

# Tutorial 2:

## Exercises for Block I

### Exercise 2: Maximum Density

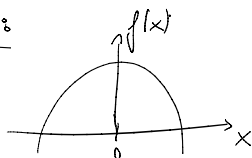
Consider a probability density  $p_x(x)$  defined over a continuous variable  $x$ , and suppose that we make a nonlinear change of variable using  $x = g(y)$ , so that the density transforms according to

$$p_y(y) = p_x(g(y)) |g'(y)| \quad (1)$$

- i) By differentiating 1, show that the location  $\hat{y}$  of the maximum of the density in  $y$  is not in general related to the location  $\hat{x}$  of the maximum of the density over  $x$  by the simple functional relation  $\hat{x} = g(\hat{y})$  as a consequence of the Jacobian factor. This shows that the maximum of a probability density (in contrast to a simple function) is dependent on the choice of variable.
- ii) Verify that, in the case of a linear transformation, the location of the maximum transforms in the same way as the variable itself.

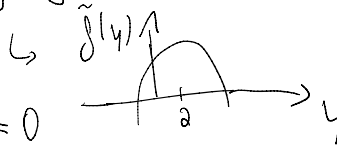
i) Normal functions:

• Take  $f(x)$ :



• Change to new variable  $y$  using  $x = g(y)$ . e.g.:  $g(a) = a - 2$

$$\hookrightarrow \tilde{f}(y) = f(g(y))$$



• Say  $f(x)$  has a mode (maximum) at  $\hat{x}$ , s.t.  $f'(\hat{x}) = 0$

$\hookrightarrow$  What is the mode  $\hat{y}$  of  $\tilde{f}(y)$ ?

$$\bullet \tilde{f}'(\hat{y}) \stackrel{!}{=} 0$$

assuming  $g'(\hat{y}) \neq 0$

$$\Leftrightarrow f'(g(\hat{y})) \cdot g'(\hat{y}) = 0 \Leftrightarrow f'(g(\hat{y})) = 0$$

$$\hookrightarrow \text{But we also know } f'(\hat{x}) = 0$$

$$\hookrightarrow \hat{x} = g(\hat{y})$$

Density:

• Consider  $p_x(x)$  with mode  $\hat{x}$ , change of variables:  $x = g(y)$

$$\hookrightarrow p_y(y) = p_x(g(y)) \cdot |g'(y)|$$

• Perform derivative w.r.t.  $y \Rightarrow$  it must be 0 for  $\hat{y}$ !

$$\begin{aligned} \frac{dp_y(y)}{dy} &= \frac{dp_x(g(y))}{dy} \cdot |g'(y)| + p_x(g(y)) \cdot \frac{d|g'(y)|}{dy} \\ &= \frac{dp_x(g(y))}{dy} \cdot \frac{dg(y)}{dy} \cdot |g'(y)| + p_x(g(y)) \cdot \frac{d|g'(y)|}{dy} \end{aligned}$$

$$= \frac{dp_x(g(y))}{dg(y)} \cdot \frac{dg(y)}{dy} \cdot |g'(y)| + \dots$$

↳ Since  $\frac{dp_x(x)}{dx} \Big|_{x=\hat{x}} = 0$ , if we chose  $\hat{x} = g(\hat{y})$ , then  $\frac{dp_x(g(y))}{dg(y)} \Big|_{y=\hat{y}} = 0$ .

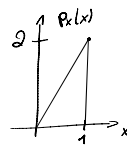
↳ What about second term? → Doesn't go away!

$$\hookrightarrow p_x(g(y)) \cdot \frac{d|g'(y)|}{dy}$$

ii) Look at this: If  $g(y) = a \cdot y + b$ , then  $|g'(y)| = |a|$  and thus  $\frac{d|g'(y)|}{dy} = 0 \Rightarrow \hat{x} = g(\hat{y})!$

Example:

Let  $p_x(x) = 2x$ ,  $x \in [0, 1]$ :



$$\hookrightarrow \hat{x} = 1$$

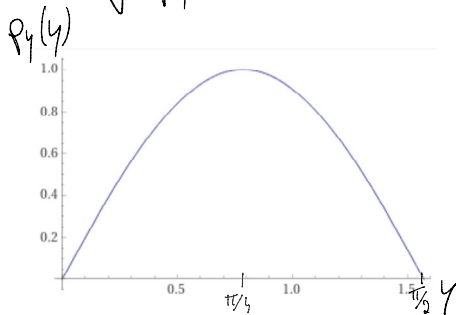
Now, set  $x = \sin(y)$ .

