

We want  $\hat{P}_e$  to be close to  $P=0.05$ .

Alternatives:

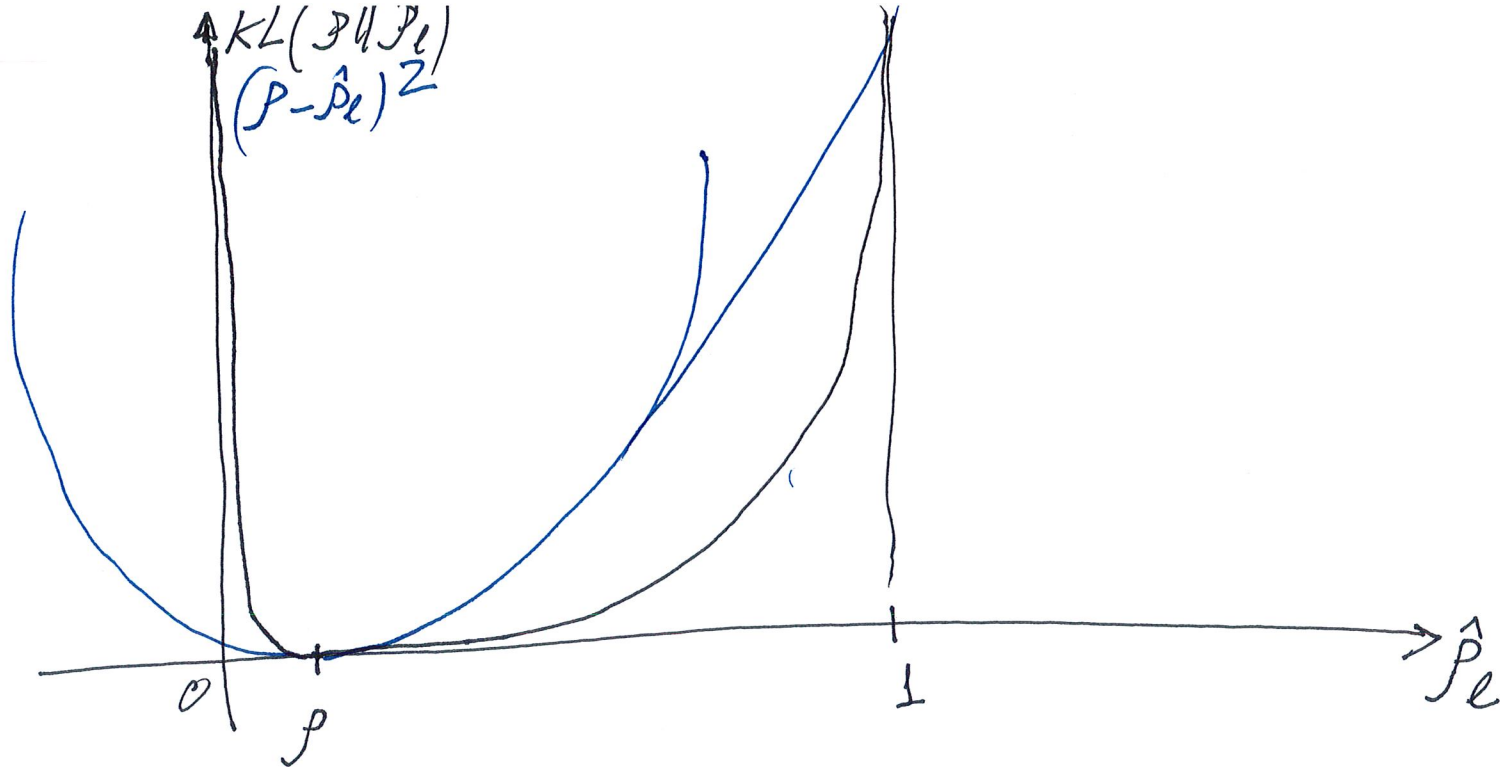
1) add  $(\hat{P}_e - P)^2$  to the loss function  $J$ .

$$2) \boxed{KL(P \parallel \hat{P}_e) = -P \log \frac{\hat{P}_e}{P} - (1-P) \log \frac{1-\hat{P}_e}{1-P}}$$

$$= +P \log \frac{P}{\hat{P}_e} + (1-P) \log \frac{1-P}{1-\hat{P}_e}$$

$$\boxed{KL(P \parallel Q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}}$$

