (x^2 - 1)g(x)$$

$$\frac{\partial}{\partial x}h(x) = (y - x)h(x)$$

$$\frac{\partial^2}{\partial x^2}h(x) = ((y - x)^2 - 1)h(x)$$

---

Using the hint we determine the first derivative of

$$f(x) = g(x)h(x),$$
$$\frac{\partial}{\partial x}f(x) = \frac{\partial}{\partial x}g(x) \cdot h(x) = -xg(x)h(x) + g(x)(y - x)h(x) = (y - 2x)f(x)$$

Setting this to zero gives

$$y - 2x = 0; \quad \text{so} \quad x = \frac{1}{2}y.$$

$$\frac{\partial}{\partial x}\ln f(x) = \frac{\frac{\partial}{\partial x}f(x)}{f(x)} = (y - 2x)$$

$$\frac{\partial^2}{\partial x^2}\ln f(x) = \frac{\partial}{\partial x}(y - 2x) = -2.$$

So, we find $A = 2$, see lecture notes, and thus

$$p_Y(y) = \int_{-\infty}^{\infty} f(x)\,dx \approx f(\frac{y}{2})\sqrt{\frac{2\pi}{A}}$$

$$= g(\frac{y}{2})h(\frac{y}{2})\sqrt{\frac{2\pi}{A}}$$

$$= \frac{1}{\sqrt{2\pi \cdot 2}}e^{-y^2/4}.$$

So $Y$ is a Gaussian with mean $m = 0$ and variance $\sigma^2 = 2$.

---

**5.**　We implement an e-mail spam filter using two features that we can extract from an e-mail. A feature can be the occurrence of a particular word or phrase in the e-mail.

Given an e-mail $E$ we denote the extracted features by $F$ and $G$.
$F = 1$ means that feature $F$ is present in the e-mail $E$.
$F = 0$ means that feature $F$ is absent. And likewise for feature $G$.
The variable $C$ indicates whether $E$ is spam ($C = 1$) or not ($C = 0$).

We are given 247 e-mails that are already classified. The following table shows how many e-mails contained certain features and the classification.

| $F$ | $G$ | $C$ | nr of e-mails |
|---|---|---|---|
| 0 | 0 | 0 | 15 |
| 0 | 0 | 1 | 28 |
| 0 | 1 | 0 | 18 |
| 0 | 1 | 1 | 25 |
| 1 | 0 | 0 | 8 |
| 1 | 0 | 1 | 75 |
| 1 | 1 | 0 | 10 |
| 1 | 1 | 1 | 68 |

a. From the table given above you can determine probability estimates using the maximum likelihood estimates. e.g. the probability $P(C = 1)$, i.e. the probability that an email will be spam, is approximated by:

$$P(C = 1) = \frac{\text{\# of e-mails with } C = 1}{\text{total \# of e-mails}} = \frac{196}{247} = 0.7935.$$

Note that the method using a beta prior would be better suited but we'll use the maximum likelihood because it is simpler.

Determine the following estimates.

$$P(F = 1|C = 0), P(F = 1|C = 1),$$
$$P(G = 1|C = 0), P(G = 1|C = 1),$$
$$P(F = 0, G = 0|C = 0), P(F = 0, G = 1|C = 0),$$
$$P(F = 1, G = 0|C = 0), P(F = 1, G = 1|C = 0),$$
$$P(F = 0, G = 0|C = 1), P(F = 0, G = 1|C = 1),$$
$$P(F = 1, G = 0|C = 1), P(F = 1, G = 1|C = 1).$$

$$P(F = 1|C = 0) = \frac{6}{17} = 0.3529 \qquad P(F = 1|C = 1) = \frac{143}{196} = 0.7296$$

$$P(G = 1|C = 0) = \frac{28}{51} = 0.5490 \qquad P(G = 1|C = 1) = \frac{93}{196} = 0.4745$$

$$P(F = 0, G = 0|C = 0) = \frac{5}{17} = 0.2941 \qquad P(F = 0, G = 1|C = 0) = \frac{6}{17} = 0.3529$$

