

# **Statistical NLP**

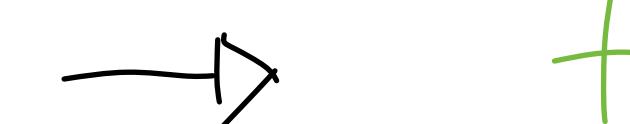
## **Manual Feature Engineering**

## **Linear Models and Classification**

Tim Rocktäschel, Sebastian Riedel

# Sentiment Analysis

This class isn't bad but I don't  
like <sup>the</sup> examples the lecturer  
is using! Bad grammar!



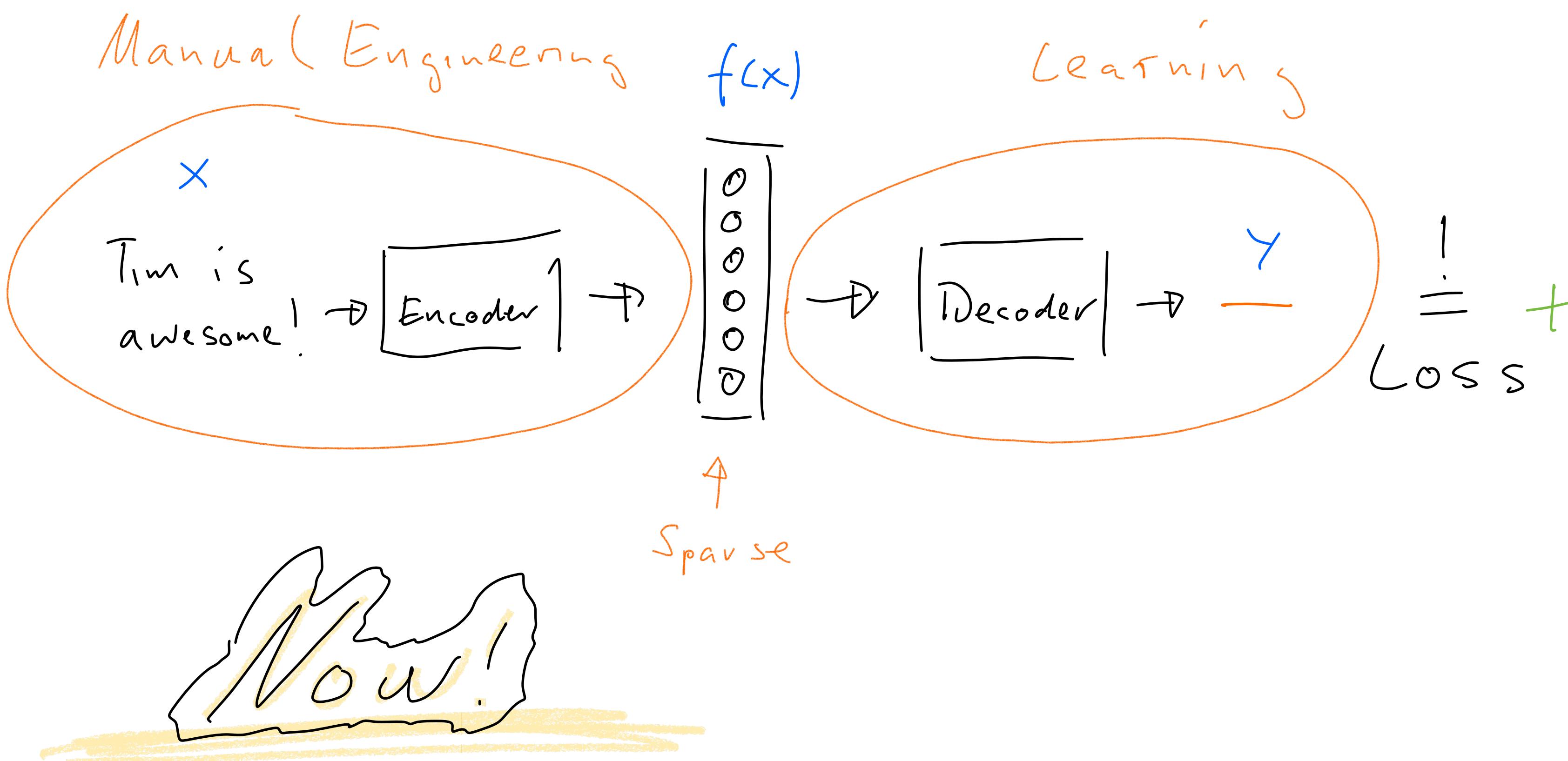
I should have gone to  
Imperial!



- Important task in research and academia (i.e. I see many papers, and get many such requests...)

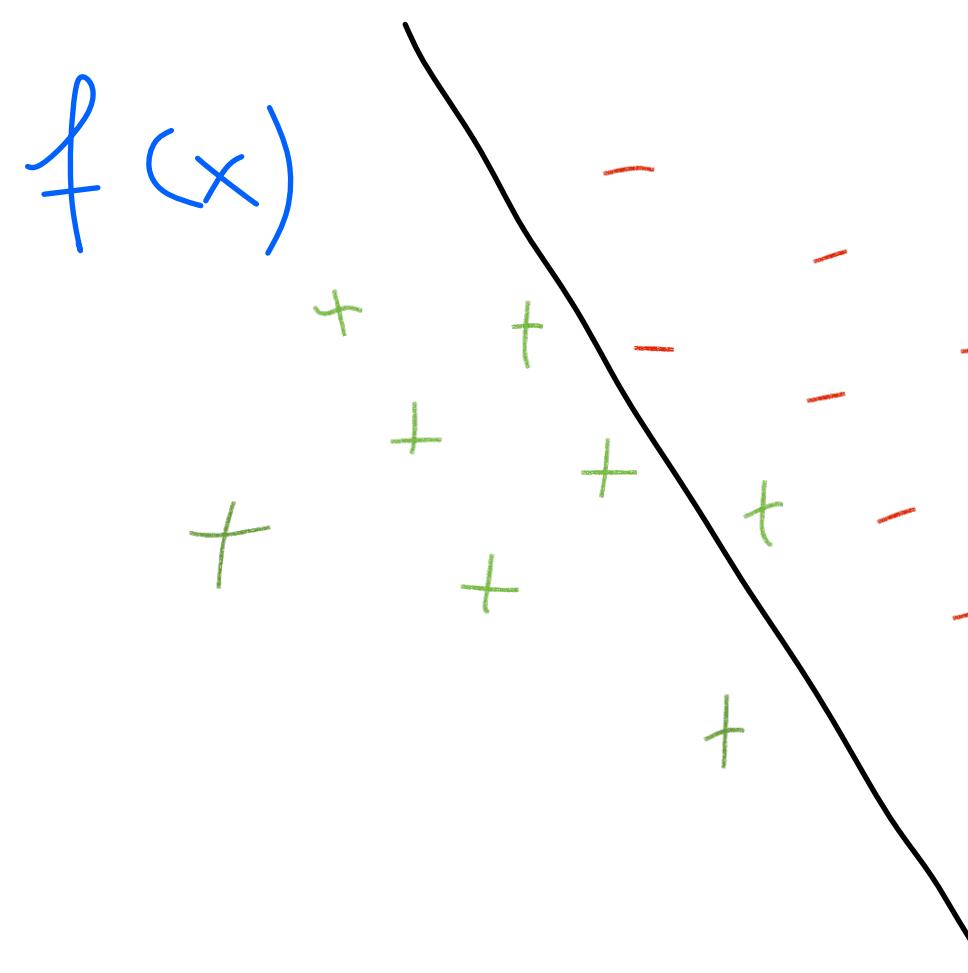
# The NLP Pipeline

This Lecture!



# It's all about the Representation

- Machine Learning is simple when you know how to represent the data
- Eg. the representation separates positive from negative text



# Tokenisation

Machine sees string as sequence  
of characters, no sense of "words"

Mr. Jones and me.

First step: break up string into tokens

"Easy" for English! Use whitespace plus a few rules...

今日もしないといけない。 for Japanese?

