

* Probabilist interpretation of linear regression

MLE: For IID data from some distribution $P(x|\theta_0)$ the MLE minimizes the KL divergence (KULLBACK-LEIBLER divergence)

$$\hat{\theta} = \arg \max_{\theta} \prod_{i=1}^m P(x^{(i)}|\theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log P(x^{(i)}|\theta)$$

$$= \arg \max_{\theta} \frac{1}{m} \sum_{i=1}^m \log P(x^{(i)}|\theta)$$

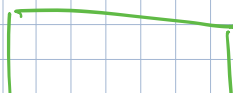
$$- \frac{1}{m} \sum_{i=1}^m \log P(x^{(i)}|\theta_0)$$

$$= \arg \max_{\theta} \frac{1}{m} \sum_{i=1}^m \log P(x^{(i)}|\theta) / P(x^{(i)}|\theta_0)$$

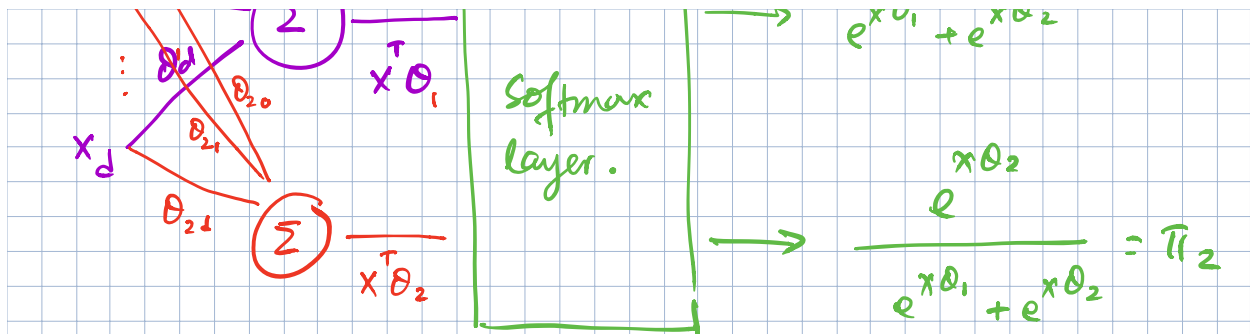
$$= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \log P(x^{(i)}|\theta_0) / P(x^{(i)}|\theta)$$

$$\underset{m \rightarrow \infty}{\theta} \arg \min_{\theta} \int \log \frac{P(x|\theta_0)}{P(x|\theta)} p(x|\theta_0) dx$$

Softmax Formulation



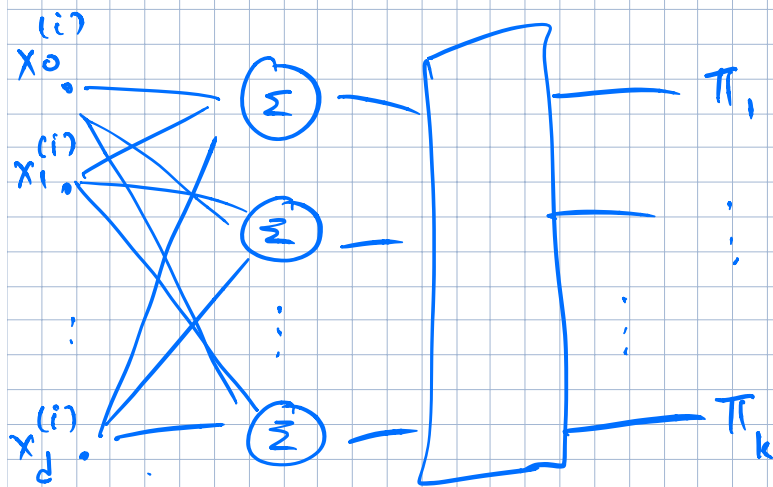
$$\frac{e^{x\theta_1}}{e^{x\theta_0} + e^{x\theta_1}} = \pi_1$$