$$\hookrightarrow p_y(y) = 2 \cdot \sin(y) \cdot |\cos(y)|, y \in [0, \frac{\pi}{2}]$$

$$= 2 \sin(y) \cdot \cos(y) = \sin(2y)$$

$$\sin(\alpha + \beta) = \sin(\alpha)\cos(\beta) + \sin(\beta)\cos(\alpha)$$

Plot of  $p_y(y)$ :



↳ Obviously,  $\hat{x} \neq \sin(\hat{y})$ .

$$\mathcal{D} = \{x_i\}_{i=1}^N$$

### Exercise 7: Maximum likelihood estimates

Verify by setting the derivatives of the log likelihood

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

with respect to  $\mu$  and  $\sigma^2$  equal to zero:

i)  $\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$  (see [1] 1.55)

ii)  $\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$  (see [1] 1.56)

iii)  $\sigma^2 \propto \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$

$$\begin{cases} p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ p(x_1, \dots, x_N|\mu, \sigma^2) = \prod_{i=1}^N p(x_i|\mu, \sigma^2) \end{cases}$$

ii)  $\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$  (see [1] 1.56)

$$\begin{aligned}
 \text{i)} \quad \frac{\partial \ln(p(x|\mu, \sigma^2))}{\partial \mu} &= -\frac{1}{\sigma^2} \sum (x_n - \mu) \cdot (-1) \\
 &= \frac{1}{\sigma^2} \sum (x_n - \mu) \stackrel{!}{=} 0 \\
 &\Leftrightarrow \frac{1}{\sigma^2} \sum x_n = \frac{1}{\sigma^2} \cdot N \cdot \mu \\
 &\Leftrightarrow \mu_{ML} = \frac{1}{N} \sum x_n
 \end{aligned}$$

$$\text{ii)} \quad \frac{\partial \ln(p(x|\mu, \sigma^2))}{\partial \sigma^2} = \frac{1}{2} (\sigma^2)^{-2} \cdot \sum (x_n - \mu)^2 - \frac{N}{2} \cdot \frac{1}{\sigma^2} \stackrel{!}{=} 0$$

$$\begin{aligned}
 &\Leftrightarrow \frac{1}{N} \sum (x_n - \mu)^2 = \sigma^2 \stackrel{!}{=} 0 \\
 &\Leftrightarrow \sigma_{ML}^2 = \frac{1}{N} \sum (x_n - \mu_{ML})^2
 \end{aligned}$$

## Exercise 8: True variance

Suppose that the variance of a Gaussian is estimated using  $\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$  (see [1] 1.56) but with the maximum likelihood estimate  $\mu_{ML}$  replaced with the true value  $\mu$  of the mean. Show that this estimator has the property that its expectation is given by the true variance  $\sigma^2$  unbiasedness

$$\begin{aligned}
 &\mathbb{E}_{\{x_n\}} \left[ \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2 \right] \\
 &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{x_n} [x_n^2 - 2x_n\mu + \mu^2] \\
 &= \frac{1}{N} \sum_{n=1}^N \left( (\mu^2 + \sigma^2) - 2\mu\mu + \mu^2 \right) \\
 &= \sigma^2
 \end{aligned}$$

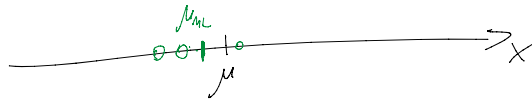
For any two r.v.'s  $X, Y$ :  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] + \text{Cov}(X, Y)$   
 $\hookrightarrow \mathbb{E}[X^2] = \mu^2 + \sigma^2$

Bonus: What if we use  $\mu_{ML}$ ?

$$\hookrightarrow \mathbb{E}_{\{x_n\}} \left[ \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2 \right] = \dots = \frac{N-1}{N} \sigma^2$$

$$\hookrightarrow \mathbb{E} \left[ \frac{N}{N-1} \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2 \right] = \dots = \frac{N}{N-1} \frac{N-1}{N} \sigma^2$$

↳ Intuition:



## Exercise 10: Misclassification bound

Consider two nonnegative numbers  $a$  and  $b$ , and show that, if  $a \leq b$ , then  $a \leq \sqrt{ab}$ . Use this result to show that, if the decision regions of a two-class classification problem with classes  $C_1, C_2$  are chosen to minimize the probability of misclassification, this probability will satisfy

$$p(\text{error}) \leq \int \{p(x, C_1)p(x, C_2)\}^{1/2} dx$$

i) Start with  $a \leq b$ . Since  $\sqrt{\cdot}$  is monotonic,

$$\Leftrightarrow \sqrt{a} \leq \sqrt{b}$$

$$\stackrel{\cdot \sqrt{a}}{\Leftrightarrow} a \leq \sqrt{a \cdot b}$$

ii)  $p(\text{error}) = P(\hat{Y} \neq Y)$  →  $\begin{cases} 1 \\ 0 \end{cases}$  if  $\hat{Y} = -1$ , else

$$= \int P(Y=1|x) \mathbb{1}_{\hat{Y}=-1} p(x) dx + \int P(Y=-1|x) \mathbb{1}_{\hat{Y}=1} p(x) dx$$