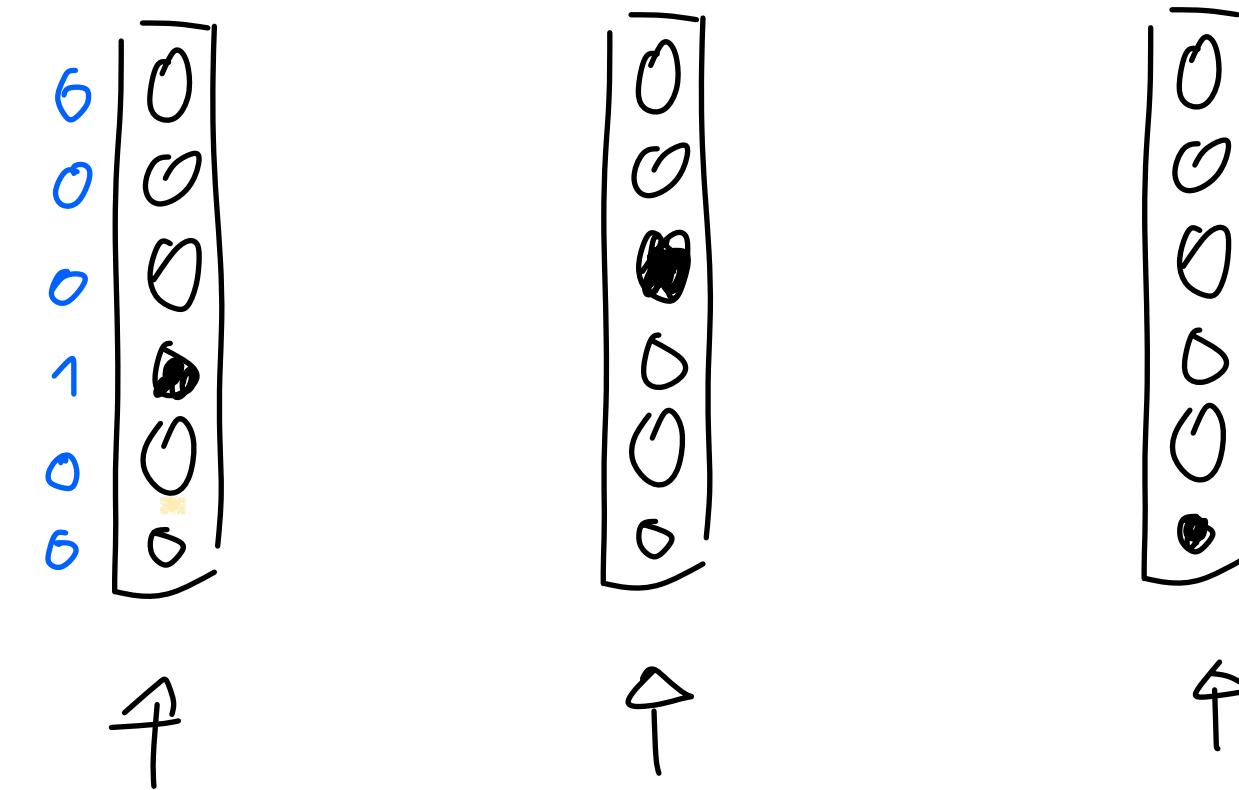
today too don't do it if doesn't work

# In other Domains

We developed a nanocarrier system of  
herceptin-conjugated nanoparticles  
of d-alpha-tocopheryl-co-poly(ethylene  
glycol) 1000 succinate (TPGS)-cisplatin  
prodrug ...

Need rules and/or machine learning

# Representing Words with One-Hot Vectors



I like UCL

word  
↓

$$f(w) \in \mathbb{R}^m \quad f_i(w) = \begin{cases} 1 & \text{if } w = v_i \\ 0 & \text{otherwise} \end{cases}$$

$$v = \begin{bmatrix} 0 & a \\ 1 & the \\ 2 & at \\ \vdots & \vdots \\ m & York \end{bmatrix}$$

Vocabulary

# Representing Sentences

Sum  
Pooling

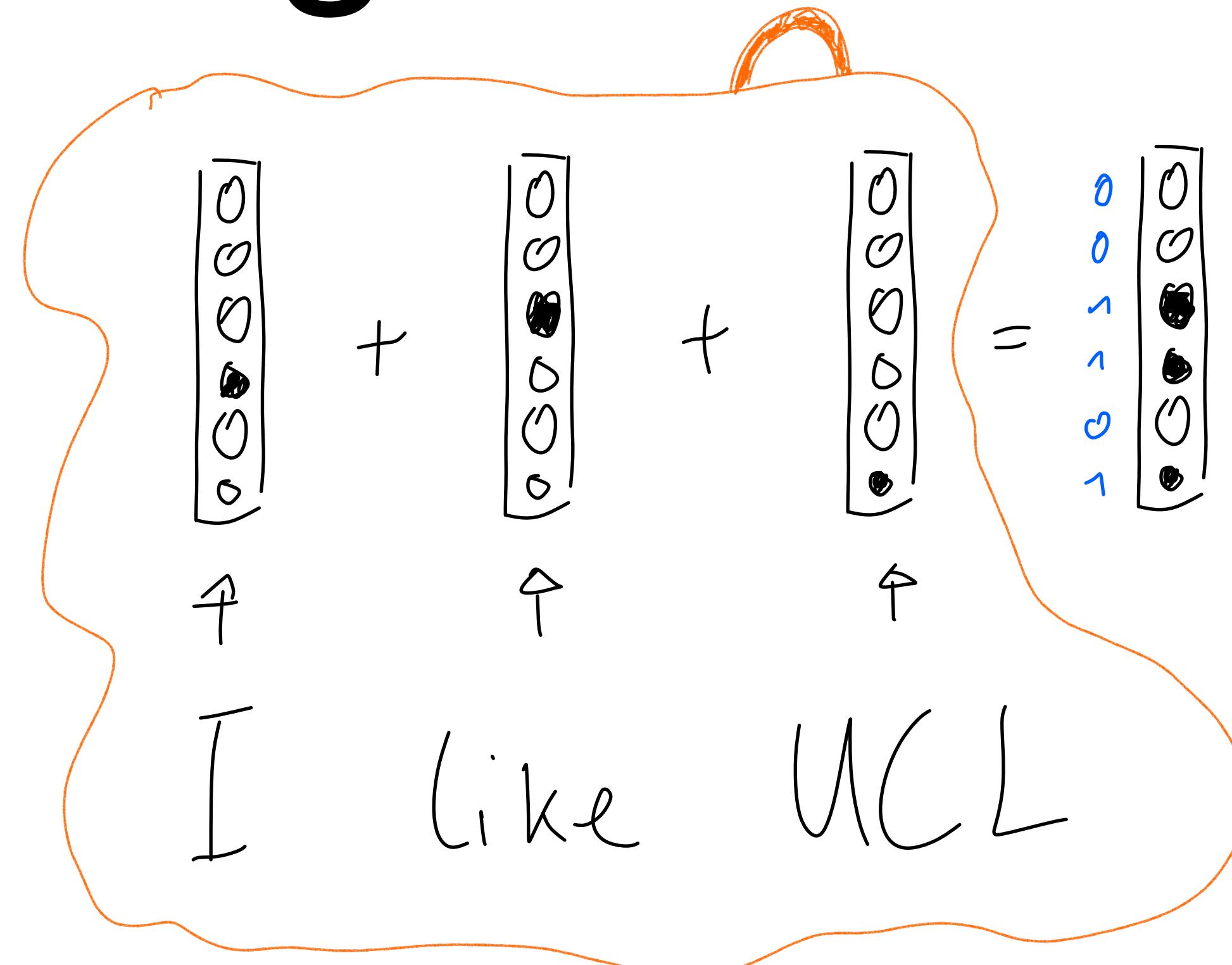
$$\begin{array}{c} \begin{array}{c} | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \end{array} + \begin{array}{c} | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \end{array} + \begin{array}{c} | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \end{array} = \begin{array}{c} | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \\ | \\ 0 \end{array} \\ \downarrow \quad \uparrow \quad \uparrow \end{array}$$

