

# Lecture 3

## NN Classification Notation:

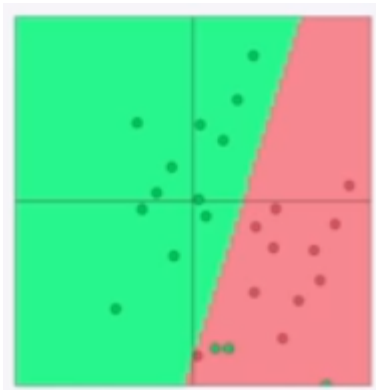
### 2. Classification setup and notation

- Generally we have a **training dataset** consisting of **samples**

$$\{x_i, y_i\}_{i=1}^N$$

- $x_i$  are **inputs**, e.g. words (indices or vectors!), sentences, documents, etc.
  - Dimension  $d$
- $y_i$  are **labels** (one of  $C$  classes) we try to predict, for example:
  - classes: sentiment, named entities, buy/sell decision
  - other words
  - later: multi-word sequences

Example of a learned line:



$X_i$  : tells its position

$Y_i$  : tells if it is green or red

## Negative log Probability with Softmax function:

In logistic regression with the above example, we have the  $X_i$  which we multiply with a learned Weights and put the results into a sigmoid function and if the probability is less than 0.5 we classify it as case 1 else case 2. Similarly we extend this to  $N$  cases with  $N-1$  weights.

On the other hand the softmax probability function has weights for every case separately. For example, in the above red-green example with two cases, in softmax probability method we don't calculate one value and classify the instance as one of the cases depending on if the calculated number is closer to zero or one, instead we calculate two probabilities with two weights each corresponding to one of the cases and classify it as the case with a higher probability.

$$p(y|x) = \frac{\exp(W_y \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)} = \frac{\exp(f_y)}{\sum_{c=1}^C \exp(f_c)} = \text{softmax}(f_y)$$

Where the given probability is the probability for the case y and C is the total number of cases.

Note: here  $W_y$  is a **row** the weight matrix w corresponding to a single case.

We take negative log of this probability to minimize because

- Negative so that we can minimize instead of maximize
- Log just to squish the probability around.

Note: **NLL** : negative log likelihood.

## Cross entropy model:

### Cross Entropy loss:

Let the true probability be p

Let the probability calculated by our model be q.

We use a function H called cross-Entropy function to calculate the error in calculated probabilities across all the classes for a case.

$$H(p, q) = - \sum_{c=1}^C p(c) \log q(c)$$

Notice how when both p and q are either together zero or one, H function disappears. Hence H function beautifully calculates the difference between p and q.

Notice:

- Assuming a ground truth (or true or gold or target) probability distribution that is 1 at the right class and 0 everywhere else:  $p = [0, \dots, 0, 1, 0, \dots, 0]$  then:
- Because of one-hot  $p$ , the only term left is the negative log probability of the true class

Therefore:

$$H(p, q) = -\log(q_i)$$

Where  $i$  is the true case / the only case where  $p$  is 1.

The above explanation is for a single point, if we have a full dataset,

$$H(p, q) = ( -\sum \log(q_i) ) / N : \text{summed over the entire dataset,}$$

With substitution of softmax probability  $q_i$  formula,

Matrix optimization:

- Instead of

$$f_y = f_y(x) = W_{y \cdot} x = \sum_{j=1}^d W_{yj} x_j$$

We will write  $f$  in matrix notation:

$$f = Wx$$

- For general machine learning  $\theta$  usually only consists of columns of  $W$ :

$$\theta = \begin{bmatrix} W_{\cdot 1} \\ \vdots \\ W_{\cdot d} \end{bmatrix} = W(:,) \in \mathbb{R}^{Cd}$$

- So we only update the decision boundary via

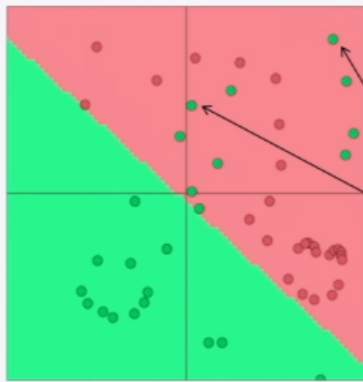
$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{\cdot 1}} \\ \vdots \\ \nabla_{W_{\cdot d}} \end{bmatrix} \in \mathbb{R}^{Cd}$$

Here each weight  $w_i$  is a row vector consisting of number of elements corresponding to the number of elements in the input  $x$ .

