

## Exercise Sheet 11

### 1 Activation Maximization

$$a) \quad \frac{\partial}{\partial x} \max_x w^T x + b - \lambda \|x\|^2 \stackrel{!}{=} 0$$

$$\Leftrightarrow w^T - 2\lambda x \stackrel{!}{=} 0 \quad \Leftrightarrow x^* = \frac{w^T}{2\lambda}$$

$$b) \quad \max_x w^T x + b + \log(p(x))$$

$$= \max_x w^T x + b + \log \left( \frac{\exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))}{(2\pi)^{\frac{d}{2}} |\Sigma|} \right)$$

$$= \max_x w^T x + b + (\log(\exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))) - \log((2\pi)^{\frac{d}{2}} |\Sigma|))$$

$$\Rightarrow \frac{\partial}{\partial x} w^T x - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \stackrel{!}{=} 0$$

$$\Leftrightarrow w^T - \Sigma^{-1}(x-\mu) = 0$$

$$\Leftrightarrow w^T - \Sigma^{-1}x + \Sigma^{-1}\mu = 0$$

$$\Leftrightarrow x = \Sigma(\Sigma^{-1}\mu + w^T)$$

$$\Leftrightarrow x^* = \mu + \Sigma w^T$$

$$c) \quad \frac{\partial}{\partial z} \max_z w^T(Az+c) + b - \lambda \|z\|^2$$

$$\Rightarrow \frac{\partial}{\partial z} w^T(Az+c) + b - \lambda \|z\|^2 \stackrel{!}{=} 0$$

$$\Leftrightarrow w^T A - 2\lambda z = 0$$

$$\Leftrightarrow z = \frac{w^T A}{2\lambda} \quad | \quad x = Az + c \Leftrightarrow z = \frac{x-c}{A}$$

$$\Leftrightarrow \frac{x-c}{A} = \frac{w^T A}{2\lambda}$$

$$\Leftrightarrow x = \frac{w^T A A}{2\lambda} + c$$



## 2 Layer - Info Relevance Propagation

$$a) \quad y = a_3 - a_4$$

$\underbrace{\quad}_{a_1} \quad \underbrace{\quad}_{\max(0, a_1 - a_2)}$

ReLU harm. effect  
i.e. as in positive

$$\text{if } a_1 \geq a_2 : y = a_1 - (a_1 - a_2) = a_2$$

$$\text{if } a_1 \leq a_2 : y = a_1 - 0 = a_1$$

$$h) \quad i) \quad R_4 = a$$

$$R_3 = a$$

$$R_1 = a$$

$$ii) \quad R_4 = a$$

$$R_3 = 0$$

$$R_1 = 0$$

## 3 Minimization

$$a) \quad O(x) = -\frac{1}{p} \log \left( \sum_{i=1}^M \alpha_i e^{-\mu \|x - \mu_i\|^2} \right)$$

$$= -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu \|x - \mu_i\|^2} \right) - \frac{1}{p} \log(\alpha_i)$$

$$h) \quad \min_{i=1}^M \{h_i\} = m$$

$$O(x) - m = -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu h_i} \right) + \frac{1}{p} \log e^{-\mu m}$$

$$= -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu h_i} e^{\mu m} \right)$$

$$= -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu (h_i - m)} \right)$$

$$= -\frac{1}{p} \log e^{-\mu \cdot 0} + \sum_{i=1}^M e^{-\mu (h_i - m)}$$

$$= -\frac{1}{p} \log (1 + \sum_{i=1}^M e^{-\mu (h_i - m)})$$

$$O(x) - m = 0 \quad \text{for } \mu \rightarrow \infty$$

$$\Leftrightarrow O(x) = m \rightarrow \text{min - pooling}$$