I like UCL

Sentence  
↓  
 $f(x) \in \mathbb{R}^m$

↙ token representation  
 $f(x) = \sum_{w \in x} f(w)$

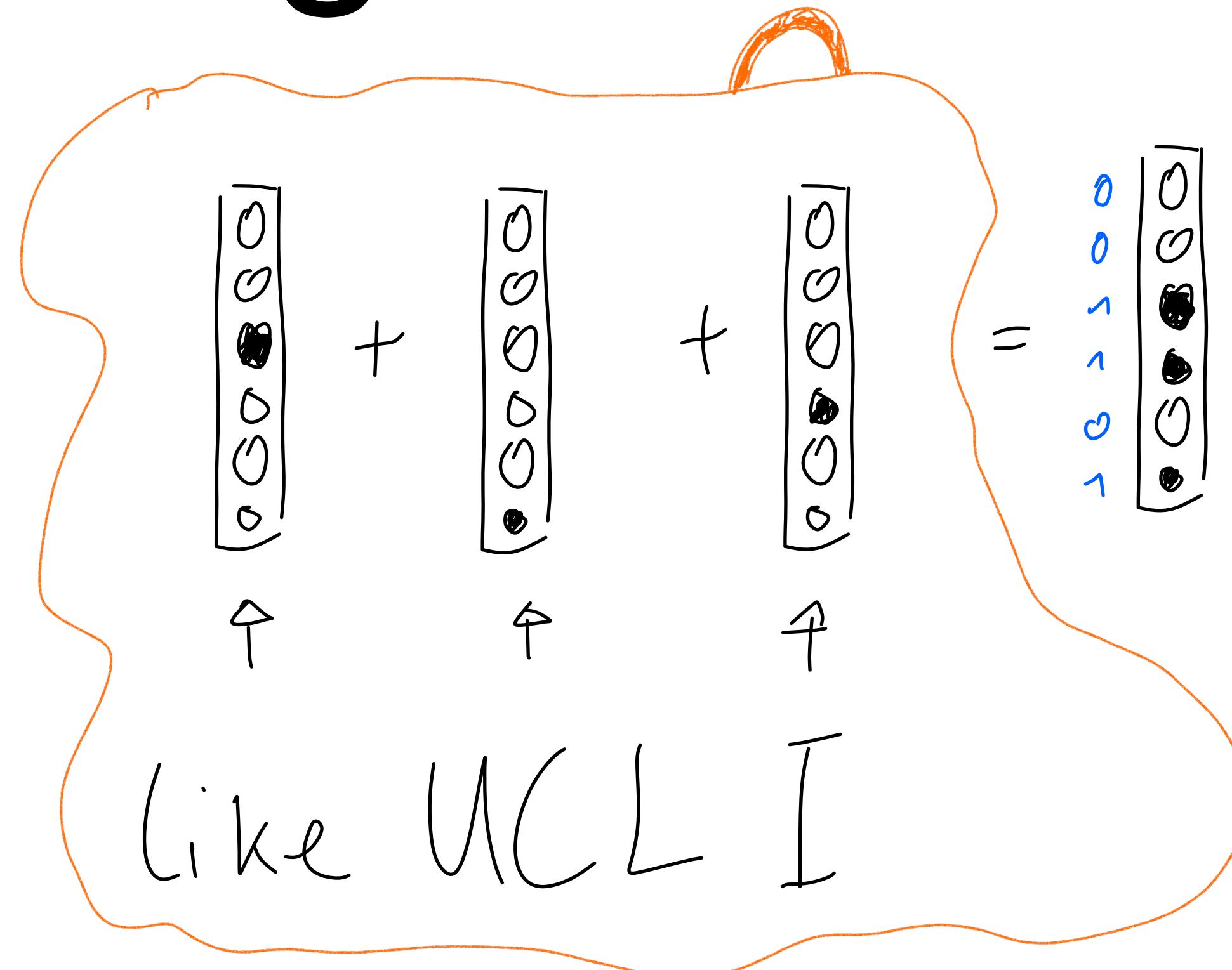
# “Bag of Words”



$$f(x) \in \mathbb{R}^m$$

$$f(x) = \sum_{w \in x} f(w)$$

# “Bag of Words”



$$f(x) \in \mathbb{R}^m$$

$$f(x) = \sum_{w \in x} f(w)$$

# Aggregation

$$\begin{array}{c} \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \\ + \end{array} \quad \begin{array}{c} \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \\ + \end{array} \quad \begin{array}{c} \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \\ + \end{array} \quad \begin{array}{c} \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{array} \\ = \end{array} \quad \begin{array}{c} \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{array} \\ + \end{array}$$

I, like, like UCL sensitive  
+ to sentence length

# Sum Poolings

# Aggregation

$$\begin{array}{c} \begin{array}{ccccc} \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} & + & \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} & + & \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} & + & \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} & = & \begin{array}{c} | \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \\ \downarrow & \uparrow & \uparrow & \uparrow & & & & & \end{array} \\ \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} & & \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} & & \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} & & \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} & & \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} \\ \times & & \frac{1}{4} & & & & & & \end{array}$$

I, like, like UCL

Mean Pooling

# Aggregation

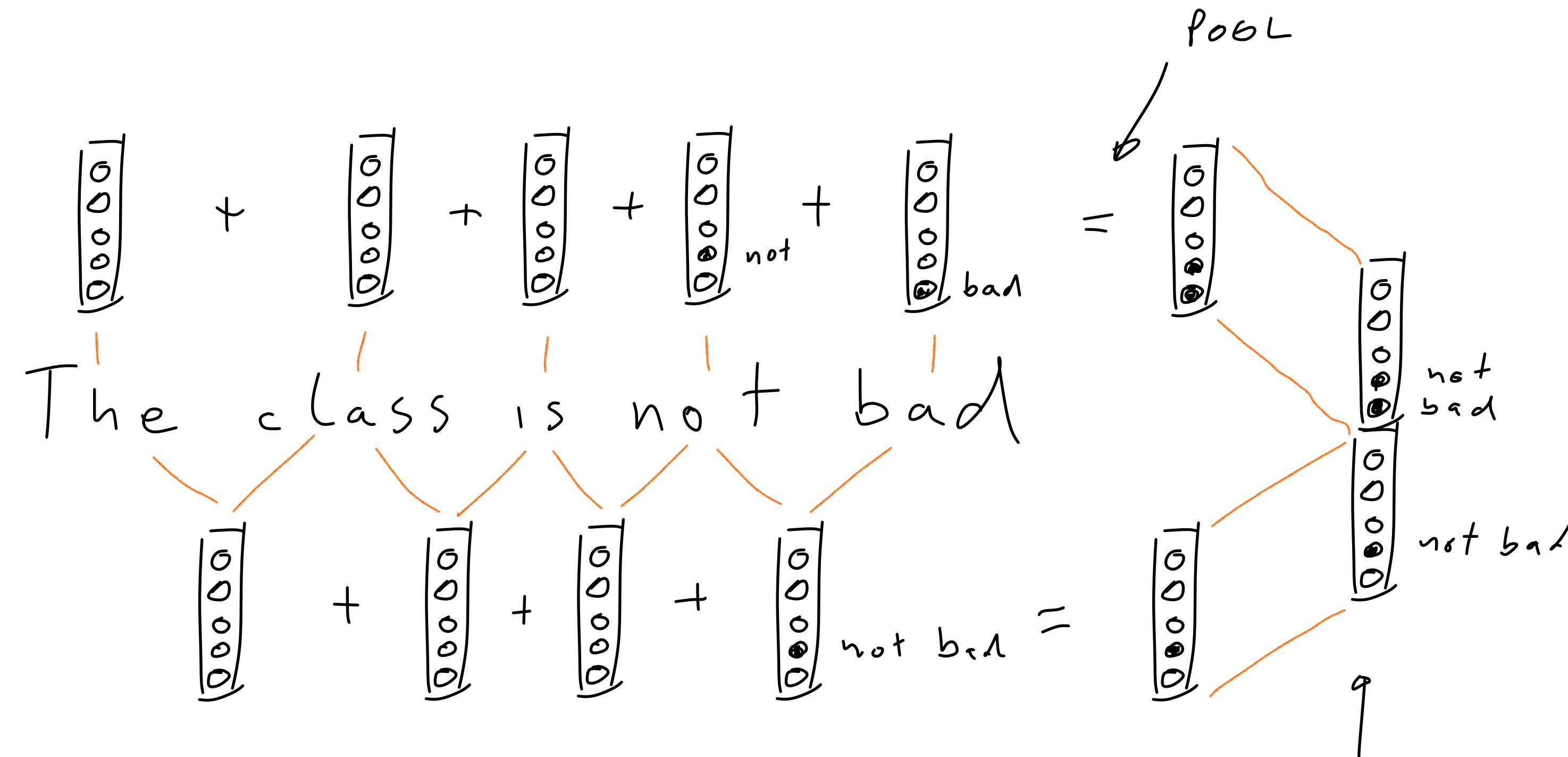
$$\max \left( \begin{matrix} \begin{array}{|c|c|c|c|} \hline & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline \end{array} & , & \begin{array}{|c|c|c|c|} \hline & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline \end{array} & , & \begin{array}{|c|c|c|c|} \hline & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline \end{array} & , & \begin{array}{|c|c|c|c|} \hline & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline & 0 & 0 & 0 \\ \hline \end{array} \end{matrix} \right) = \begin{array}{|c|c|c|c|} \hline & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \textcolor{red}{1} & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline 1 & 0 & 0 & 0 \\ \hline \end{array}$$