Now do normal / stochastic gradient decs etc.

## Why do we need neural network classifiers ?

- Softmax ( $\approx$  logistic regression) alone not very powerful
- Softmax gives only linear decision boundaries



This can be quite limiting

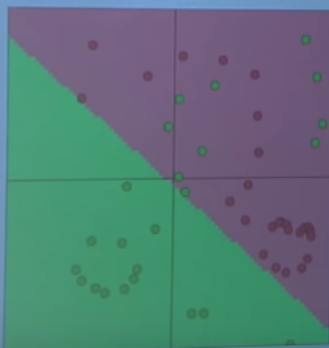
→ Unhelpful when a problem is complex

Wouldn't it be cool to get these correct?

### Neural Nets for the Win!

- Neural networks can learn much more complex functions and nonlinear decision boundaries!

• In original space



## Why do we need Non-Linearity functions ?

Non-Linearity functions are functions like sigmoid function that are applied to  $z_i$  to make them  $a_i$ . We call the sigmoid squishification function an activation function because in analogy it does the work an activation potential does in a neuron (The all or null effect).

They are needed because without them these multilayer neural networks will basically still just produce linear functions as extra layers can just be compiled down to single linear function.  $w_1 w_2 x = Wx$ .

## Named Entity Recognition:

As the name suggests it is used to find named entities like name of organisations , people, countries etc. it has many uses like keeping track of how often a company is mentioned, or in answering questions.

Common way of doing it is using a classification nn.

Tricky cases:

- Hard to work out boundaries of entity

### **First National Bank Donates 2 Vans To Future School Of Fort Smith**

POSTED 3:43 PM, JANUARY 11, 2015, BY 5NEWS WEB STAFF

Is the first entity “First National Bank” or “National Bank”

- Hard to know if something is an entity  
Is there a school called “Future School” or is it a future school?
- Hard to know class of unknown/novel entity:

- Entity class is ambiguous and depends on context

“Charles Schwab” is PER  
not ORG here! 🙅

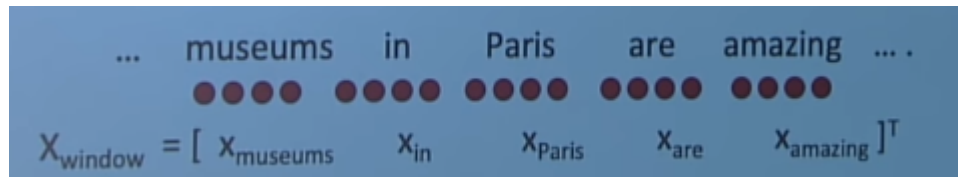
where Larry Ellison and Charles Schwab can  
live discreetly amongst wooded estates. And

- Interesting problems like ambiguity arise in context!
- Example: auto-antonyms:
  - "To sanction" can mean "to permit" or "to punish"
  - "To seed" can mean "to place seeds" or "to remove seeds"
- Example: resolving linking of ambiguous named entities:
  - Paris → Paris, France vs. Paris Hilton vs. Paris, Texas
  - Hathaway → Berkshire Hathaway vs. Anne Hathaway

## Solution:

The solution to the dependency on context for the class of a named entity is passing nearby word's word vector along with word vector in question as arguments to the neural network.

Example, if we want to classify paris in a particular context, lets just pass 2 words b4 and after it to the neural net.



This  $X_{\text{window}}$  is the input to the NN.

## Important points:

1. In deep learning in NLP we not only learn the weights but we can also learn the word vectors together.