

## Exercise Sheet 3

### Exercise 1: Maximum Entropy Distributions (5 + 10 + 5 P)

Consider a discrete random variable  $X$  defined over positive numbers  $\mathbb{N}$ , each of them occurring with respective probability  $\Pr(X = k)$ . The entropy  $H(X)$  is given by

$$H(X) = - \sum_{k=0}^{\infty} \Pr(X = k) \log \Pr(X = k)$$

We would like to find the probability distribution that maximizes the entropy  $H(X)$  under the expectation constraint

$$\mathbb{E}[X] = \sum_{k=0}^{\infty} k \Pr(X = k) = 1.$$

and forcing probabilities to be strictly positive, i.e.  $\Pr(X = k) > 0$  for all  $k \in \mathbb{N}$ . The latter can be enforced using the reparameterization  $\Pr(X = k) = \exp(s_k)$ . To enforce that probability values sum to one, we can impose the additional constraint

$$\sum_{k=0}^{\infty} \Pr(X = k) = 1.$$

The problem of finding the maximum-entropy distribution can therefore be modeled as a constrained optimization problem, and such problems can be solved using the method of Lagrange multipliers.

- (a) Write the Lagrange function  $\mathcal{L}((s_k)_{k \in \mathbb{N}}, \lambda_1, \lambda_2)$  corresponding to the constrained optimization problem above, where  $\lambda_1, \lambda_2 \in \mathbb{R}$  are used to incorporate the sum-to-one and the expectation constraints respectively.
- (b) Show that the probability distribution that maximizes the objective  $H(X)$  has the form:

$$\Pr(X = k) = \alpha \cdot \beta^k$$

- (c) Show that  $\alpha = 0.5$  and  $\beta = 0.5$ , i.e. the probability distribution is geometric with failure and success probabilities 0.5, and  $k$  denoting the number of failures before success. (Hint: you can make use of the geometric series and Gabriel's staircase series).

### Exercise 2: Independent Components in Two Dimensions (5 + 10 + 10 P)

High entropy can be seen as the result of superposing many independent low-entropy sources. Independent component analysis (ICA) aims to recover the independent sources from the data by finding projections of the data that have low entropy, i.e. that diverge the most from a Gaussian probability distribution.

In the following, we consider a simple two-dimensional problem where we are given a joint probability distribution  $p(x, y) = p(x)p(y|x)$  with

$$p(x) \sim \mathcal{N}(0, 1),$$
$$p(y|x) = \frac{1}{2} \delta(y - x) + \frac{1}{2} \delta(y + x),$$

where  $\delta(\cdot)$  denotes the Dirac delta function. A useful property of linear component analysis for two-dimensional probability distributions is that the set of all possible directions to look for in  $\mathbb{R}^2$  is directly given by

$$\left\{ \begin{pmatrix} \cos \theta \\ \sin \theta \end{pmatrix}, \quad 0 \leq \theta < 2\pi \right\}.$$

The projection of the random vector  $(x, y)$  on a particular component can therefore be expressed as a function of  $\theta$ :

$$z(\theta) = x \cos(\theta) + y \sin(\theta).$$

As a result, ICA in the two-dimensional space is reduced to finding the values of the parameter  $\theta \in [0, 2\pi[$  that maximize a certain objective  $J(z(\theta))$ .

- Sketch* the joint probability distribution  $p(x, y)$ , along with the projections  $z(\theta)$  of this distribution for angles  $\theta = 0$ ,  $\theta = \pi/8$  and  $\theta = \pi/4$ .
- Find* the principal components of  $p(x, y)$ . That is, find the values of the parameter  $\theta \in [0, 2\pi[$  that maximize the variance of the projected data  $z(\theta)$ .
- Find* the independent components of  $p(x, y)$ , more specifically, find the values of the parameter  $\theta \in [0, 2\pi[$  that maximize the non-Gaussianity of  $z(\theta)$ . We use as a measure of non-Gaussianity the excess kurtosis defined as

$$\text{kurt}[z(\theta)] = \frac{\mathbb{E}[(z(\theta) - \mathbb{E}[z(\theta)])^4]}{(\text{Var}[z(\theta)])^2} - 3.$$

### Exercise 3: Deriving a Special Case of FastICA (5 + 5 + 5 + 10 P)

We consider our input data  $x$  to be in  $\mathbb{R}^d$  and coming from some distribution  $p(x)$ . We assume we have applied as an initial step the centering and whitening procedures so that:

$$\mathbb{E}[x] = 0 \quad \text{and} \quad \mathbb{E}[xx^\top] = I.$$

To extract an independent component, we would like to find a unit vector  $w$  (i.e.  $\|w\|^2 = 1$ ) such that the excess kurtosis of the projected data:

$$\text{kurt}[w^\top x] = \frac{\mathbb{E}[(w^\top x - \mathbb{E}[w^\top x])^4]}{(\text{Var}[w^\top x])^2} - 3.$$

is maximized.