↑      ↑      ↑      ↑

I, like, like UCL

Max Pooling

# More Engineering



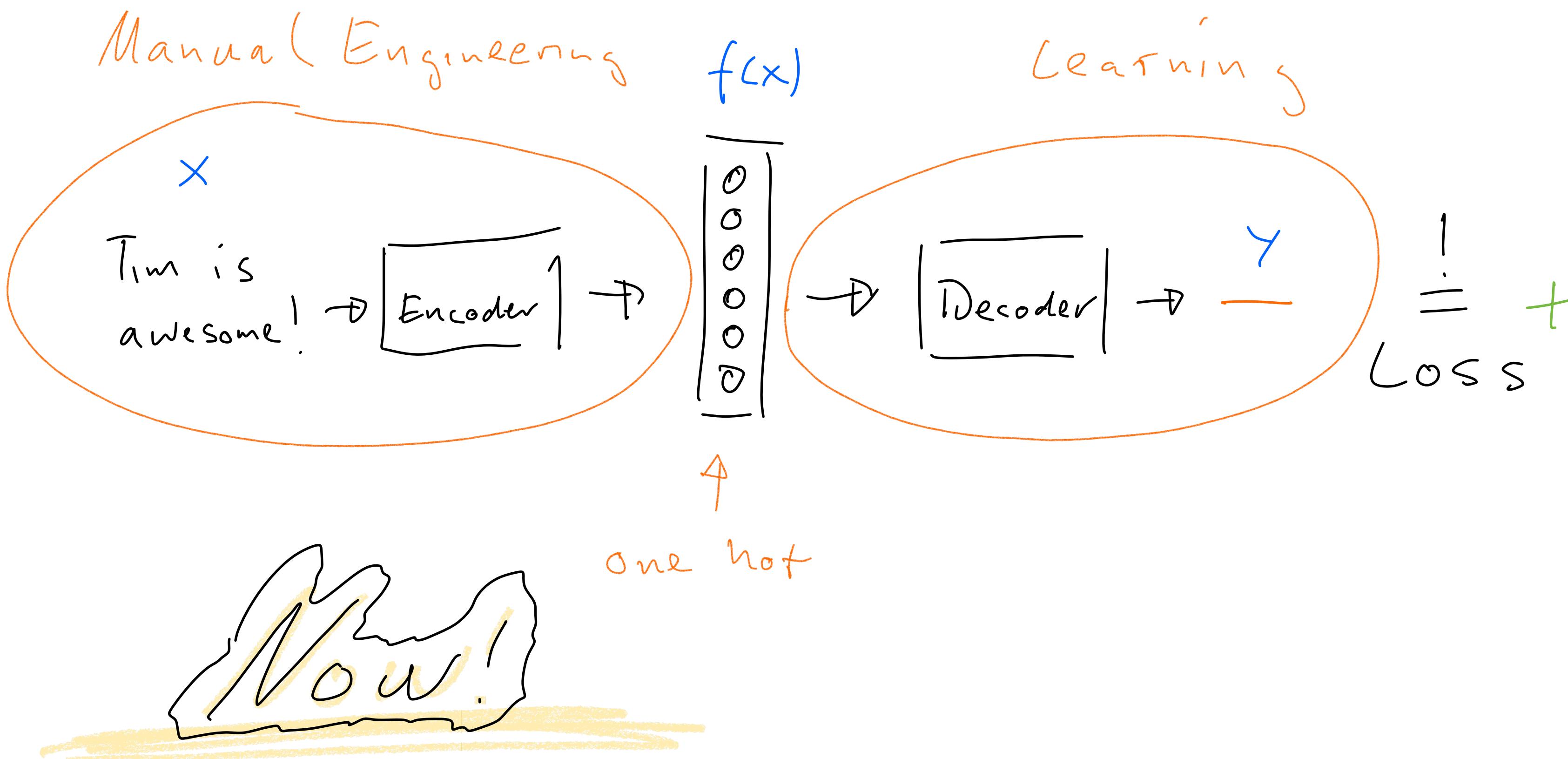
uniagram features may not be enough  
engineer more features! E.g. **bigramm!**

# More Ideas?

- ① Use dictionaries ?
  - ② Use syntax ?
  - ③ Preprocessing ?
- ⋮

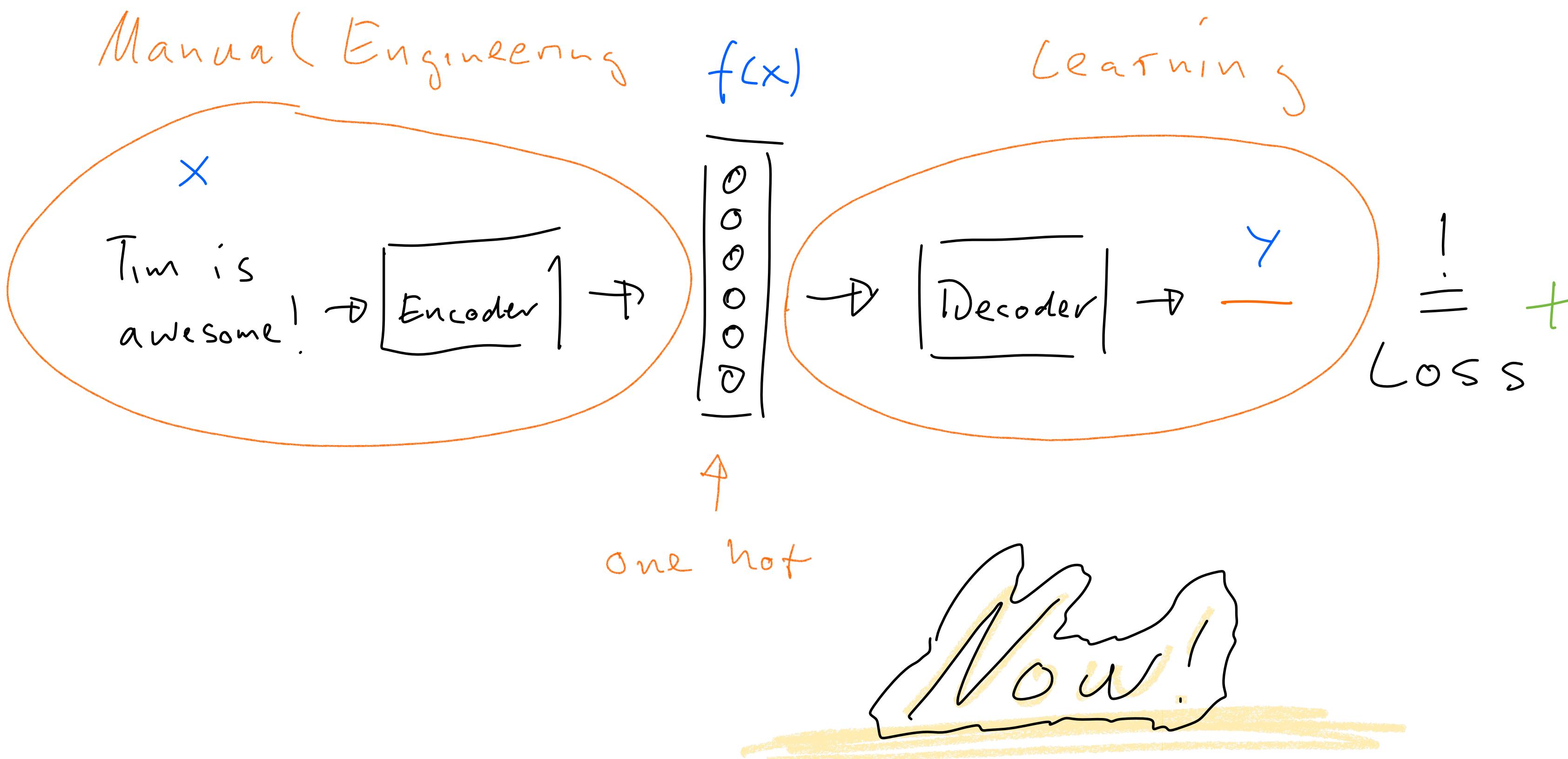
# The NLP Pipeline

This Lecture!



# The NLP Pipeline

This Lecture!