→ KL divergence between  $P \equiv$  Bernoulli distrib w/param  $P$   
 $Q \equiv$  — u — w/param  $\hat{P}_e$



Input  $a^{(1)} = x$  to  $a^{(2)}$  to  $a^{(3)} = h(x)$

$$\hat{p}_e = \frac{1}{m} \sum_{k=1}^m a^{(2)}(x^{(k)})$$

$W^{(2)}$

$$J = \frac{1}{2} \|a^{(3)} - x\|^2 + \frac{1}{2} \|W\|^2 + \beta \sum_{l=1}^{s_2} KL(\hat{p}_l)$$

Need to compute gradient of new (sparsity) term  $\beta \sum_{l=1}^L \frac{KL(\mathcal{P} || \hat{\mathcal{P}}_l)}{KL(\hat{\mathcal{P}}_l)}$   
 w.r.t.  $W_{ij}^{(1)}$

$$KL(\mathcal{P}, \hat{\mathcal{P}}_l) = \mathcal{P} \log \frac{\mathcal{P}}{\hat{\mathcal{P}}_l} + (1-\mathcal{P}) \log \frac{(1-\mathcal{P})}{(1-\hat{\mathcal{P}}_l)}$$

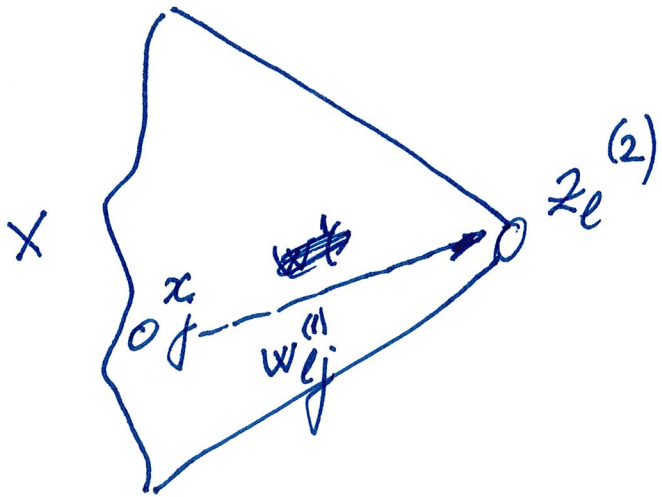
$$\hat{\mathcal{P}}_l = \frac{1}{m} \sum_{k=1}^m a_l^{(2)}(x^{(k)}) = \frac{1}{m} \sum_{k=1}^m f\left(\frac{|x^{(k)}|}{a-1} \sum_{a=1}^a W_{la}^{(1)} \cdot x_a^{(k)} + b_l^{(1)}\right)$$

$$\frac{\partial KL(\hat{\mathcal{P}}_l)}{\partial W_{ij}^{(1)}} = \frac{\partial KL(\hat{\mathcal{P}}_l)}{\partial \hat{\mathcal{P}}_l} \cdot \frac{\partial \hat{\mathcal{P}}_l}{\partial W_{ij}^{(1)}}$$

$$KL(\mathcal{P}, \hat{\mathcal{P}}_l) = -\mathcal{P} \log \frac{\hat{\mathcal{P}}_l}{\mathcal{P}} - (1-\mathcal{P}) \log \frac{1-\hat{\mathcal{P}}_l}{1-\mathcal{P}}$$

$$\Rightarrow \frac{\partial KL}{\partial \hat{\mathcal{P}}_l} = -\mathcal{P} \cdot \frac{\mathcal{P}}{\hat{\mathcal{P}}_l} \cdot \frac{1}{\mathcal{P}} - (1-\mathcal{P}) \cdot \frac{1-\mathcal{P}}{1-\hat{\mathcal{P}}_l} \cdot -\frac{1}{1-\mathcal{P}}$$

$$= -\frac{\mathcal{P}}{\hat{\mathcal{P}}_l} + \frac{1-\mathcal{P}}{1-\hat{\mathcal{P}}_l}$$



$$\frac{\partial z_l}{\partial w_{ij}^{(1)}} = \begin{cases} x_j \text{ (also } a_j^{(1)}) & \text{if } l=i \\ 0 & \text{if } l \neq i \end{cases}$$