$$\theta_1 = \begin{pmatrix} \theta_{10} \\ \theta_{11} \\ \vdots \\ \theta_{1d} \end{pmatrix}$$

$$\theta_2 = \begin{pmatrix} \theta_{20} \\ \theta_{21} \\ \vdots \\ \theta_{2d} \end{pmatrix}$$

$$\pi_1 + \pi_2 = 1$$



Likelihood

$$\begin{aligned}
 P(y|x;\theta) &= \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta) \\
 &= \prod_{i=1}^n \pi_{i0}(y^{(i)}) \prod_{i2} \pi_{i1}(y^{(i)})
 \end{aligned}$$

$$\pi_{i1} = P(y=0 | x; \theta)$$

$$\pi_{i2} = P(y=1 | x; \theta)$$

$$\mathbb{1}_0(y^{(i)}) = 1 \quad \text{if } y^{(i)} = 0$$

$$\mathbb{1}_1(y^{(i)}) = 1 \quad \text{if } y^{(i)} = 1$$

$$\pi_{i1} = \frac{e^{x^T \theta_1}}{e^{x^T \theta_1} + e^{x^T \theta_2}}$$

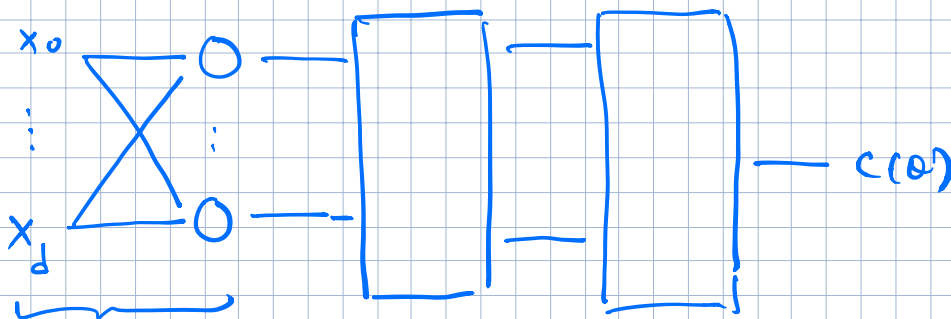
$$\pi_{i2} = \frac{e^{x^T \theta_2}}{e^{x^T \theta_1} + e^{x^T \theta_2}}$$

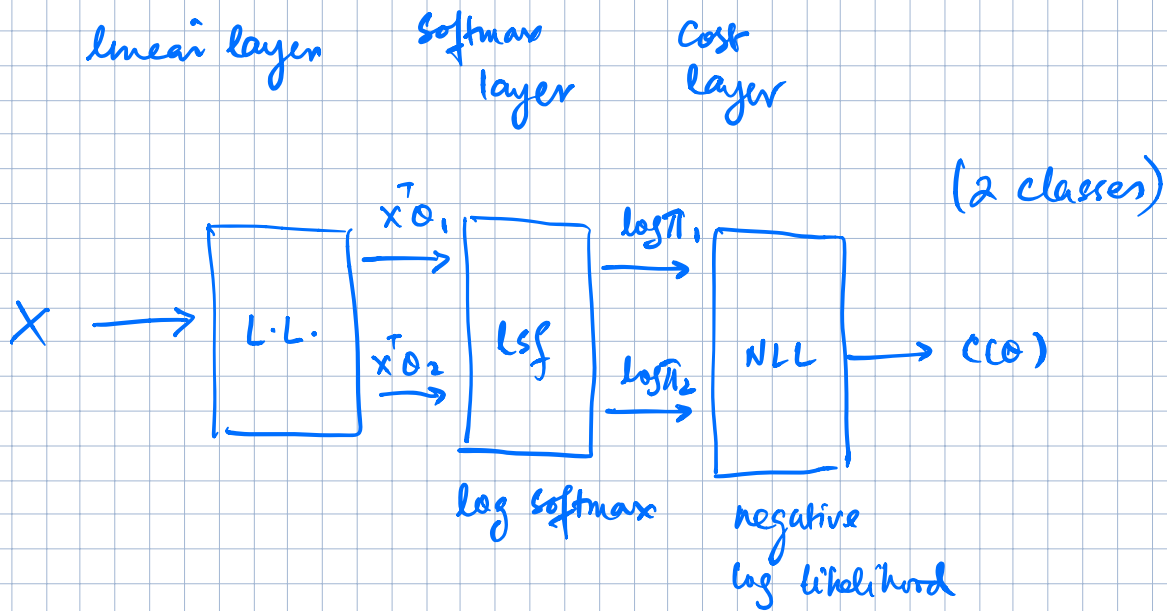
I am interested in MLE.

