

Plan

Review

Logistic Regression

↳ Sigmoid

↳ Softmax

Cross Entropy Loss

Logistics

- Zoom!
- check in form (10/15)
- scribed notes ☺
- 2-3 work and struggle

Linear Regression

$$(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$$

$$x^{(i)} \in \mathbb{R}^d \quad y^{(i)} \in \mathbb{R}$$

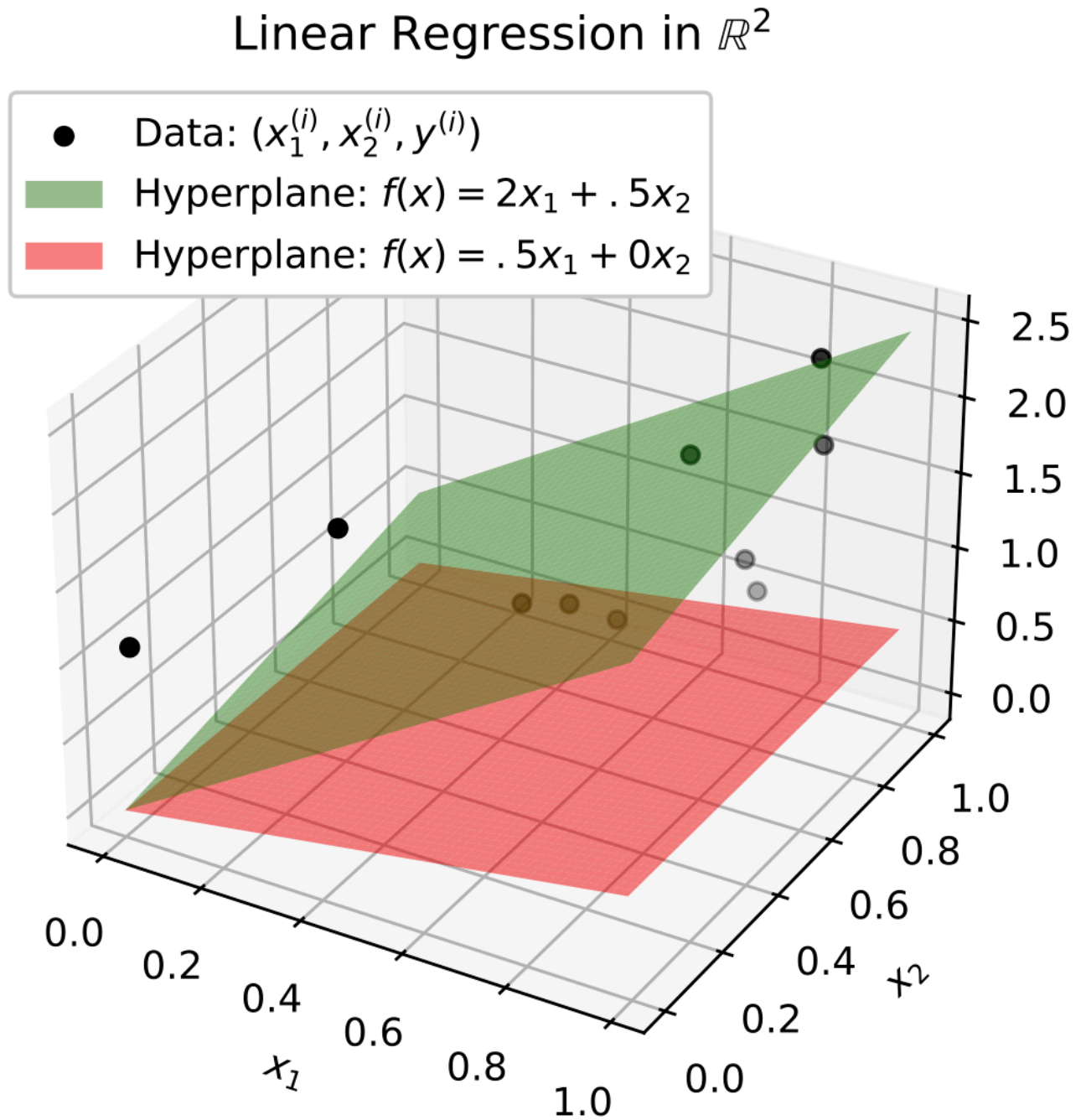
① Model: $f(x) = w \cdot x$
for $w \in \mathbb{R}^d$

