# Machine Learning 2 — Homework 2

Maurice Frank
11650656
maurice.frank@posteo.de

September 16, 2019

## Problem 1.

**1.**

We have the discrete random variables $X, Y, Z$. Mutual information than is:

$$I(X;Y) = \mathcal{KL}(p(x,y)||p(x)p(y))$$
$$= \mathcal{H}(X) - \mathcal{H}(X|Y)$$

The conditional mutual information is:

$$I(X;Y|Z) = \mathbb{E}_{p(z)}[\mathcal{KL}(p(x,y|z)||p(x|y)p(y|z))]$$
$$= \mathcal{H}(X|Z) - \mathcal{H}(X|Y,Z)$$

We see that the conditional mutual information measures the expected mutual information between $X$ and $Y$ given $Z$.

**2.**

We have $x, y, z \in \{0, 1\}$ with $p(x, y, z)$. First we will write down $p(x, y)$ (see Table 1), $p(x)$ and $p(y)$ (see Table 2).

| X | Y | p(x,y) |
|---|---|--------|
| 0 | 0 | 0.336 |
| 0 | 1 | 0.264 |
| 1 | 0 | 0.256 |
| 1 | 1 | 0.144 |

Table 1: Marginalizing Z.

| Y X | 0 | 1 | p(x) |
|---|---|---|---|
| 0 | 0.336 | 0.264 | 0.6 |
| 1 | 0.256 | 0.144 | 0.4 |
| p(y) | 0.592 | 0.408 | 1 |

Table 2: Marginalizing X and Y.

| y x | 0 | 1 | p(x\|z=0) |
|---|---|---|---|
| 0 | 0.4 | 0.1 | 0.5 |
| 1 | 0.4 | 0.1 | 0.5 |
| p(y\|z=0) | 0.8 | 0.2 | 1.0 |

Table 3: p(x,y | z=0) and p(x | z=0), p(y | z=0)

$$I(X;Y) = \mathcal{KL}(p(x,y)||p(x)p(y))$$
$$= -\sum_{x,y} p(x,y) \ln\left(\frac{p(x)p(y)}{p(x,y)}\right)$$
$$= 3.197 \cdot 10^{-3}$$
$$> 0$$

As the mutual information between $X$ and $Y$ is bigger than zero (it is symmetric!) we showed that having knowledge about one of the values we gain knowledge about the possible distribution of values of the second variable. We see this with $I(X;Y) = \mathcal{H}(X) - \mathcal{H}(X|Y) > 0 \implies \mathcal{H}(X|Y) < \mathcal{H}(X)$, the entropy of one given the other is lower than without this conditional information.

**3.**

See Table 3 and Table 4 for intermediate calculation tables.

| y x | 0 | 1 | p(x\|z=1) |
|---|---|---|---|
| 0 | 0.277 | 0.415 | 0.692 |
| 1 | 0.123 | 0.185 | 0.308 |
| p(y\|z=1) | 0.4 | 0.6 | 1.0 |

Table 4: p(x,y | z=1) and p(x | z=1), p(y | z=1)

| $x$ | $p(x)$ | $z$ | $p(z\|x)$ | $y$ | $p(y\|z)$ | $p(x) \cdot p(z\|x) \cdot p(y\|z)$ |
|---|---|---|---|---|---|---|
| 0 | 0.6 | 0 | 0.4 | 0 | 0.8 | 0.192 |
| 0 | 0.6 | 0 | 0.4 | 1 | 0.2 | 0.048 |
| 0 | 0.6 | 1 | 0.6 | 0 | 0.4 | 0.144 |
| 0 | 0.6 | 1 | 0.6 | 1 | 0.6 | 0.216 |
| 1 | 0.4 | 0 | 0.6 | 0 | 0.8 | 0.192 |
| 1 | 0.4 | 0 | 0.6 | 1 | 0.2 | 0.048 |
| 1 | 0.4 | 1 | 0.4 | 0 | 0.4 | 0.064 |
| 1 | 0.4 | 1 | 0.4 | 1 | 0.6 | 0.096 |