$$P(F = 1, G = 0|C = 0) = \frac{8}{51} = 0.1569 \qquad P(F = 1, G = 1|C = 0) = \frac{10}{51} = 0.1961$$

$$P(F = 0, G = 0|C = 1) = \frac{1}{7} = 0.1429 \qquad P(F = 0, G = 1|C = 1) = \frac{25}{196} = 0.1276$$

$$P(F = 1, G = 0|C = 1) = \frac{75}{196} = 0.3827 \qquad P(F = 1, G = 1|C = 1) = \frac{17}{49} = 0.3469$$

Model $M_1$ for e-mail does not consider any feature. So $P(C)$ can be used to estimate the probability that the next e-mail will be spam or not. We will write that as $P(C|M_1)$.

b. Model $M_2$ considers only feature $F$ to predict whether the next e-mail will be spam or not.
Use Bayes rule and the probability estimates determined in the previous question to determine an estimate for $P(C|M_2) = P(C|F)$.

Bayes $\qquad P(C = 1|F) = \dfrac{P(C = 1)P(F|C = 1)}{P(F)}$,

where $\qquad\qquad P(F) = P(C = 0)P(F|C = 0) + P(C = 1)P(F|C = 1)$.

We get $\qquad P(C = 1|F = 0) = \dfrac{53}{86} = 0.6163$

and $\qquad P(C = 1|F = 1) = \dfrac{143}{161} = 0.8882$.

Model $M_3$ considers feature $G$ only and model $M_4$ considers both $F$ and $G$. Model $M_5$ also considers both $F$ and $G$ but assumes that $F$ and $G$ are independent given the classification $C$.

c.   Use Bayes rule again to show how you would calculate $P(C|M_5)$.

Bayes $\qquad\qquad P(C = 1|M_5) = \dfrac{P(C = 1)P(F|C = 1)P(G|C = 1)}{P(F, G)}$,

where $\qquad\qquad P(F, G) = P(C = 0)P(F|C = 0)P(G|C = 0)+$
$$+ P(C = 1)P(F|C = 1)P(G|C = 1).$$

We get $\qquad P(C = 1|F = 0, G = 0) = \dfrac{92803}{142391} = 0.6517$

$P(C = 1|F = 0, G = 1) = \dfrac{83793}{144161} = 0.5812$

$P(C = 1|F = 1, G = 0) = \dfrac{250393}{277441} = 0.9025$

and $\qquad P(C = 1|F = 1, G = 1) = \dfrac{75361}{86337} = 0.8729$.

d.   The models $M_1, M_2, \ldots, M_5$ all have a certain number of free parameters. Determine the number of free parameters for each of the five models.

Model 1: No features so we consider the joint probability $P\{C\}$ only. Thus we have one free parameter, for instance $P\{C = 1\}$.
Answer: 1 free parameter

Model 2: The joint probability is $P\{C, F\} = P\{C\}P\{F|C\}$. The free parameters are $P\{C = 1\}$, $P\{F = 1|C = 0\}$, $P\{F = 1|C = 1\}$.
Answer: 3 free parameters

Model 3: The joint probability is $P\{C, G\} = P\{C\}P\{G|C\}$. The free parameters are $P\{C = 1\}$, $P\{G = 1|C = 0\}$, $P\{G = 1|C = 1\}$.
Answer: 3 free parameters

Model 4: The joint probability is $P\{C, F, G\} = P\{C\}P\{F|C\}P\{G|C, F\}$. The free parameters are $P\{C = 1\}$, $P\{F = 1|C = 0\}$, $P\{F = 1|C = 1\}$, $P\{G = 1|C = 0, F = 0\}$, $P\{G = 1|C = 0, F = 1\}$, $P\{G = 1|C = 1, F = 0\}$, $P\{G = 1|C = 1, F = 1\}$.
Answer: 7 free parameters

Model 5: The joint probability is $P\{C, F, G\} = P\{C\}P\{F|C\}P\{G|C\}$. The free parameters are $P\{C = 1\}$, $P\{F = 1|C = 0\}$, $P\{F = 1|C = 1\}$, $P\{G = 1|C = 0\}$, $P\{G = 1|C = 1\}$.
Answer: 5 free parameters

e.  Given the training set of the 247 e-mail as shown in the table above, which of the five models would you prefer? Use an MDL argument in your answer.