- Show* that for any  $w$  subject to  $\|w\|^2 = 1$ , the projection  $z = w^\top x$  has mean 0 and variance 1.
- Show using the method of Lagrange multipliers that the projection  $w$  that maximizes  $\text{kurt}[w^\top x]$  is a solution of the equation

$$\lambda w = \mathbb{E}[x \cdot (w^\top x)^3]$$

where  $\lambda$  can be resolved by enforcing the constraint  $\|w\|^2 = 1$ .

- The solution of the previous equation cannot be found analytically, and must instead be solved iteratively using e.g. the Newton's method. The Newton's method assumes that the equation is given in the form  $F(w) = 0$  and then defines the iteration as  $w^+ = w - J(w)^{-1}F(w)$  where  $w^+$  denotes the next iterate and where  $J$  is the Jacobian associated to the function  $F$ . *Show* that the Newton iteration in our case takes the form:

$$w^+ = w - (\mathbb{E}[3xx^\top(w^\top x)^2 - \lambda I])^{-1} \cdot (\mathbb{E}[x \cdot (w^\top x)^3] - \lambda w)$$

- Show* that when making the decorrelation approximation  $\mathbb{E}[xx^\top(w^\top x)^2] = \mathbb{E}[xx^\top] \cdot \mathbb{E}[(w^\top x)^2]$ , the Newton iteration further reduces to

$$w^+ = \gamma \cdot (\mathbb{E}[x \cdot (w^\top x)^3] - 3w)$$

where  $\gamma$  is some constant factor. (The constant factor does not need to be resolved since we subsequently apply the normalization step  $w^+ \leftarrow w^+/\|w^+\|$ .)

### Exercise 4: Programming (30 P)

Download the programming files on ISIS and follow the instructions.

# 1 Maximum Entropy Distribution

$$\begin{aligned}
 a) \quad \mathcal{L}((s_k)_{k \in \mathbb{N}}, \lambda_1, \lambda_2) &= - \sum_{k=0}^{\infty} e^{s_k} \cdot s_k + \lambda_1 \left( \sum_{k=0}^{\infty} k e^{s_k} - 1 \right) \\
 &\quad + \lambda_2 \left( \sum_{k=0}^{\infty} e^{s_k} - 1 \right) \\
 &= \sum_{k=0}^{\infty} e^{s_k} (-s_k) + \sum_{k=0}^{\infty} e^{s_k} (\lambda_1 k) - \lambda_1 + \sum_{k=0}^{\infty} e^{s_k} (\lambda_2) - \lambda_2 \\
 &= \sum_{k=0}^{\infty} e^{s_k} (\lambda_1 k + \lambda_2 - s_k) - \lambda_1 - \lambda_2
 \end{aligned}$$

$$1) \quad \frac{\partial}{\partial s_k} \mathcal{L}(\cdot) \stackrel{!}{=} 0 \quad \forall k \in \mathbb{N}$$

$$\frac{\partial}{\partial s_k} \mathcal{L}((s_k)_{k \in \mathbb{N}}, \lambda_1, \lambda_2) = e^{s_k} (\lambda_1 k + \lambda_2 - s_k) - e^{s_k} = 0$$

$$\Leftrightarrow e^{s_k} (\lambda_1 k + \lambda_2 - s_k - 1) = 0 \Leftrightarrow s_k = \lambda_1 k + \lambda_2 - 1$$

$$Pr(X=k) = e^{\lambda_1 k + \lambda_2 - 1} = e^{\lambda_2 - 1} e^{\lambda_1 k}$$

$$= \underbrace{e^{\lambda_2 - 1}}_{\alpha} \underbrace{e^{\lambda_1}}_{\beta}^k = \alpha \beta^k$$

$$c) \quad I \quad \frac{\partial \mathcal{L}(\cdot)}{\partial s_k} \stackrel{!}{=} 0 \Leftrightarrow Pr(X=k) = \alpha \beta^k$$

$$II \quad \frac{\partial \mathcal{L}(\cdot)}{\partial \lambda_2} \stackrel{!}{=} 0 \Leftrightarrow \sum_{k=0}^{\infty} Pr(X=k) = 1 \stackrel{(I)}{\Leftrightarrow} \sum_{k=0}^{\infty} \alpha \beta^k = 1$$

$$\Leftrightarrow \frac{\alpha}{1-\beta} = 1 \Leftrightarrow \alpha = 1-\beta$$

$$III \quad \frac{\partial \mathcal{L}(\cdot)}{\partial \lambda_1} \stackrel{!}{=} 0 \Leftrightarrow \sum_{k=0}^{\infty} k Pr(X=k) = 1 \stackrel{(I)}{\Leftrightarrow} \sum_{k=0}^{\infty} k \alpha \beta^k = 1$$