Table 5: Computation of $p(x) \cdot p(z|x) \cdot p(y|z)$



Figure 1: The directed graph to the factorization $p(x)p(z|x)p(y|z)$

$$
\begin{aligned}
I(X;Y|Z) = \ & p(z=0) \cdot \mathcal{KL}(p(x,y|z=0)||p(x|z=0)p(y|z=0)) \\
& + p(z=1) \cdot \mathcal{KL}(p(x,y|z=1)||p(x|z=1)p(y|z=1)) \\
= & - \sum_{z \in Z} p(Z=z) \cdot \sum_{x,y} p(x,y|z) \cdot \log\left(\frac{p(x|z)p(y|z)}{p(x,y|z)}\right) \\
= & -(0.48 \cdot 0 + 0.52 \cdot 0) \\
= & \ 0
\end{aligned}
$$

That now the conditional mutual information is zero tells us that given that we know about the value of $Z$ than having information about $X$ or $Y$ will not tell us anything about the respective third variable. $I(X;Y|Z) = 0 \implies \mathcal{H}(X|Z) = \mathcal{H}(X|Y,Z) \wedge \mathcal{H}(Y|Z) = \mathcal{H}(Y|X,Z)$.

**4.**

For the Computation of $p(x) \cdot p(z|x) \cdot p(y|z) = p(x,y,z)$ see Table 5. The directed graph is show in Figure 1.

## Problem 2.

We strictly plot not the permutated directed graphs. See Figures 3 until 7 for the plotted clusters. Any not written conditional relationship is independent. As for the definition of permutation consider the graph in Figure 2. This is the same graph as in Figure 4 as we can permute the variable names: $X \to Y \wedge Z \to X \wedge Y \to Z$. This is a permutation and proves that those two clusters are the same.
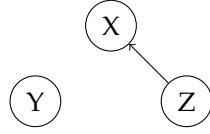
Figure 2: Example DAG



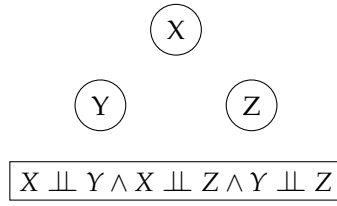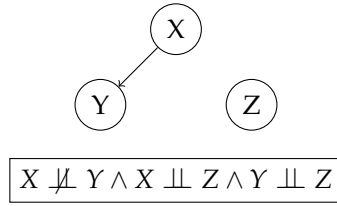$$X \perp\!\!\!\perp Y \wedge X \perp\!\!\!\perp Z \wedge Y \perp\!\!\!\perp Z$$

Figure 3: Cluster I



$$X \not\perp\!\!\!\perp Y \wedge X \perp\!\!\!\perp Z \wedge Y \perp\!\!\!\perp Z$$

Figure 4: Cluster II



$$X \not\perp\!\!\!\perp Z \wedge X \not\perp\!\!\!\perp Y \wedge Y \not\perp\!\!\!\perp Z|\emptyset \wedge Y \perp\!\!\!\perp Z|X$$

Figure 5: Cluster III



$$X \not\perp\!\!\!\perp Y \wedge X \not\perp\!\!\!\perp Z \wedge Y \perp\!\!\!\perp Z \wedge Y \not\perp\!\!\!\perp Z|X$$
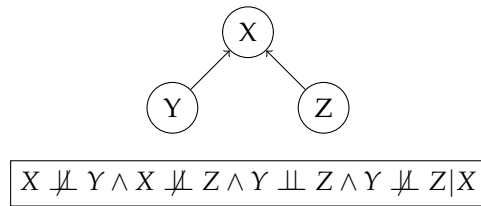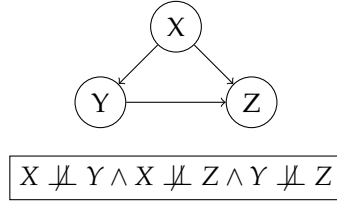
Figure 6: Cluster IV

Figure 7: Cluster V

## Problem 3.

### 1.

We have given $p(x) = \mathcal{N}(x|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $q(x) = \mathcal{N}(x|\boldsymbol{m}, \boldsymbol{L})$. $x \in \mathbb{R}^k$. We are using the result of Section 2. in this computation.