$$C(\theta) = -\log \text{likelihood}$$

$$= -\sum_{i=1}^n \left(\mathbb{1}_0(y^{(i)}) \log \pi_{i1} + \mathbb{1}_1(y^{(i)}) \log \pi_{i2} \right)$$

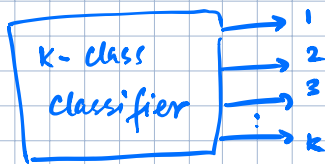
Layered view of Softmax





Cross-Entropy

Problem: How do we compare two vectors?



x_1	x_2	Labels	
3	7	3	0010000
18	47	1	1000000
2	4	4	0001000
42	1	7	0000001

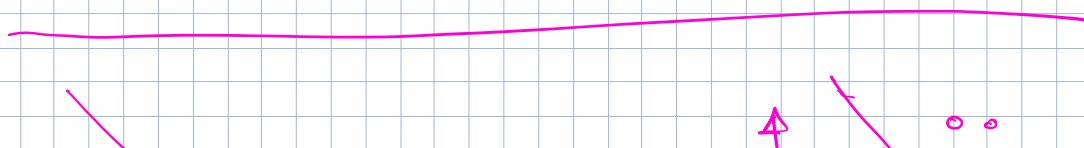
\hat{y}
(predicted)

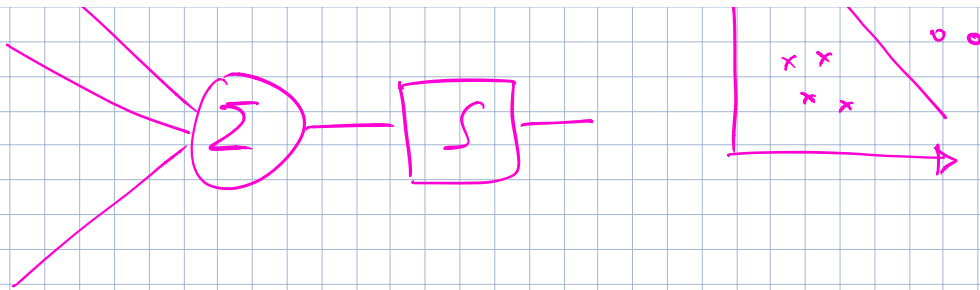
y
(targets)

Compare \hat{y} and y using cross-entropy.

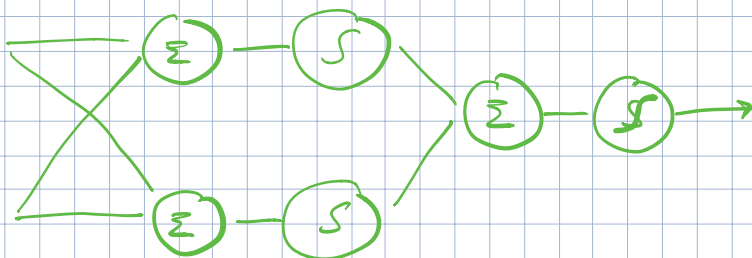
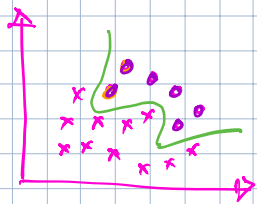
$$D(\hat{y}, y) = - \sum_{i=1}^7 y_i \log \hat{y}_i$$

