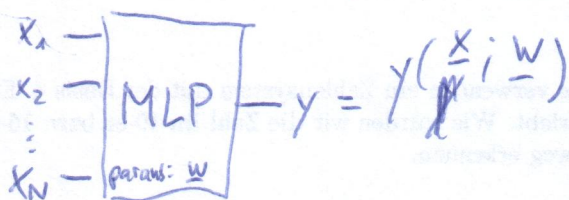


# Multilayer Perceptron (MLP) : Cost Functions

$(*) \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 = \frac{1}{2}$  vs.  $1 = \frac{1}{2} + \frac{1}{2}$   
 $\left(\frac{1}{4}\right)^2 + \left(\frac{3}{4}\right)^2 = \frac{10}{16} = \frac{5}{8}$  vs.  $1 = \frac{1}{4} + \frac{3}{4}$   
 for distance  $d$  between labels

a) What will be the effect of the error measure choice?



$$\underline{w} \in \underset{\underline{w}}{\operatorname{argmin}} E^T(\underline{w}) = \underset{\underline{w}}{\operatorname{argmin}} \frac{1}{P} \sum_{\alpha=1}^P e^{(\alpha)}(\underline{w})$$

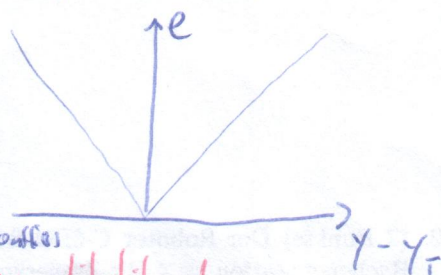
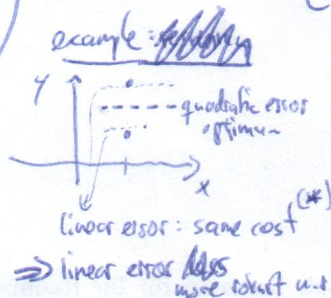
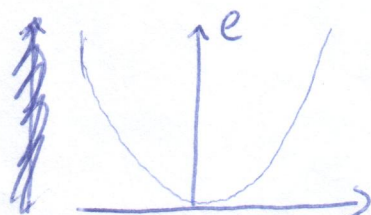
## COST FUNCTIONS FOR REGRESSION

1) quadratic error function:

$$e^{(\alpha)} = \frac{1}{2} \left( y(x^{(\alpha)}; \underline{w}) - y_T^{(\alpha)} \right)^2$$

2) linear error

$$e^{(\alpha)} = \left| y(x^{(\alpha)}; \underline{w}) - y_T^{(\alpha)} \right|$$



b) Relation between quadratic error and Gaussian labels  $\rightarrow$  Maximum likelihood.

Remark:

if labels are conditionally Gaussian,  
 i.e.,  $P(y|x) = \mathcal{N}(y | y(x; \underline{w}), \sigma^2)$

then:  $\underline{w}_{\text{opt}} \in \underset{\underline{w}}{\operatorname{argmin}} E^T(\underline{w})$

$$\Leftrightarrow \underline{w}_{\text{opt}} \in \underset{\underline{w}}{\operatorname{argmax}} \prod_{\alpha=1}^P \mathcal{N}(y^{(\alpha)} | y(x^{(\alpha)}; \underline{w}), \sigma^2)$$

$\Rightarrow$  min. squared error corresponds to max. likelihood solution for Gaussian labels. (hint: log-likelihood...)

This property makes ~~square~~ the quadratic cost function the standard choice for regression!

$$\underline{w}_{\text{opt}} \in \underset{\underline{w}}{\operatorname{argmin}} \left( -\ln [\mathcal{L}(\underline{w})] \right)$$



## MLP cost functions (cont'd)

c) Derive the cross entropy cost function for binary classification.

Data:  $\{(\underline{x}^{(a)}, y_T^{(a)})\}_{a=1}^P$  with  $\underline{x}^{(a)} \in \mathbb{R}^N$ ,  $y_T^{(a)} \in \{0, 1\}$   
generated by  $(\underline{x}^{(a)}, y_T^{(a)}) \stackrel{iid}{\sim} p_{\text{data}}(\underline{x}, y)$

Model: MLP, output interpreted as probability that  $y=1$ :

$$y(\underline{x}; \underline{w}) =: p_{\text{model}}(y=1 | \underline{x}) \Rightarrow p_{\text{model}}(y=0 | \underline{x}) = 1 - y(\underline{x}; \underline{w})$$

Cost function derivation via minimization of Kullback-Leibler divergence between data and model distribution. ( $D_{KL} = 0$  for identical and  $D_{KL} > 0$  for different distributions, not symmetric!)

$$D_{KL}(p_{\text{data}}(\underline{x}, y) || p_{\text{model}}(\underline{x}, y))$$

→ see lecture.