② Loss: $\mathcal{L}(w) = \frac{1}{n} \|Xw - y\|_2^2$

③ Optimization: $\nabla_w \mathcal{L}(w^*) = 0$

$$\Leftrightarrow w^* = \underbrace{(X^T X)^{-1} X^T y}_{d \times 1}$$

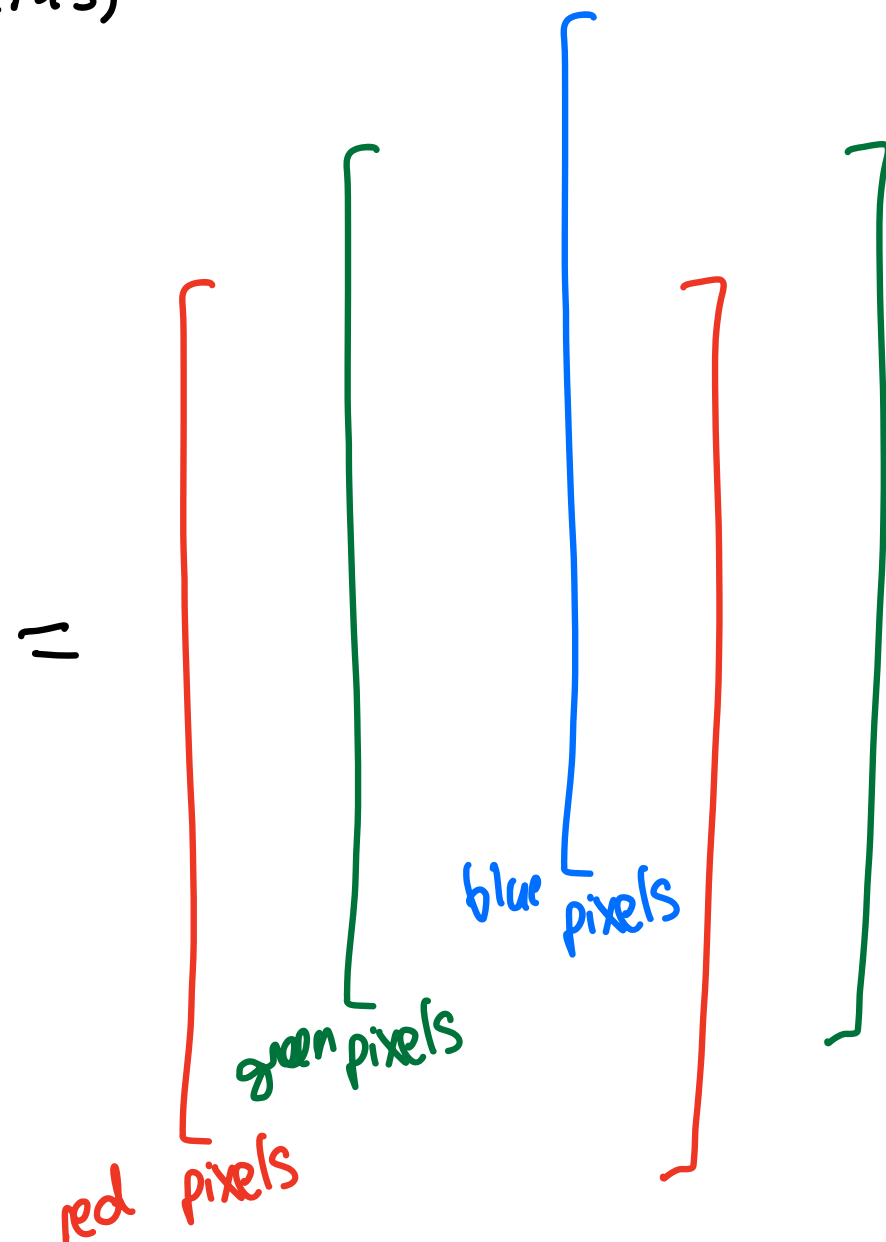
$\begin{matrix} d \times n & n \times d & d \times n & n \times 1 \end{matrix}$



