

(\*) Both imply a factorisation of  $p(x|\theta)$

(\*) Neyman factorisation  $\rightarrow$  frequentist definition of sufficiency

(\*) Sufficiency in this context means  $T(x)$  is sufficient for  $\theta$  if:-

$$\theta \perp\!\!\!\perp x \mid T(x)$$

Jordan (2003)

- 8.1.8. ML and KL divergence

- A general rel. between ML and KL divergence (not spec. to exp.)

- Necessary for later lec. material 16, 17

- Statistical interp of KL divergence to illus. rel. between KL and exp. family

(\*) empirical distn:  $\hat{p}(x)$

- Places a point mass at each data point  $x_n$  in  $\mathcal{D}$  (dataset) (discrete)

(\*) empirical distn: 
$$\hat{p}(x) := \frac{1}{N} \sum_{n=1}^N \delta(x, x_n) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}(x = x_n)$$

(\*) Sum/integrate

$\hat{p}(x)$  against a function of  $x$ ; we evaluate

at each point  $x_n$

(\*) log likelihood: (also, cross entropy of  $\hat{p}(x)$  and  $p(x|\theta)$ )

$$\sum_x \hat{p}(x) \log p(x|\theta) = \sum_x \left( \frac{1}{N} \sum_{n=1}^N \delta(x, x_n) \log p(x_n|\theta) \right)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_x \delta(x, x_n) \log p(x_n|\theta)$$

$$= \frac{1}{N} \sum_{n=1}^N \log p(x_n|\theta)$$

$$= \frac{1}{N} \ell(\theta|\mathcal{D})$$

$\delta$ -Kronecker  
delta in cont.  
case.

↙ de-trickifying

(\*) Note; the scaled log likelihood (by factor  $1/N$ ) is equivalent to the cross-entropy between the empirical distri and the model  
 $(\hat{p}(x))$   $\neq p(x|\theta)$

- same result for continuous

- KL divergence between empirical; model :-

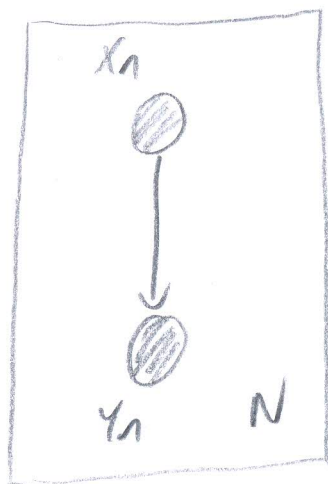
$$\begin{aligned} D(\hat{p}(x) \| p(x|\theta)) &= \sum_x \hat{p}(x) \log \frac{\hat{p}(x)}{p(x|\theta)} \\ &= \sum_x \hat{p}(x) \log \hat{p}(x) - \sum_x \hat{p}(x) \log p(x|\theta) \\ &= \underbrace{\sum_x \hat{p}(x) \log \hat{p}(x)}_{(1)} - \underbrace{\frac{1}{N} \ell(\theta|D)}_{(2)} \end{aligned}$$

(1) - independent of  $\theta$

- value of  $\theta$  that minimises LHS is the value of  $\theta$  that maximises the RHS

(\*) Minimising KL divergence between the empirical distri and model distri is equivalent to maximising the likelihood

(\*) Generalised linear Models (GLM) - linear regression/classification  
 covers linear reg./discriminative linear classification.



(\*) Both LR/UC  $\rightarrow$  both assume a rep. for conditional expectation of  $y$ .

$$(*) \mu = f(\theta^T x) = \mathbb{E}_{y \sim p(y|x)} [y]$$

(\*) LR:  $f(\cdot)$  - identity

UC:  $f(\cdot)$  - sigmoid (logistic)

(\*) Also: endow  $Y$  with a particular cond. prob. distribution, with  $\mu$  as a parameter.

(\*) Remember JP  $\rightarrow$  ColumbiaX ML (prob. intep of ML for LR!)

(\*) LR - Gaussian LC - Bernoulli / Multinomial.

(\*) Generalised Linear Model Framework

- 3 assumptions on  $p(y|x)$  :-

1. observed input  $x$  enters into model via linear comb.  $\xi = \theta^T x$
2. conditional mean  $\mu$  rep as a function  $f(\xi)$  of the linear combination  $\xi$  where  $f$  is known as the response function
3. observed output  $y$  is assumed to be characterised by exp. family with conditional mean  $\mu$ .