# Linear Classification

① Assign score to +

$$\theta^T \times f(x) = s_\theta(x)$$

class weights      representation      score/logits

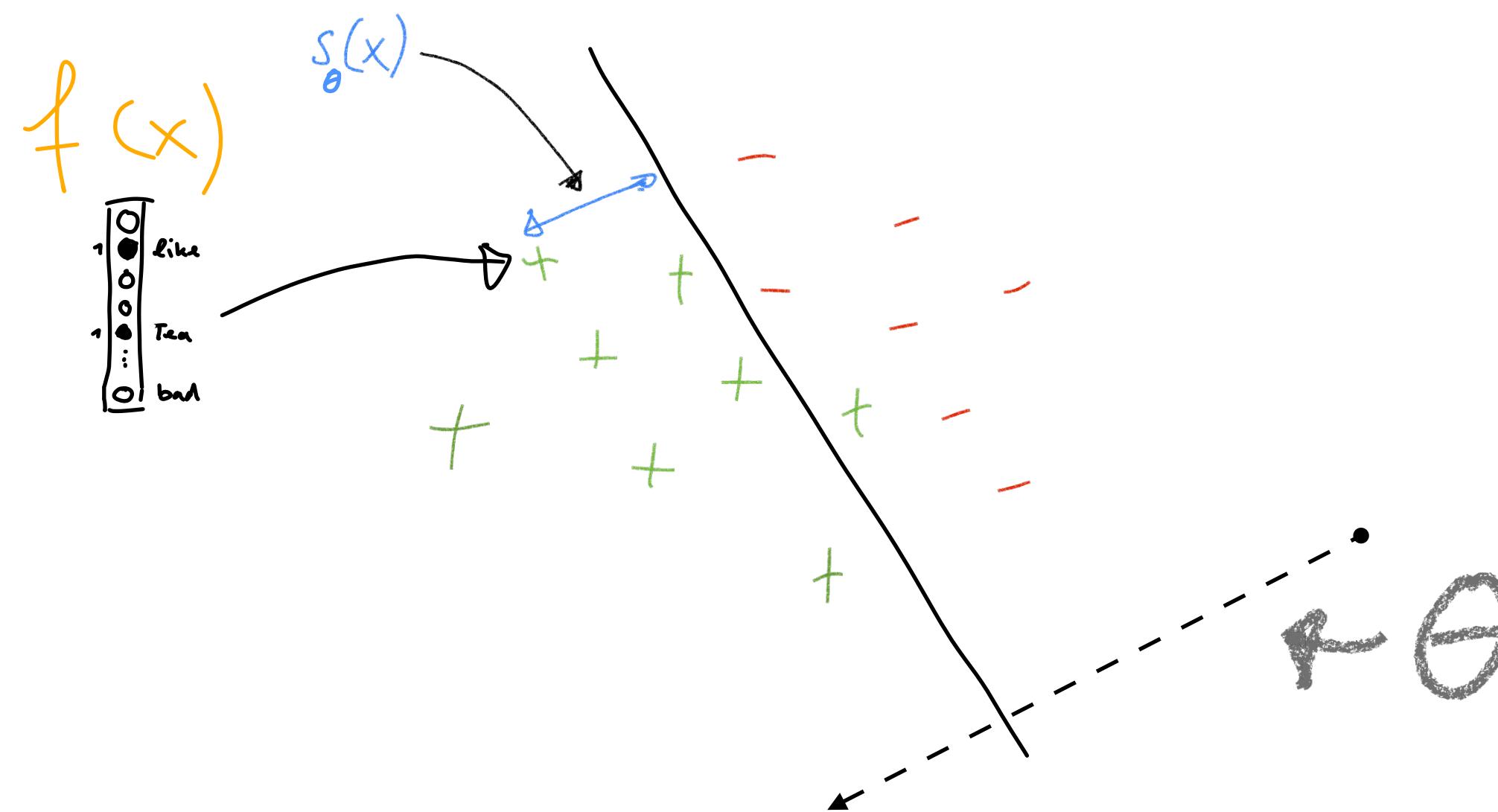
$$\begin{pmatrix} \text{like} & \text{Tea} & \text{bad} \\ \hline 1.5 & 0.1 & -1.1 \\ 0 & 0 & \cdots & 0 \end{pmatrix} \times$$

$$\begin{array}{c|ccc} & 0 & 1 & \dots & n \\ \hline 1 & 0 & 1 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{array}$$

like  
Tea  
bad

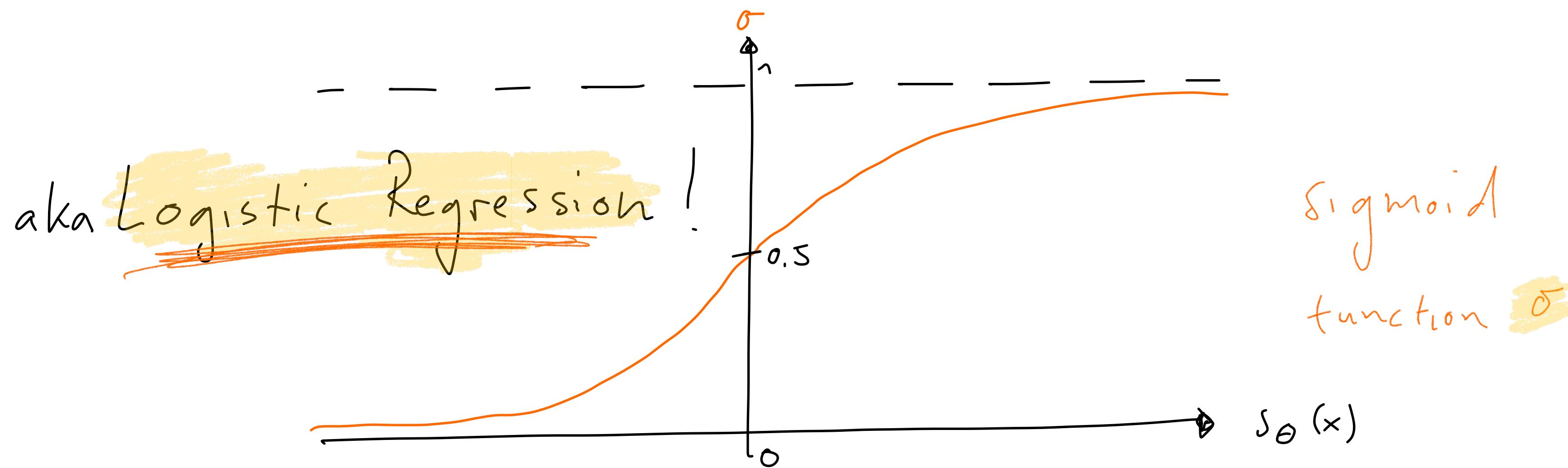
$$= 1.5 \times 1 + 0.1 \times 1 = 1.6$$

# Decision Boundary



# Probabilistic Output

② Assign probabilities to +, -



$$P_\theta(y=+|x) = \sigma(s_\theta(x)) = \frac{e^{s_\theta(x)}}{1 + e^{s_\theta(x)}}$$

$$P_\theta(y=-|x) = 1 - P_\theta(y=+|x)$$

# Search

- ③ Choose label with highest probability

Easy for two classes:

- ① calculate  $p_{\theta}(+ | x)$
- ② "  $p_{\theta}(- | x)$

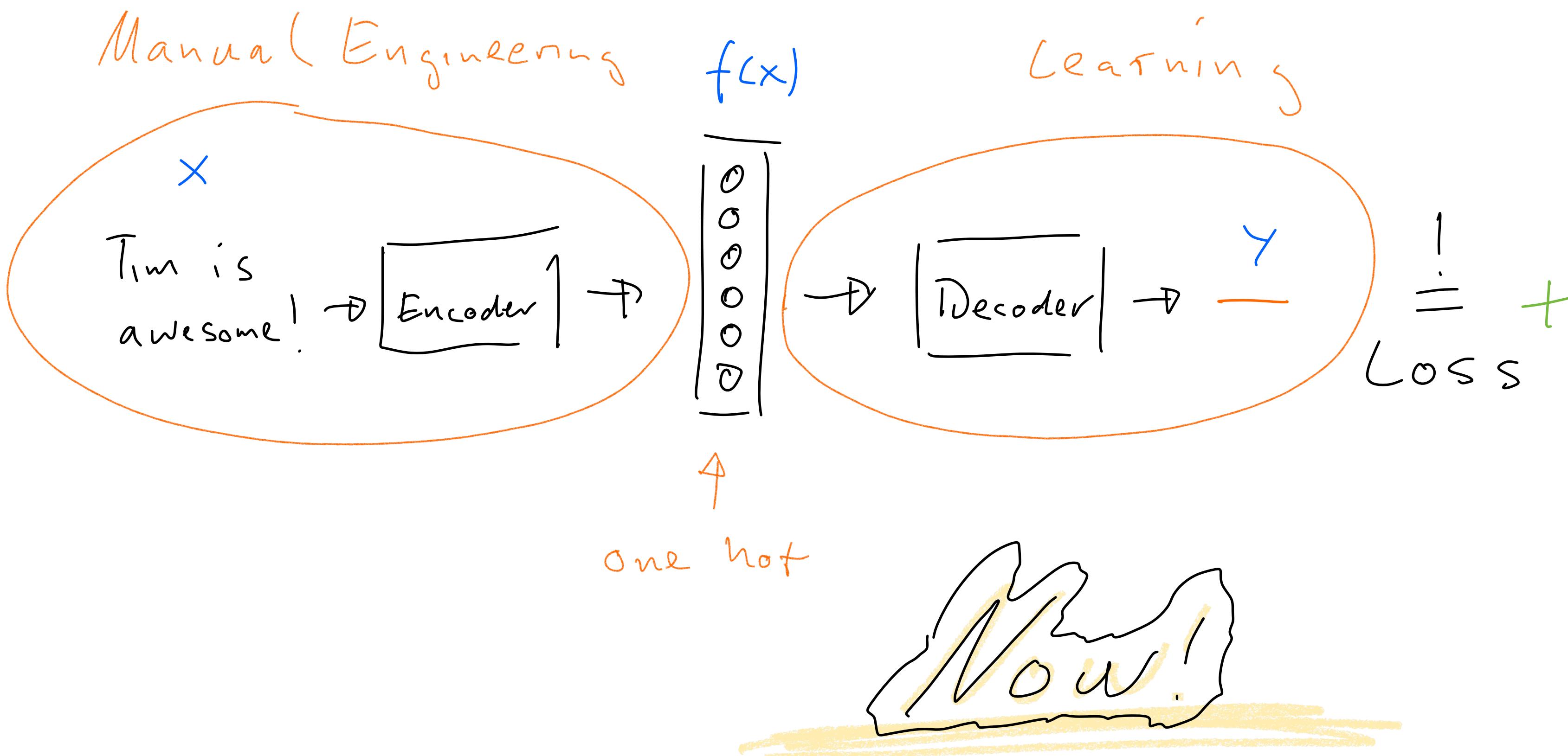
- ③ Choose higher one

$$y^* = \operatorname{argmax}_{y \in \{+, -\}} p_{\theta}(y | x)$$

But consider  
Machine Translation

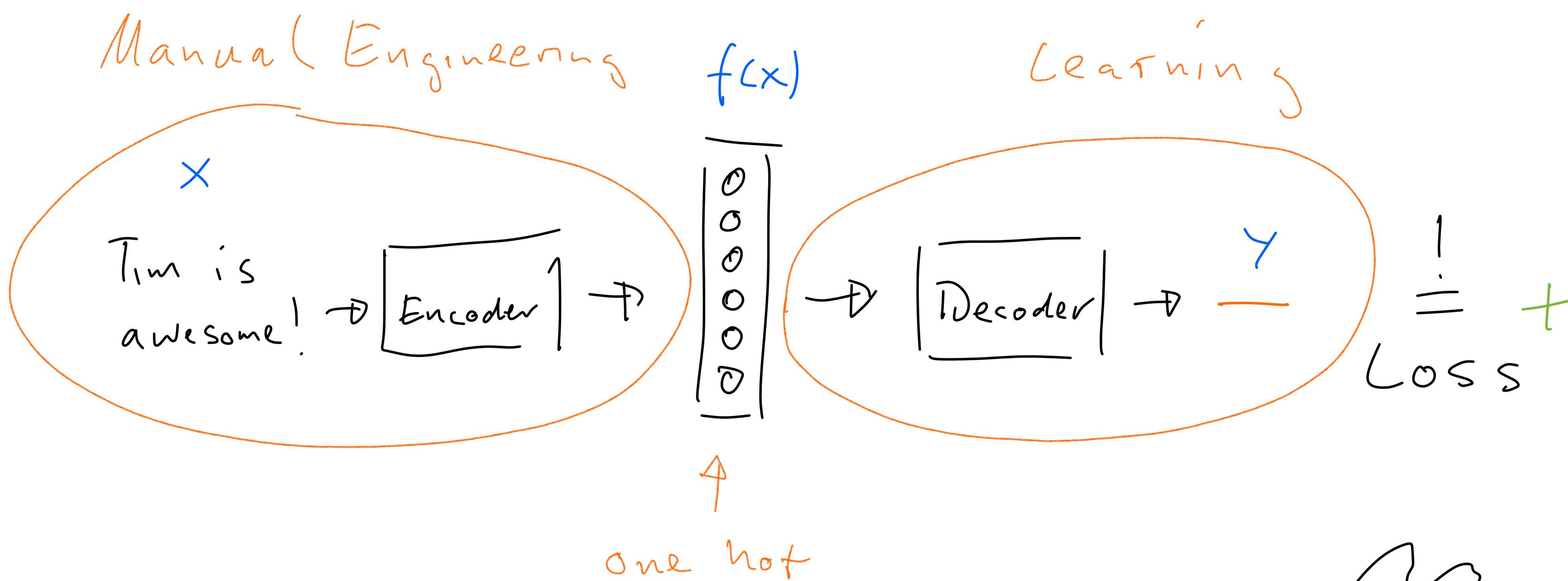
# The NLP Pipeline

This Lecture!



# The NLP Pipeline

This Lecture!



# Training Loss

Dataset

$$D = (x_1, y_1) \dots (x_n, y_n)$$

Loss

$$L(D, \theta) = \frac{1}{n} \sum_{i=1}^n l(x_i, y_i, \theta)$$

per instance loss

$$l(x, y, \theta) = -\log p_\theta(y | x)$$

Conditional  
Loglikelihood

# Cross Entropy Loss

let's think of label as distribution  $\tilde{p}_i(y)$

$$\tilde{p}_i(y) = \begin{cases} 1 & \text{if } y_i = + \\ 0 & \text{otherwise} \end{cases}$$

then

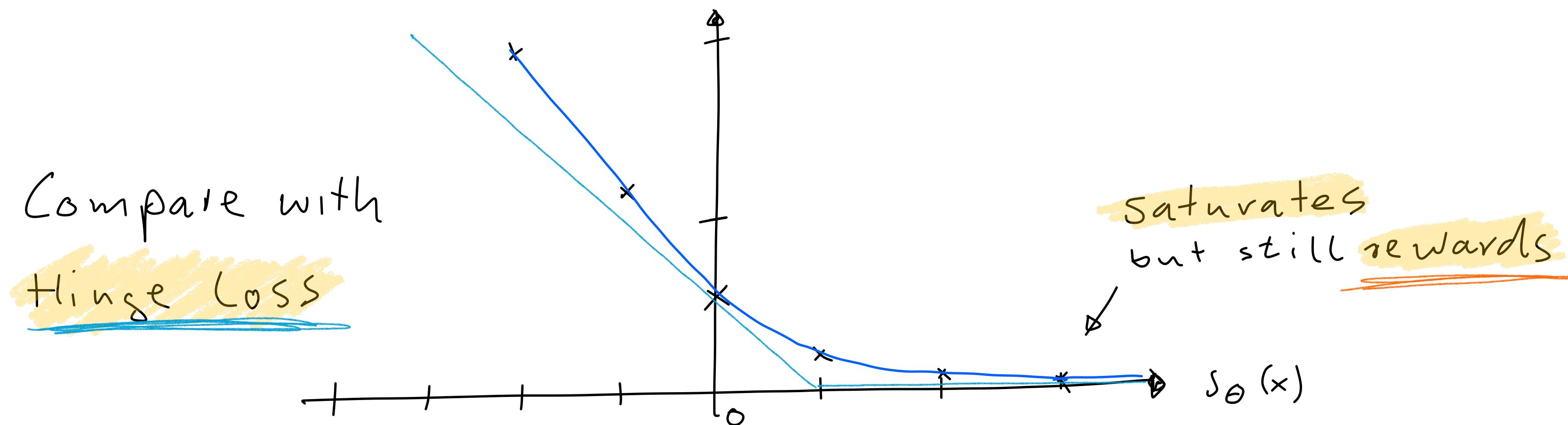
$$l(x_i, \tilde{p}_i, \theta) = - \sum_{y \in \Sigma^+ \cup \{-\}} \tilde{p}_i(y) \cdot \log p_\theta(y|x)$$

can be generalized  
to soft labels!