$$\begin{aligned} &= \int_{\mathbb{R}^N} d\underline{x} \sum_{y \in \{0,1\}} p_{\text{data}}(\underline{x}, y) \ln \left( \frac{p_{\text{data}}(\underline{x}, y)}{p_{\text{model}}(\underline{x}, y)} \right) = \int_{\mathbb{R}^N} d\underline{x} \sum_{y \in \{0,1\}} p_{\text{data}}(\underline{x}) p_{\text{data}}(y | \underline{x}) \ln \frac{p_{\text{data}}(\underline{x}) p_{\text{data}}(y | \underline{x})}{p_{\text{model}}(\underline{x}) p_{\text{model}}(y | \underline{x})} \\ &= \underbrace{\int_{\mathbb{R}^N} d\underline{x} p_{\text{data}}(\underline{x}) \sum_{y \in \{0,1\}} p_{\text{data}}(y | \underline{x}) \ln [p_{\text{data}}(y | \underline{x})]}_{\text{independent of model parameters } \underline{w}} - \underbrace{\int_{\mathbb{R}^N} d\underline{x} p_{\text{data}}(\underline{x}) \sum_{y \in \{0,1\}} p_{\text{data}}(y | \underline{x}) \ln [p_{\text{model}}(y | \underline{x})]}_{\text{cross entropy between (conditional) data and model distribution (for given } \underline{x})} \end{aligned}$$

discriminative not generative model  
data

Empirical risk minimization:

$$\begin{aligned} \bar{E}^T &= -\frac{1}{P} \sum_{a=1}^P \sum_{y \in \{0,1\}} p_{\text{data}}(y | \underline{x}^{(a)}) \ln [p_{\text{model}}(y | \underline{x}^{(a)})] \\ &= \frac{1}{P} \sum_{a=1}^P - \left( p_{\text{data}}(y=1 | \underline{x}^{(a)}) \ln [p_{\text{model}}(y=1 | \underline{x}^{(a)})] + p_{\text{data}}(y=0 | \underline{x}^{(a)}) \ln [p_{\text{model}}(y=0 | \underline{x}^{(a)})] \right) \\ &= \frac{1}{P} \sum_{a=1}^P \left( -y_T^{(a)} \ln [y(\underline{x}^{(a)}; \underline{w})] - (1 - y_T^{(a)}) \ln [1 - y(\underline{x}^{(a)}; \underline{w})] \right) \end{aligned}$$

→ cross entropy cost function,  $\bar{E}^{(a)}$

multi-class version: softmax output / multiclass cross entropy, cf. lecture



# NLP cost fcts (cont'd)

d)

## Cost functions & output layers

→ show after backprop derivatives

task	data	output layer	cost function
scalar regression	$\{(\underline{x}^{(a)}, y_T^{(a)})\}$ $\prod_{a=1}^n \mathbb{R}^N \quad \prod_{a=1}^n \mathbb{R}$	linear $y = \sum_{k=0}^N \underbrace{w_{1k}^{L,L-1} s_k^{L-1}}_{=h_1^L}$	mean squared error $e^{(a)} = \frac{1}{2} [y_T^{(a)} - y(\underline{x}^{(a)}; \underline{w})]^2$
vectorial regression	$\{(\underline{x}^{(a)}, \underline{y}_T^{(a)})\}$ $\prod_{a=1}^n \mathbb{R}^N \quad \prod_{a=1}^n \mathbb{R}^M$	linear $\underline{y}_i = \sum_{k=0}^N \underbrace{w_{ik}^{L,L-1} s_k^{L-1}}_{=h_i^L}$	mean squared Euclidean distance $e^{(a)} = \frac{1}{2} \ \underline{y}(\underline{x}^{(a)}; \underline{w}) - \underline{y}_T^{(a)}\ _2^2$
binary classification	$\{(\underline{x}^{(a)}, y_T^{(a)})\}$ $\prod_{a=1}^n \{0,1\}$	logistic sigmoidal $y = \frac{1}{1 + \exp(-h_1^L)}$ <del>linear</del>	cross entropy $e^{(a)} = -y_T^{(a)} \ln y(\underline{x}^{(a)}; \underline{w}) - (1 - y_T^{(a)}) \ln [1 - y(\underline{x}^{(a)}; \underline{w})]$
classification with M classes	$\{(\underline{x}^{(a)}, \underline{y}_T^{(a)})\}$ $\prod_{a=1}^n \{0,1\}^M$ <u>1-hot-code</u> e.g. $\underline{y}_T = (0,0,1,0,0,0)^T \Rightarrow$ class 3 of 6	softmax $y_i = \frac{\exp(h_i^L)}{\sum_{k=1}^M \exp(h_k^L)} \Rightarrow \sum_{k=1}^M y_k = 1$	cross entropy (multi class case) $e^{(a)} = - \sum_{i=1}^M y_{T,i}^{(a)} \ln y_i(\underline{x}^{(a)}; \underline{w})$

M-class case: next lecture(s)

Remarks  
 → all classes explicitly included:  $y_i = P(y=i | \underline{x})$   
 → binary classification implicit:  $y = P(y=1 | \underline{x})$   
 $1-y = P(y=0 | \underline{x})$