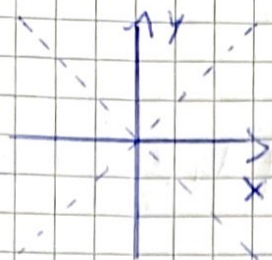
$$\Leftrightarrow \alpha \sum_{k=0}^{\infty} k \beta^k = 1 \Leftrightarrow \alpha \frac{\beta}{(1-\beta)^2} = 1 \stackrel{(II)}{\Leftrightarrow} \frac{(1-\beta)\beta}{(1-\beta)^2} = 1$$

$$\Leftrightarrow \beta = (1-\beta) \Leftrightarrow \beta = 0,5$$



## 2. Independent Components in Two Dimensions

a)



$p(x) \rightarrow$  sample over  $x$  according to gaussian dist.

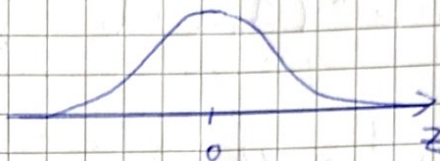
$p(y|x) \rightarrow$  "p of  $y$  given  $x$ "

$$p(y|x) = \frac{1}{2} \delta(y-x) + \frac{1}{2} \delta(y+x)$$

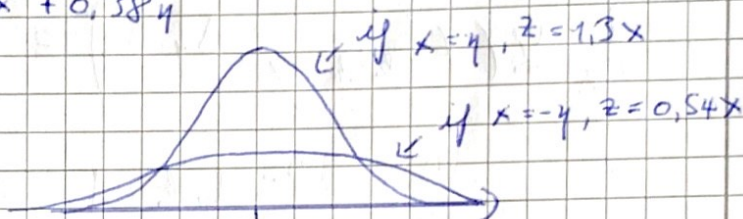
$y=x$                        $y=-x$

$\downarrow$   
50% of samples

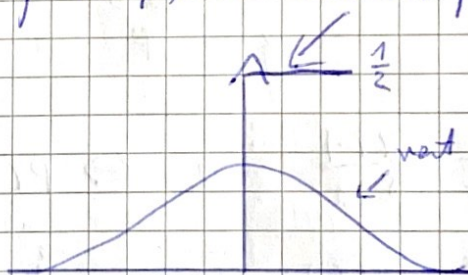
If  $\vartheta = 0: z(\vartheta) = x$



If  $\vartheta = \frac{1}{8}\pi: z = 0,92x + 0,38y$



If  $\vartheta = \frac{\pi}{4}: z = \frac{\sqrt{2}}{2}(x+y) \rightarrow$  if  $x=y, z = \frac{\sqrt{2}}{2} \cdot 2x = \sqrt{2}x$   
 $\rightarrow$  gaussian is "lower"  $\rightarrow$  more variance  
 $\rightarrow$  if  $x=-y, z=0 \rightarrow$  direct peak





$$\begin{aligned}
 1) E[z] &= E[x \cos \theta + y \sin \theta] = \underbrace{\cos \theta E[x]}_{=0} + \sin \theta E[y] \\
 &= \sin \theta E[E[y]] = \sin \theta E_x[0] = 0
 \end{aligned}$$

$$\text{Var}[z] = E[(z - \underbrace{E[z]}_{=0})^2] = E[z^2]$$

$$\begin{aligned}
 &= E[x^2 \cos^2 \theta + xy \cos \theta \sin \theta + y^2 \sin^2 \theta] \\
 &= \underbrace{\cos^2 \theta E[x^2]}_{=1} + \underbrace{\cos \theta \sin \theta E[xy]}_{=0} + \underbrace{\sin^2 \theta E[y^2]}_{=1} \\
 &= \cos^2 \theta + \sin^2 \theta = 1
 \end{aligned}$$