Cross Entropy  $H(\tilde{p}, p)$

# Intuition

$$\ell(s_\theta(x), +) = -\log \left( \frac{e^{s_\theta(x)}}{1 + e^{s_\theta(x)}} \right)$$

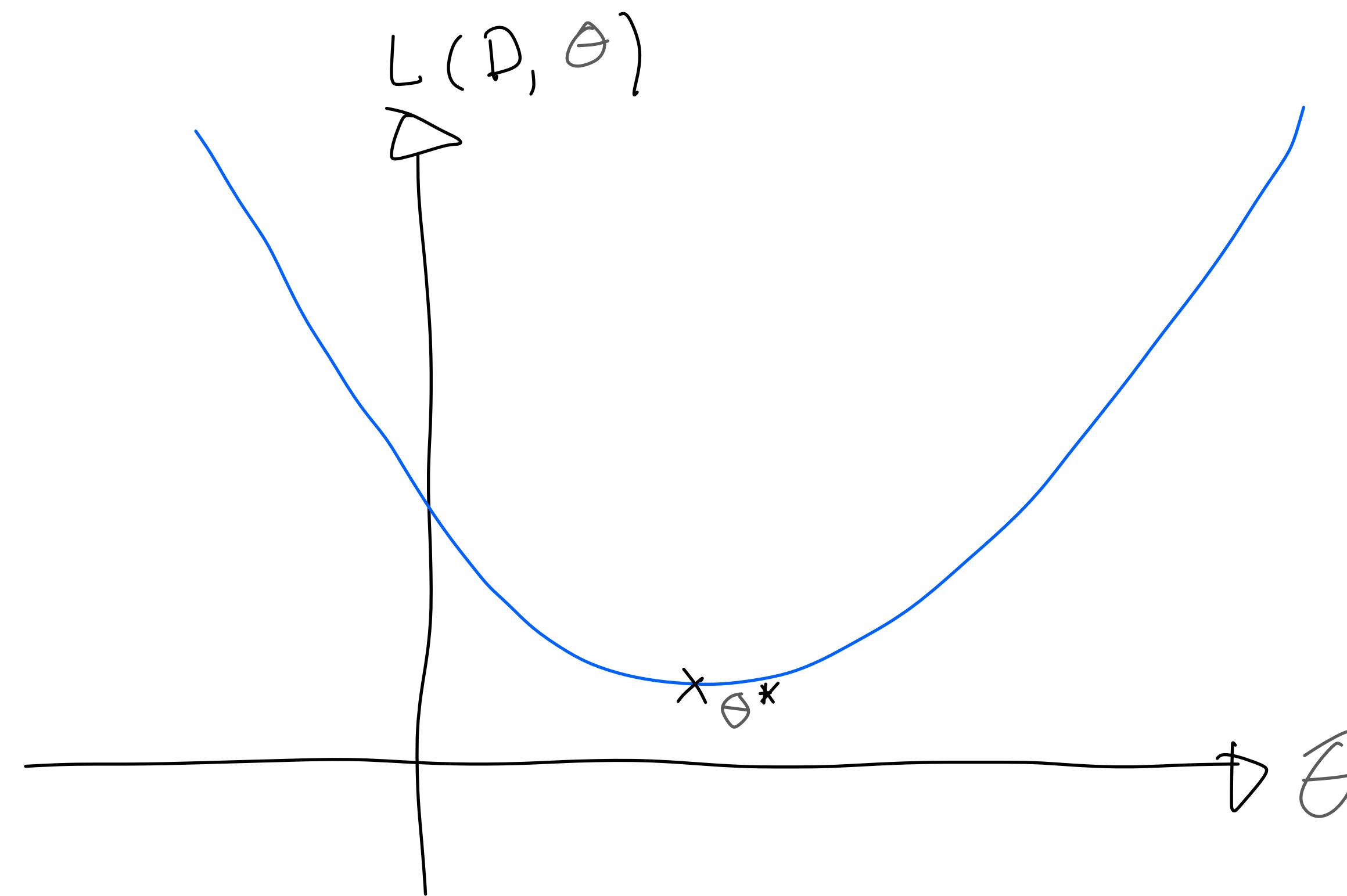


# Training

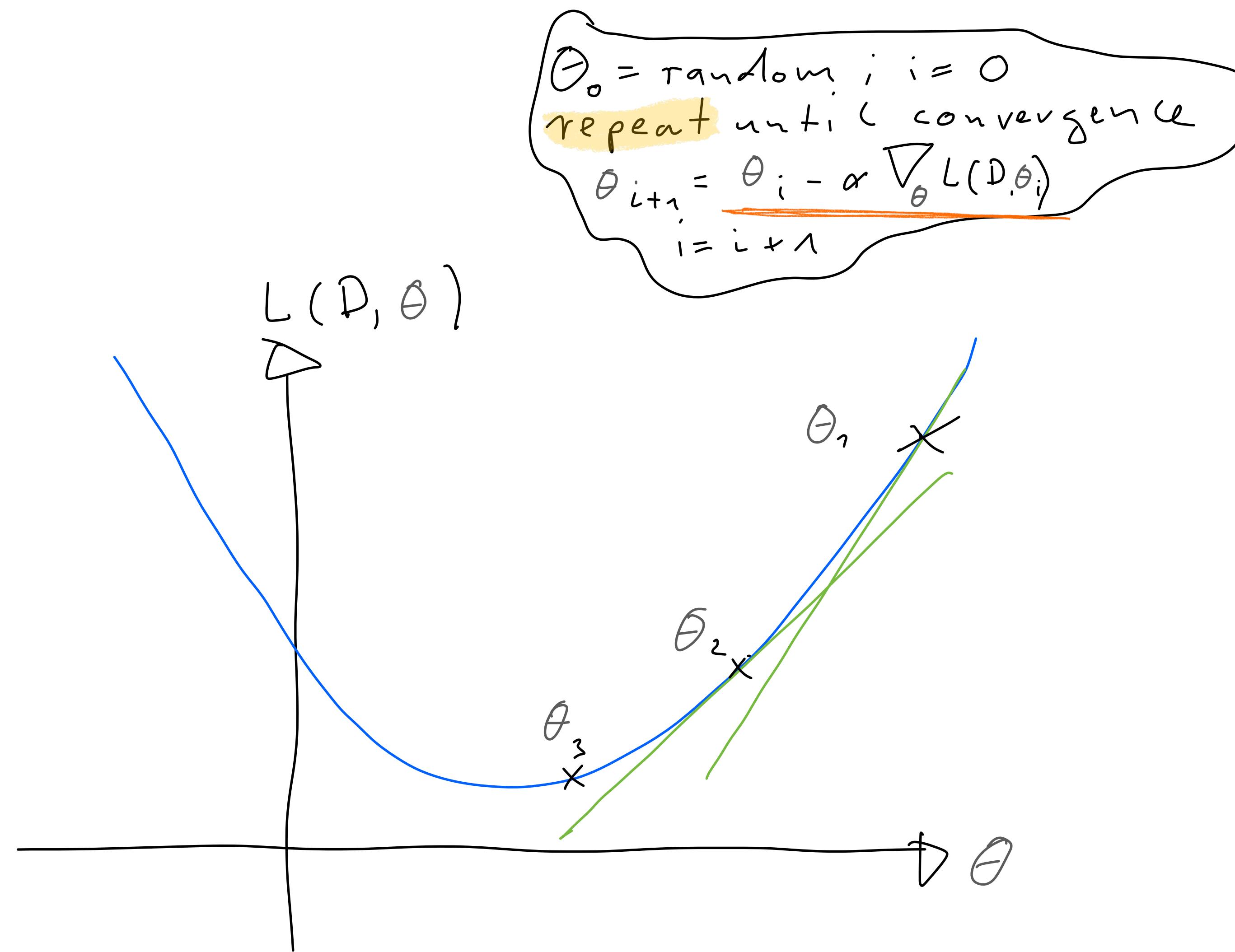
$$\theta^* = \underset{\theta \in \mathbb{R}^m}{\operatorname{arg\,min}} L(D, \theta)$$

best parameter  $\theta^*$  →

Loss ↑



# Gradient Descent



# Stochastic Gradient Descent

$$\nabla_w L(\rho, \theta) = \nabla_w \frac{1}{m} [l(x_1, y_1, \theta) + \dots + l(x_n, y_n, \theta)]$$