$$D(\hat{y}, y) \neq D(y, \hat{y})$$

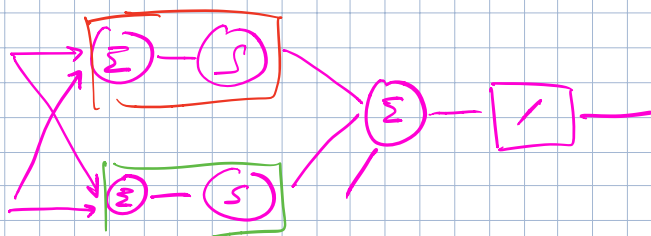
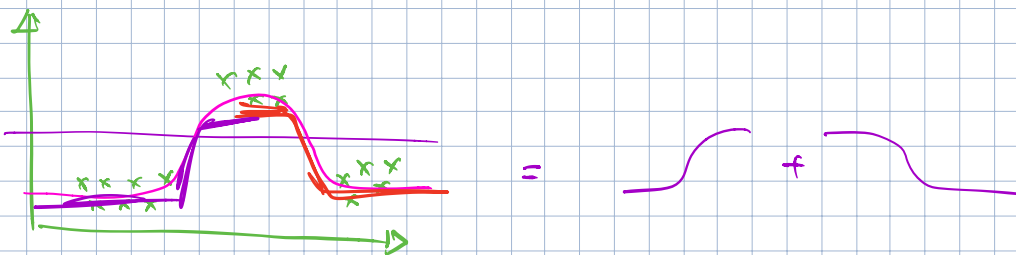




How do I classify the following:



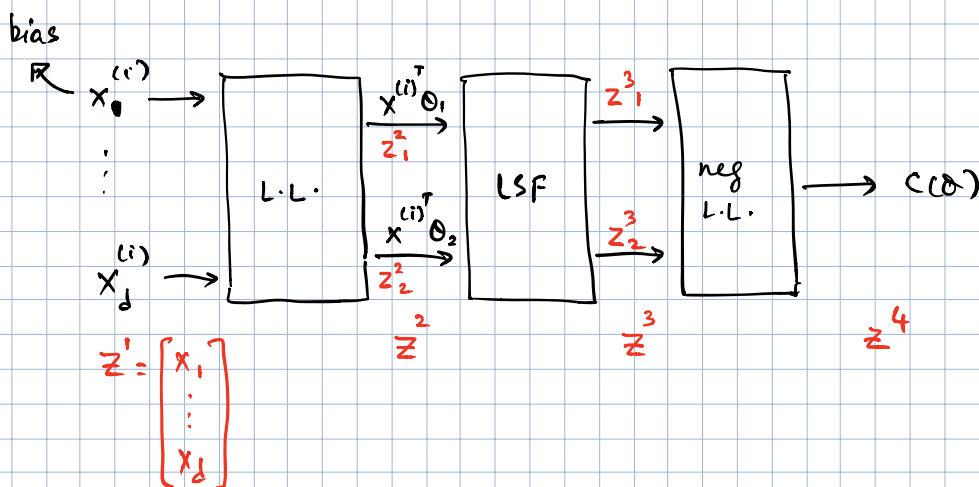
The more neurons the more complex (wobbly) the discriminant function.



Backpropagation

What information does a layer need to know?

- How to compute output from its input
- How to propagate derivatives backwards
- How to compute derivatives w.r.t. learning parameters.



$$z_1^2 = x^T \theta_1$$

$$z_2^2 = x^T \theta_2$$

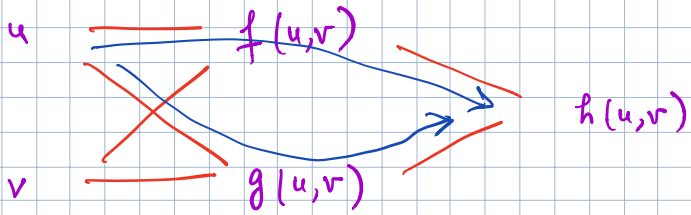
$$z_1^3 = \log \frac{e^{z_1^2}}{e^{z_1^2} + e^{z_2^2}}, \quad z_2^3 = \log \frac{e^{z_2^2}}{e^{z_1^2} + e^{z_2^2}}$$