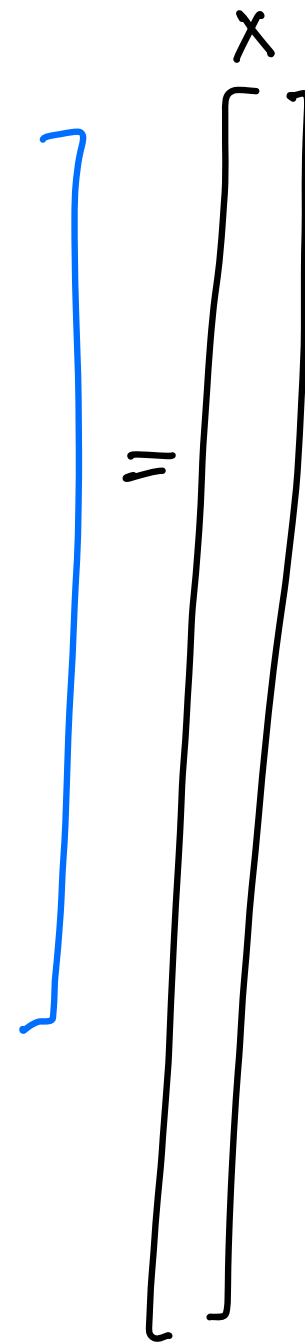
Motivation

- What if labels are classes (rather than values)?

(image, cat status)



- What happens when we can't find the exact optimal?



$y \in \{0, 1\}$

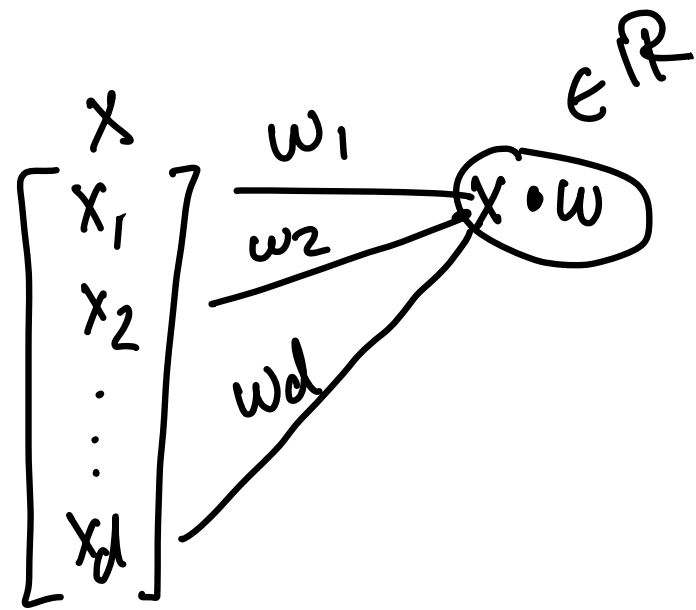
Supervised Binary Classification

... $(x^{(i)}, y^{(i)})$...