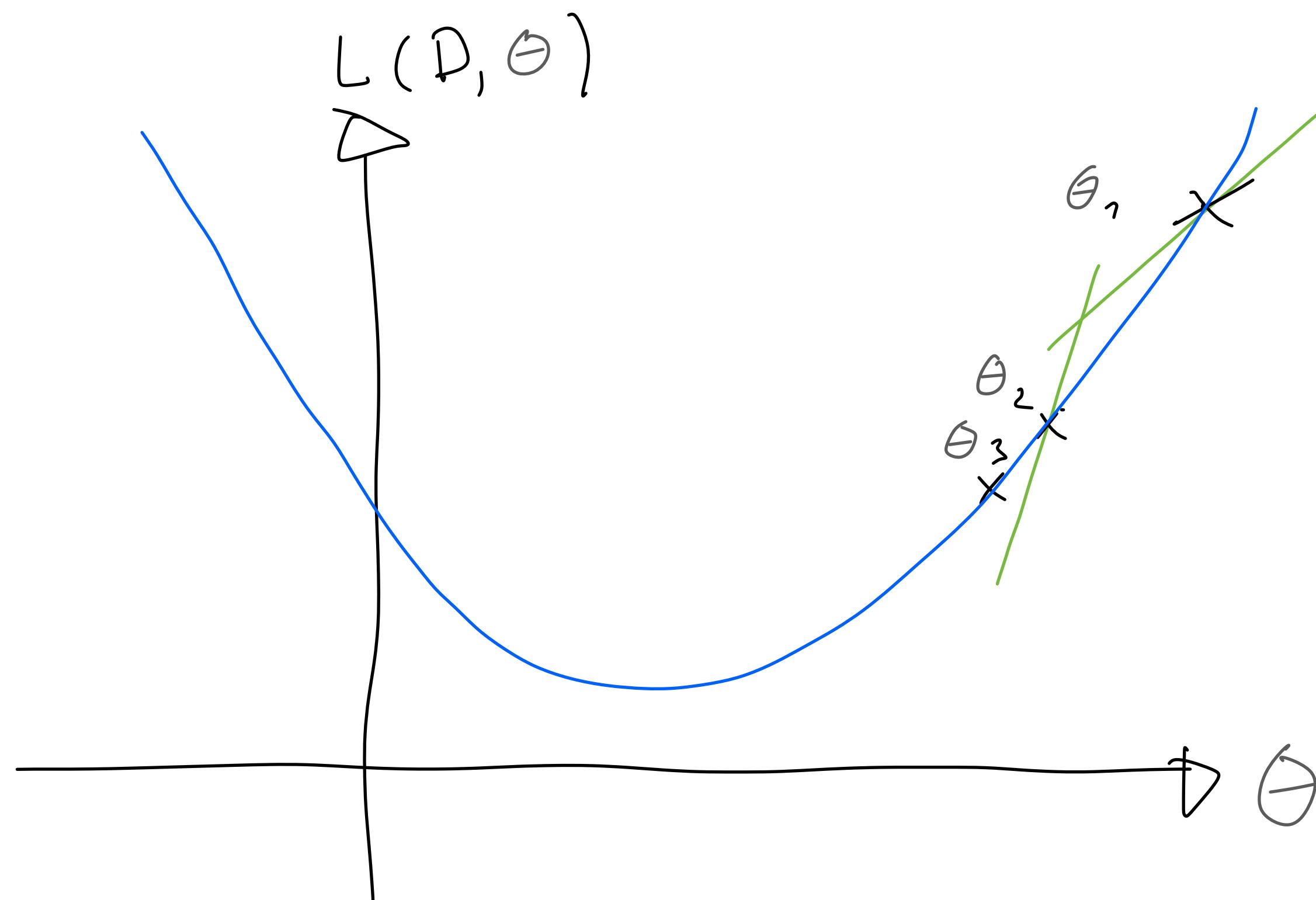
long time spent w/o update  
using **stale** weights

Instead: approximate gradient via sampling **single**  
**instance j**

$$\nabla_w L(\rho, \theta) \approx \nabla_w l(x_j, y_j, \theta)$$

# Stochastic Gradient Descent

$$\theta_{i+1} = \theta_i - \alpha \nabla_{\theta} l(x_j, y_j, \theta_i)$$



# Regularisation

$\theta_1^*$	$\theta_2^*$	
1	0	good
-1	-1	bad
0.5	0	like
:	:	:
0	1	good class
0	1	good module
0	1	good lecturer
:	:	:
0	1	this was a good class
0	1	this is a good class

which is better?

$$L(\theta) \quad 0.02$$

$$\|\theta\|_2^2 \quad 2.25$$

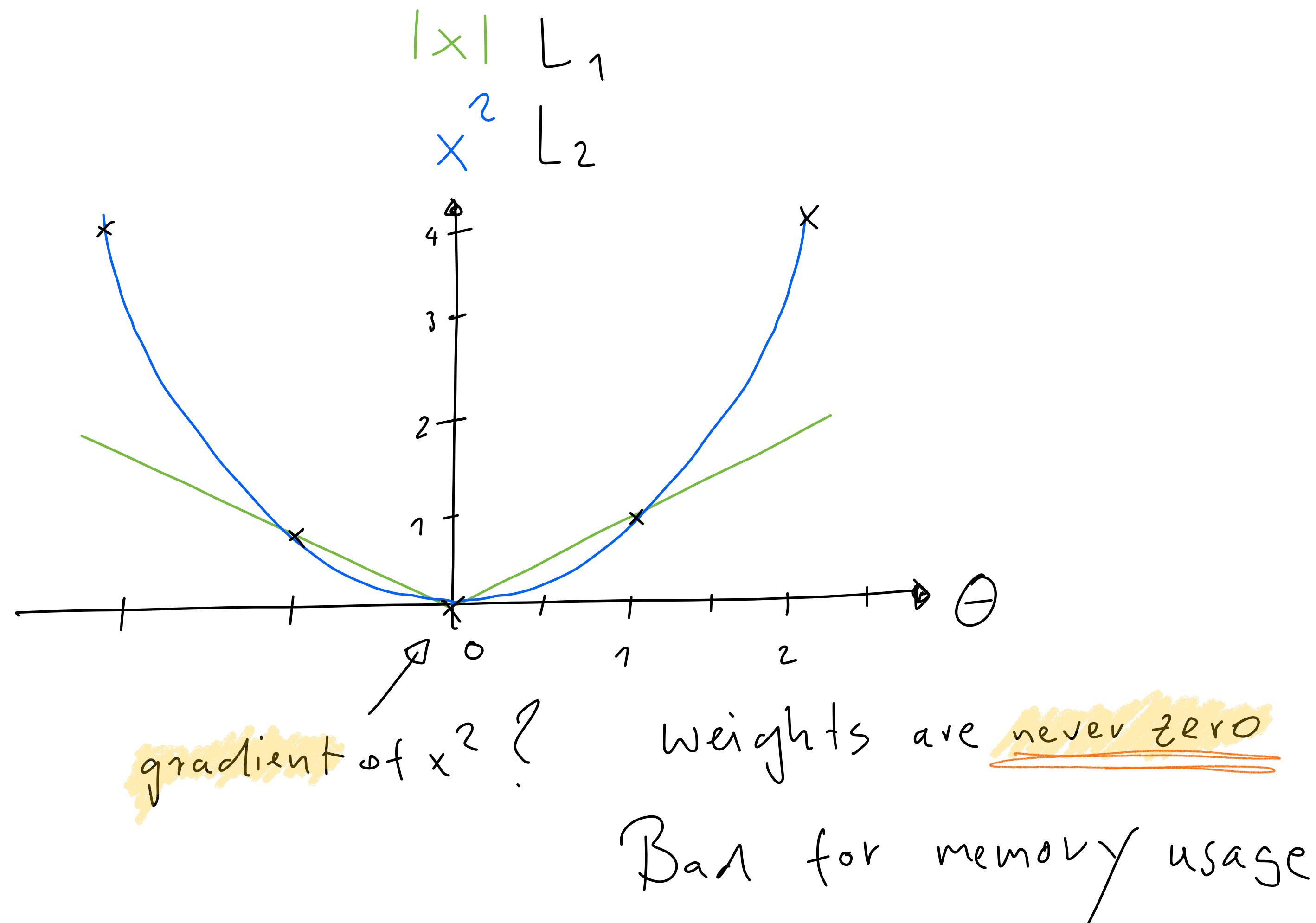
$$0.02$$

$$30.55$$

$$\rightarrow \hat{L}_{\alpha}(\theta) = L(\theta) + \alpha \|\theta\|_2^2$$

$L_2$  regularization

# L2 vs L1



# Material

- Stat-NLP-Book
  - Tokenization
  - Text Classification
- Goldberg, chapter 3
- JM, chapter 2, 5