$$z^4 = - \sum_{i=1}^n \mathbb{1}_0(y_i) z_1^3 + \mathbb{1}_1(y_i) z_2^3$$

$$c(\theta) = z^4 \left[z_1^3 \left\{ z_1^2(\theta_1, z^1), z_2^2(\theta_2, z^1) \right\}, z_2^3 \left\{ z_1^2(\theta_1, z^1), z_2^2(\theta_2, z^1) \right\} \right]$$

$$\frac{\partial c(\theta)}{\partial \theta} = ?$$

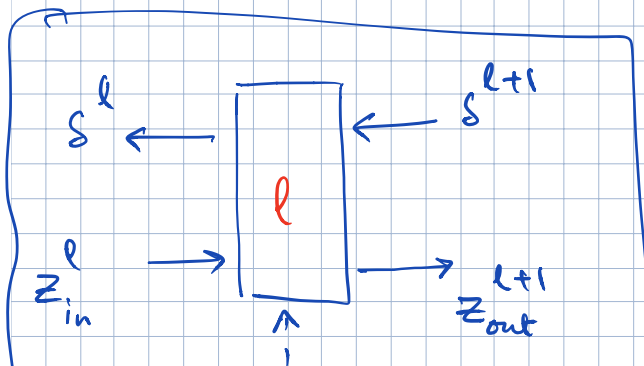
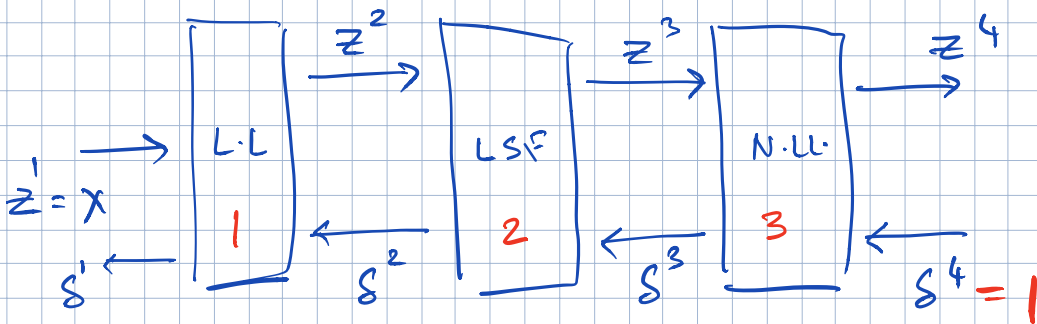
Chain Rule:



$$h(u, v) = \{ f(u, v), g(u, v) \}$$

$$\frac{\partial h}{\partial u} = \frac{\partial h}{\partial f} \frac{\partial f}{\partial u} + \frac{\partial h}{\partial g} \frac{\partial g}{\partial u}$$

$$\frac{\partial z^4}{\partial \theta_1} = \frac{\partial z^4}{\partial z_1^3} \frac{\partial z_1^3}{\partial \theta_1} + \frac{\partial z^4}{\partial z_2^3} \frac{\partial z_2^3}{\partial \theta_1}$$



$$\frac{\partial c}{\partial \theta^l}$$

$$\begin{aligned}\delta_i^l &= \frac{\partial c}{\partial z_i^l} = \sum_j \frac{\partial c}{\partial z_j^{l+1}} \frac{\partial z_j^{l+1}}{\partial z_i^l} \\ &= \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial z_i^l}\end{aligned}$$

$$\delta_i^l = \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

$$z^{l+1} = f^l(z^l)$$

$$\frac{1}{1 + e^{-z'}}$$

$$= \sigma(z) (1 - \sigma(z))$$

$$\frac{\partial c}{\partial \theta^l} = \sum_j \frac{\partial c}{\partial z_j^{l+1}} \frac{\partial z_j^{l+1}}{\partial \theta}$$

↑
j

$$= \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial \theta}$$