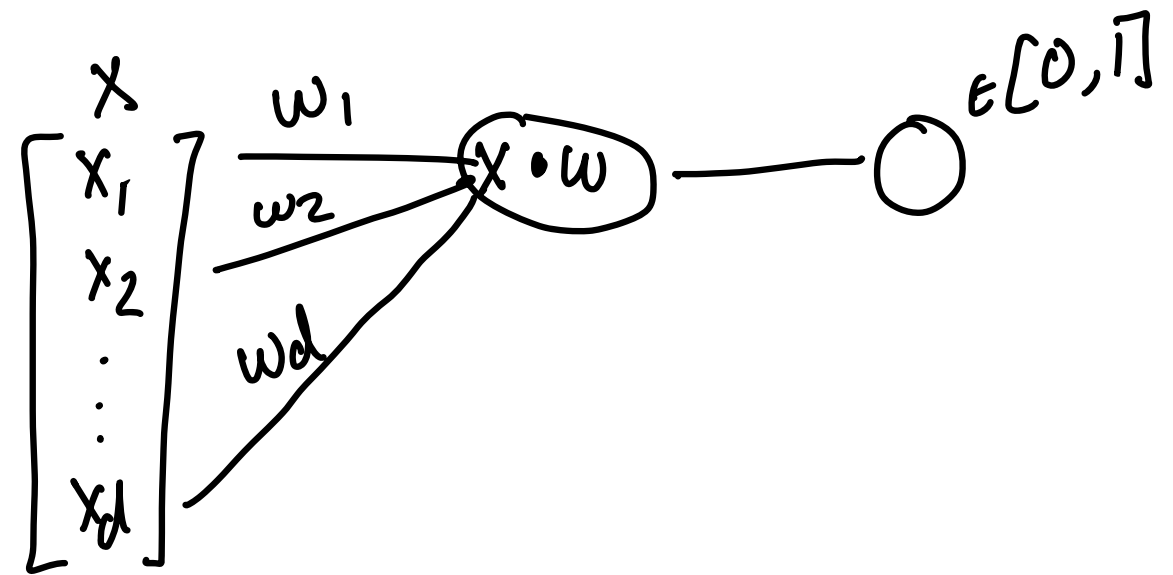
$$x^{(i)} \in \mathbb{R}^d$$

$$y^{(i)} \in \{0, 1\}$$

Goal: $f(x^{(i)}) =$ probability of positive class



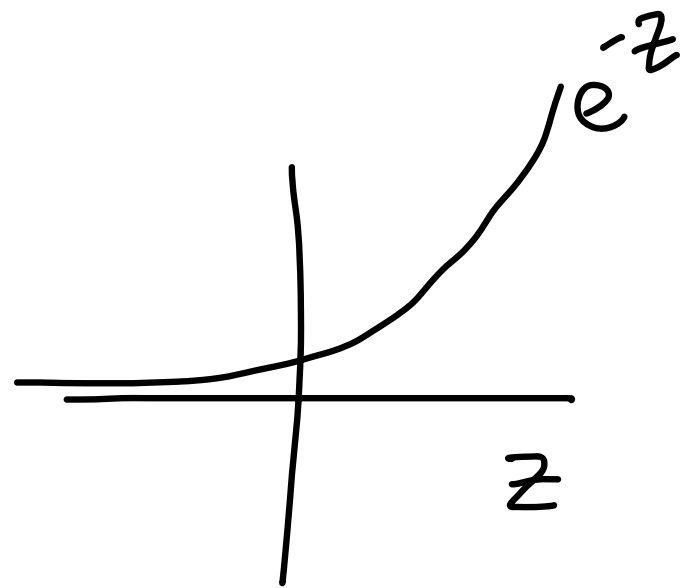
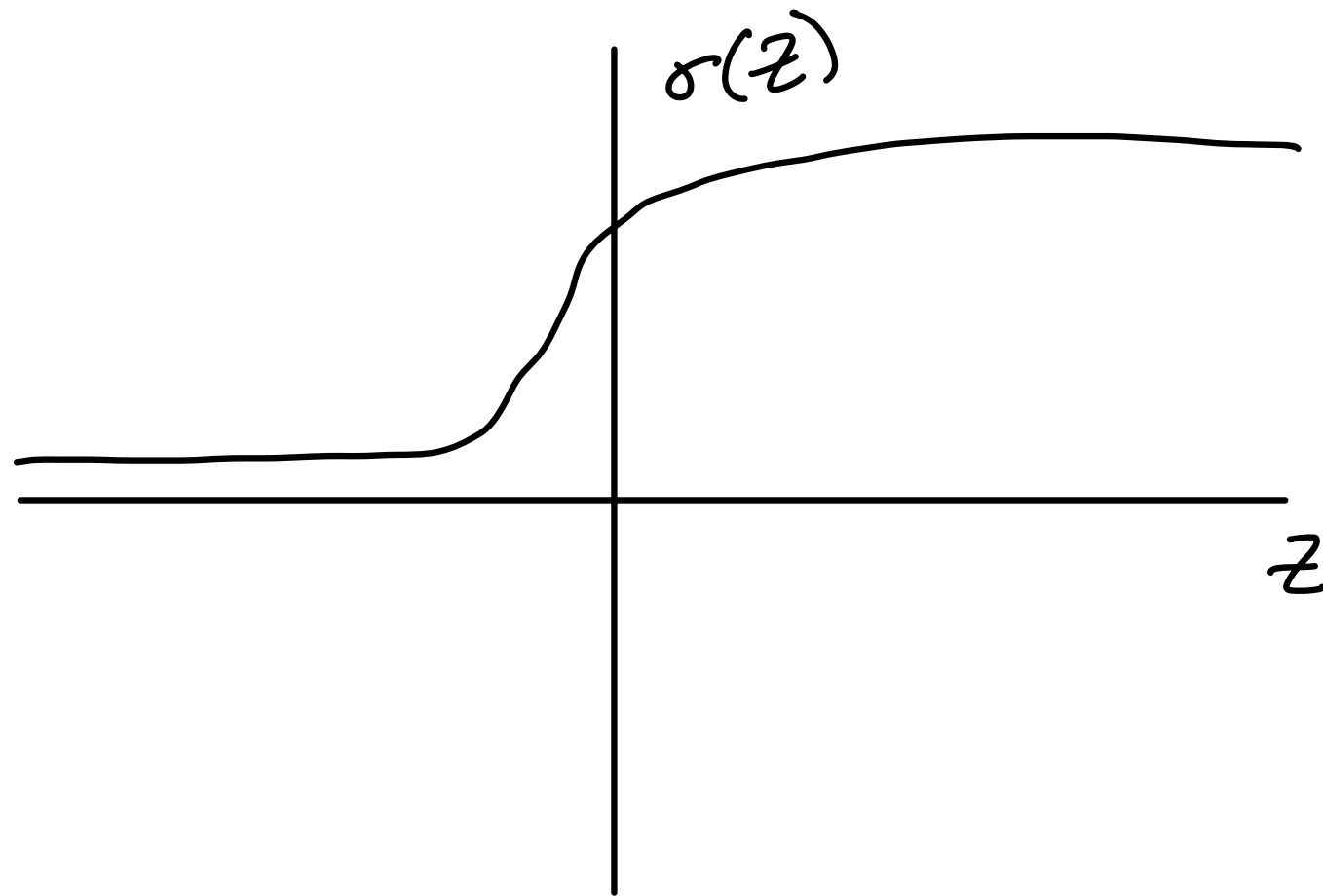
Linear Regression