HINT: You will need to calculate an estimate for the email entropy for each model. For model $M_1$ you make an estimate of $H(C)$ using the maximum likelihood estimate $P(C = 1) = 0.7935$. Likewise you calculate for $M_2$ the entropy $H(C|F)$ and thus you'll need to compute $P(C, F)$. For $M_3$ you must compute the entropy $H(C|G)$; for $M_4$ you calculate $H(C|F, G)$ and for $M_5$ also $H(C|F, G)$ although this will be a different calculation than for $M_4$.

---

As Stochastic complexity we use the formula:

$$\mathrm{sc} = \#\mathrm{param}/2 \log N + N H(C|\text{relevant params}),$$

Where N is the number of e-mails (247).

---

First we compute the five entropies using the general formula

$$H(C|X) = \sum_{x \in \mathcal{X}} P_X(x) H(C|X = x)$$

$$H(C|X = x) = \sum_{c=0}^{1} -P_{C|X}(c|x) \log_2 P_{C|X}(c|x)$$

$$P_{C|X}(c|x) = \frac{P_{C,X}(c, x)}{P_X(x)}$$

$$P_X(x) = \sum_{c'=0}^{1} P_{C,X}(c', x)$$

$H(M_1) = H(C) = 0.7347,$       $X$ is empty, so $P(C, X) = P(C)$

$H(M_2) = H(C|F) = 0.6639,$     $X = F$, so $P(C, X) = P(C)P(F|C)$

$H(M_3) = H(C|G) = 0.7321,$     $X = G$, so $P(C, X) = P(C)P(G|C)$

$H(M_4) = H(C|F, G) = 0.6614,$   $X = F, G$, so $P(C, X) = P(C)P(F|C)P(G|C, F)$

$H(M_5) = H(C|F, G) = 0.6615,$   $X = F, G$, so $P(C, X) = P(C)P(F|C)P(G|C)$

---

The stochastic complexities are:

$$\mathrm{sc}(M_1) = 1 \times \frac{1}{2} \log_2 N + N H(C) = 185.444$$

$$\mathrm{sc}(M_2) = 3 \times \frac{1}{2} \log_2 N + N H(C|F) = 175.894$$

$$\mathrm{sc}(M_3) = 3 \times \frac{1}{2} \log_2 N + N H(C|G) = 192.743$$

$$\mathrm{sc}(M_4) = 7 \times \frac{1}{2} \log_2 N + N H(C|F, G) = 191.175$$

$$\mathrm{sc}(M_5) = 5 \times \frac{1}{2} \log_2 N + N H(C|F, G) = 183.259$$

Based on this we should prefer $M_2$, which has the lowest stochastic complexity.

Note that eventually, we would prefer $M4$ because for very long sequences the entropy is the most important term and $H(M_4)$ is the smallest of the five.

## Appendix: formula sheet

The *Bernoulli distribution* is a discrete distribution having two possible outcomes labeled by $x = 0$ and $x = 1$ in which $x = 1$ ("success") occurs with probability $\theta$ and $x = 0$ ("failure") occurs with probability $1 - \theta$. It therefore has probability function

$$p(x|\theta) = \theta^x (1 - \theta)^{1-x} \tag{A.1}$$

The *Gaussian distribution* with mean $\mu$ and variance $\sigma^2$ is defined as

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

---

Points that can be scored per question:

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

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

Question 3:   a) 3 points; b) 3 points. Total 6 points.

Question 4:   a) 3 points; b) 3 points. Total 6 points.

Question 5:   a) 1 point; b) 1 point; c) 2 points; d) 2 points; e) 2 points. Total 8 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.