$$\ln q(x) = -\frac{k}{2} \ln (2\pi) - \frac{1}{2} \ln |\boldsymbol{L}| - \frac{1}{2}(x - \boldsymbol{m})^T \boldsymbol{L}^{-1}(x - \boldsymbol{m})$$

$$\mathcal{KL}(p||q) = -\int p(x) \ln \left(\frac{q(x)}{p(x)}\right) dx$$

$$= -\int p(x) \left[\ln q(x) - \ln p(x)\right] dx$$

$$= \int p(x) \ln p(x) dx - \int p(x) \ln q(x) dx$$

$$= -\mathcal{H}(p(X)) - \int p(x) \ln q(x) dx$$

$$= -\mathcal{H}(p(X)) + \frac{k}{2} \ln (2\pi) + \frac{1}{2} \ln |\boldsymbol{L}| + \frac{1}{2} \int p(x)(x - \boldsymbol{m})^T \boldsymbol{L}^{-1}(x - \boldsymbol{m}) dx$$

$$= -\frac{k}{2} - \frac{k}{2} \ln (2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}| + \frac{k}{2} \ln (2\pi) + \frac{1}{2} \ln |\boldsymbol{L}| + \frac{1}{2} \int p(x)(x - \boldsymbol{m})^T \boldsymbol{L}^{-1}(x - \boldsymbol{m}) dx$$

$$= -\frac{k}{2} + \frac{1}{2} \ln \frac{|\boldsymbol{L}|}{|\boldsymbol{\Sigma}|} + \frac{1}{2} \int p(x)(x - \boldsymbol{m})^T \boldsymbol{L}^{-1}(x - \boldsymbol{m}) dx$$

$$= -\frac{k}{2} + \frac{1}{2} \ln \frac{|\boldsymbol{L}|}{|\boldsymbol{\Sigma}|} + \frac{1}{2} \mathbb{E}[(x - \boldsymbol{m})^T \boldsymbol{L}^{-1}(x - \boldsymbol{m})]_{p(x)} \qquad \text{(using (1))}$$

$$= -\frac{k}{2} + \frac{1}{2} \ln \frac{|\boldsymbol{L}|}{|\boldsymbol{\Sigma}|} + \frac{1}{2} [(\boldsymbol{\mu} - \boldsymbol{m})^T \boldsymbol{L}^{-1}(\boldsymbol{\mu} - \boldsymbol{m}) + \text{Tr}\,(\boldsymbol{L}^{-1}\boldsymbol{\Sigma})]$$

$$= \frac{1}{2} \left[\ln \frac{|\boldsymbol{L}|}{|\boldsymbol{\Sigma}|} - k + (\boldsymbol{\mu} - \boldsymbol{m})^T \boldsymbol{L}^{-1}(\boldsymbol{\mu} - \boldsymbol{m}) + \text{Tr}\,(\boldsymbol{L}^{-1}\boldsymbol{\Sigma})\right]$$

**2.**

$$\ln p(x) = -\frac{k}{2}\ln(2\pi) - \frac{1}{2}\ln|\boldsymbol{\Sigma}| - \frac{1}{2}(x-\mu)^T\boldsymbol{\Sigma}^{-1}(x-\mu)$$

$$\mathcal{H}(p) = -\int p(x)\ln p(x)dx$$

$$= \frac{k}{2}\ln(2\pi) + \frac{1}{2}\ln|\boldsymbol{\Sigma}| + \frac{1}{2}\mathbb{E}_{p(x)}\left[(x-\mu)^T\boldsymbol{\Sigma}^{-1}(x-\mu)\right]$$

$$= \frac{k}{2}\ln(2\pi) + \frac{1}{2}\ln|\boldsymbol{\Sigma}| + \frac{1}{2}\left[(\mu-\mu)^T\boldsymbol{\Sigma}^{-1}(\mu-\mu) + \mathrm{Tr}\left(\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}\right)\right]$$
$$\text{(using (1))}$$

$$= \frac{k}{2}\ln(2\pi) + \frac{1}{2}\ln|\boldsymbol{\Sigma}| + \frac{1}{2}\mathrm{Tr}\left(\mathbb{1}_k\right)$$

$$= \frac{k}{2}\ln(2\pi) + \frac{1}{2}\ln|\boldsymbol{\Sigma}| + \frac{k}{2}$$