Logistic Regression

Sigmoid

$$\sigma(z) = \frac{1}{1+e^{-z}}$$



$$\lim_{z \rightarrow \infty} \frac{1}{1+e^{-z}} = \frac{1}{1+0} = 1$$

$$\lim_{z \rightarrow -\infty} \frac{1}{1+e^{-z}} = \frac{1}{\infty} = 0$$

Multiple Classes

... $(x^{(i)}, y^{(i)})$...

$x^{(i)} \in \mathbb{R}^d$ $y^{(i)} \in \{0, 1, \dots, k\}$

$$f: \mathbb{R}^d \rightarrow [0, 1]^k$$

Probability distribution

1. non-negative
2. Sums to 1

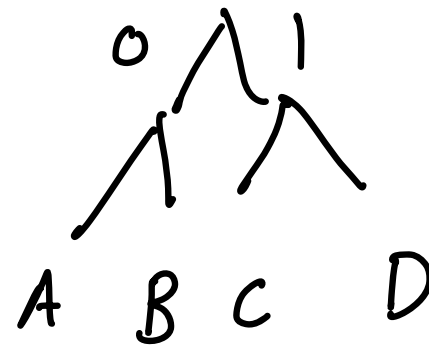
Goal: $f(x^{(i)})_l =$ prob of class l

What should the architecture be?