Since we're Bayes-optimal, we have that:

• if  $\hat{Y} = -1$  (first integral), then:

$$P(Y=1|x) \leq P(Y=-1|x)$$

• if  $\hat{Y} = 1$  (second integral), then:

$$P(Y=-1|x) \leq P(Y=1|x)$$

↳ Using  $(a \leq b \Rightarrow a \leq \sqrt{ab})$  we write:

$$p(\text{error}) \leq \int \{P(Y=1|x) \cdot P(Y=-1|x)\}^{1/2} \cdot \mathbb{1}_{\hat{Y}=-1} \cdot p(x) dx + \int \{P(Y=-1|x) \cdot P(Y=1|x)\}^{1/2} \cdot \mathbb{1}_{\hat{Y}=1} \cdot p(x) dx$$

$$\begin{aligned}
& + \int \{ P(Y=-1|x) \cdot P(Y=1|x) \}^{1/2} \cdot \frac{1}{\sqrt{2}} \cdot p(x) dx \\
& = \int \{ P(Y=1|x) \cdot P(Y=-1|x) \}^{1/2} \underbrace{\left( \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} \right)}_{=1} p(x) dx \\
& = \int \{ P(Y=1|x) \cdot p(x) \cdot P(Y=-1|x) \cdot p(x) \}^{1/2} dx \\
& = \int \{ P(Y=1, x) - P(Y=-1, x) \}^{1/2} dx
\end{aligned}$$


---

### Exercise 15: Decision boundary

Consider the following decision rule for a two-category one-dimensional problem: Decide  $C_1$  if  $x > \theta$ ; otherwise decide  $C_2$ .

i) Show that the probability of error for this rule is given by

$$P(\text{error}) = P(C_1) \int_{-\infty}^{\theta} p(x|C_1) dx + P(C_2) \int_{\theta}^{\infty} p(x|C_2) dx$$

ii) By differentiating, show that a necessary condition to minimize  $P(\text{error})$  is that  $\theta$  satisfy

$$p(\theta|C_1)P(C_1) = p(\theta|C_2)P(C_2) \quad (8)$$

iii) Does equation 8 define  $\theta$  uniquely?

iv) Give an example where a value of  $\theta$  satisfying the equation actually maximizes the probability of error.

$$i) P(\text{error}|x) = \begin{cases} P(C_2|x), & \text{if we decide } C_1 \\ P(C_1|x), & \text{if we decide } C_2 \end{cases}$$

$$\begin{aligned}
P(\text{error}) &= \int_{-\infty}^{\theta} P(\text{error}|x) p(x) dx \\
&= \int_{-\infty}^{\theta} P(C_2|x) p(x) dx \quad \xrightarrow{P(A|B) = \frac{P(A) \cdot P(B)}{P(B)}} \text{"we decide } C_2\text{"} \\
&\quad + \int_{\theta}^{\infty} P(C_1|x) p(x) dx \quad \text{"we decide } C_1\text{"} \\
&= P(C_2) \int_{-\infty}^{\theta} p(x|C_2) dx \\
&\quad + P(C_1) \int_{\theta}^{\infty} p(x|C_1) dx
\end{aligned}$$

ii) Derivative: Reminder: 2. Fundamental Theorem of Calculus:  $\frac{d}{dx} \int_a^x f(t) dt = f(x)$

$$\frac{dP(\text{error})}{d\theta} = P(C_2) \cdot p(\theta|C_2) - P(C_1) \cdot p(\theta|C_1) \stackrel{!}{=} 0$$

$$\Leftrightarrow P(C_1) p(\theta|C_1) = P(C_2) \cdot p(\theta|C_2)$$

iii)  $N_0$  can be true over a range of  $\theta$ .

iv) If  $p(x|C_1) \sim N(1,1)$  and  $p(x|C_2) \sim N(-1,1)$  and

$P(C_1) = P(C_2) = 1/2$  then, although  $\theta = 0$  satisfies the cond.,  $p(\text{error})$  has a maximum at 0.