$$c) E[z(\theta)] = 0, \text{Var}[z(\theta)] = 1$$

$$kurt[z(\theta)] = \frac{E[(z(\theta) - 0)^4]}{1^2} - 3 = E[z(\theta)^4] - 3$$

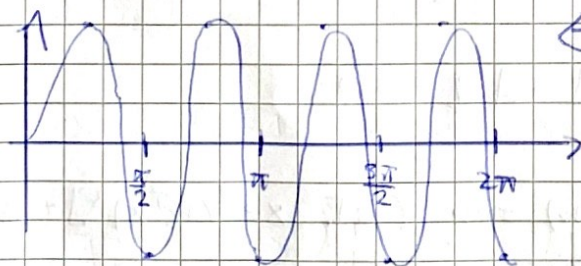
$$= E[(x \cos \theta + y \sin \theta)^4] - 3$$

$$= E[(x^2 \cos^2 \theta + 2xy \cos \theta \sin \theta + y^2 \sin^2 \theta)^2] - 3$$

$$= \cos^4 \theta E[x^4] + 4 \cos^2 \theta \sin^2 \theta E[x^2 y^2] + 4 \cos \theta \sin^3 \theta E[x y^3] + \sin^4 \theta E[y^4] - 3$$

$$= E[x^4] (\underbrace{\cos^2 \theta + \sin^2 \theta}_{=1}) + 4 \cos^2 \theta \sin^2 \theta - 3$$

$$= E[x^4] (1 + 1 - \frac{1 - \cos 4\theta}{2}) - 3 \rightarrow \cos 4\theta \text{ needs to be maximized!}$$



→ independent components:

$$\frac{\pi}{4}, \frac{3\pi}{4}, \frac{5\pi}{4}, \frac{7\pi}{4}$$



### 3 Deriving a Special Case of Fast SCA

$$a) \quad \underbrace{E[w^T x]}_{\text{mean of } w^T x} = \underbrace{w^T E[x]}_{=0} = 0$$

$$\begin{aligned} \text{Var}[w^T x] &= E[(w^T x)^2] = E[(w^T x w^T x)] E[w^T x x^T w] \\ &= \underbrace{w^T E[x x^T]}_{=I} w = w^T w = \|w\|^2 = 1 \end{aligned}$$

$$b) \quad \text{skurt}[w^T x] = \frac{E[(w^T x - E[w^T x])^4]}{(\underbrace{\text{Var}[w^T x]}_{=1})^2} - 3$$

$$= E[(w^T x)^4] - 3$$

$$\mathcal{L}(w, \lambda) \Rightarrow E[(w^T x)^4] - 3 + \lambda(1 - \|w\|^2) = 0$$

$$\frac{\partial \mathcal{L}(\cdot)}{\partial w} \Rightarrow E\left[\frac{\partial}{\partial w} (w^T x)^4\right] + \lambda \left(-\frac{\partial}{\partial w} w^T w\right) = 0$$

$$\Leftrightarrow E[4(w^T x)^3 x^T] - 4\lambda w^T = 0$$

$$\Leftrightarrow E[(w^T x)^3 x] - \lambda w$$

$$\Leftrightarrow \lambda w = E[x (w^T x)^3]$$

$$c) \quad F(w) = E[x (w^T x)^3] - \lambda w = 0$$

$$f(w) = \frac{\partial F(w)}{\partial w} = E\left[\frac{\partial}{\partial w} x (w^T x)^3\right] - \lambda = 0$$

$$= E[3x x^T (w^T x)^2] - \lambda \mathbb{1} = 0$$

$$w^+ = w - f(w)^{-1} F(w) = w - (E[3x x^T (w^T x)^2] - \lambda \mathbb{1})^{-1} (E[x (w^T x)^3] - \lambda w)$$



$$d) w^+ = w - (E[3xx^T(w^T x)^2] - 1\mathbb{1})^{-1} (E[x(w^T x)^3] - 1w)$$

$$f(w) = E[3xx^T(w^T x)^2] = 3E[xx^T]E[(w^T x)^2] - 1\mathbb{1}$$

$$= 3\mathbb{1} - 1\mathbb{1} = \mathbb{1}(3-1)$$

$$w^+ = w - \mathbb{1}^{-1}(3-1)(E[x(w^T x)^3] - 1w)$$

$$\Leftrightarrow w^+ = w - \frac{E[x(w^T x)^3] - 1w}{3-1}$$

$$\Leftrightarrow (3-1)w^+ = 3w - 1w - E[x(w^T x)^3] + 1w$$

$$\Leftrightarrow 2(3-1)w^+ = -E[x(w^T x)^3] + 3w$$

$$\Leftrightarrow w^+ = -\underbrace{\frac{1}{3-1}}_f (E[x(w^T x)^3] - 3w)$$