Cross Entropy (aside)

Communicate A, B, C, D

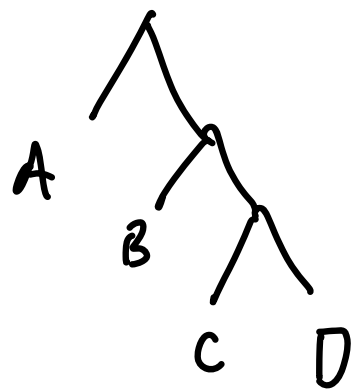
Approach #1: 00, 01, 10, 11



e.g., $A = 1/2, B = 1/4, C = 1/8, D = 1/8$
 1, 2, 3, 3

But what if A more likely?
 Then

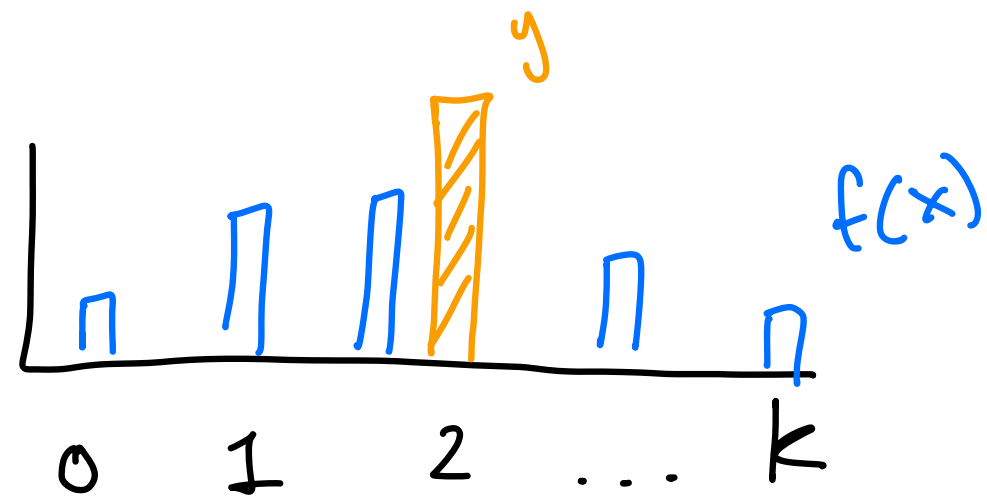
$$l_i = -\lceil \log_2 q_i \rceil$$



$$H(q) = \text{Entropy} = \mathbb{E}[\# \text{ bits}] = \mathbb{E}_{i \sim q}[l_i] = -\mathbb{E}_{i \sim q}[\log_2 q_i] = -\sum_{i=1}^k q_i \log_2(q_i)$$

$$H(p, q) = \text{Cross entropy} = \text{communicate w/ wrong distribution} = -\mathbb{E}_{i \sim p}[\log_2 q_i]$$

Loss



Goal: Measure distance between
distributions y and f(x)

Cross entropy between p and q:

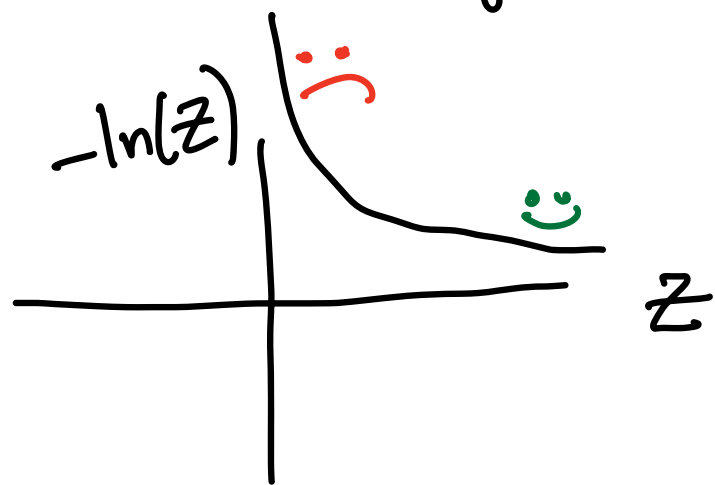
$$H(p, q) = - \mathbb{E}_{j \sim p} [-\ln(q(j))] = - \sum_{j=1}^k p(j) \ln(q(j))$$

When p is "one-hot" i.e., $\exists j^*$ s.t. $p(j^*) = 1$ then

$$H(p, q) = - p(j^*) \ln(q(j^*)) = -\ln(q(j^*))$$

Cross entropy between p and q :

$$H(p, q) = - \sum_{j=1}^k p(j) \ln(q(j)) \stackrel{\text{"one-hot"}}{=} -\ln(q(j^*))$$



Optimization

- Exact is not doable (you'll see on problem),
but yet we can still find good weights. How?