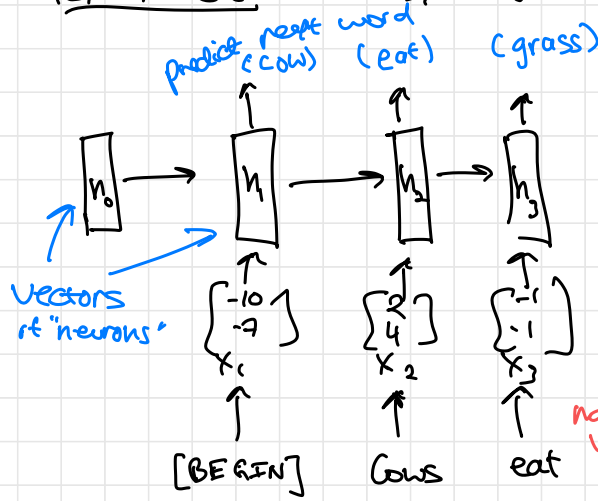


# 10/5/2023: RNNs, Review for Midterm



multiclass classification  
 $\approx$  Soft max regression

RNN recursion

$$h_t = \tanh(\underbrace{w^h}_{\text{matrix}} h_{t-1} + \underbrace{w^x}_{\text{matrix}} x_t + \underbrace{b}_{\text{vector}})$$

map each word to vector

learned by gradient descent

## Model Params

$w$ -weight vector for each class  
 #classes = #vocab words

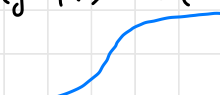
$w^h, w^x$  matrix  
 $h_0, b$  vector

## word vector parameters

- "Dogs"  $\rightarrow [2, 3]$
- "Cows"  $\rightarrow [2, 4]$
- "eat"  $\rightarrow [-1, -1]$
- "[BEGIN]"  $\rightarrow [-10, -7]$

# Review

All cases: Training data  $\{(x^{(1)}, y^{(1)}) \dots, (x^{(n)}, y^{(n)})\}$   
Map  $x \rightarrow y$

	Regression	Binary Classification	Multi-class Classification
What is $y$ ?	$y \in \mathbb{R}$	$y \in \{-1, 1\}$	$y \in \{1, 2, 3, \dots, C\}$ <i>generating text 1 word at a time</i>
Corresponding linear model	Linear regression	logistic Regression	Softmax regression
Parameters	$w \in \mathbb{R}^d$ <i>dimension of <math>x</math></i>	$w \in \mathbb{R}^1$	$w^{(1)}, \dots, w^{(C)} \in \mathbb{R}^d$ <i><math>C \times d</math> total params</i>
Probabilistic Story	$y \sim \text{Normal}(w^T x, \sigma^2)$ <i>mean</i>	$p(y=i x) = \sigma(w^T x)$ 	$p(y=j x) = \frac{\exp(w^{(j)T} x)}{\sum_{k=1}^C \exp(w^{(k)T} x)}$ <i>Normalise to a probability distribution</i>

Loss Function  
*measures how bad a choice of parameters is*

Maximum Likelihood Estimation (MLE)  
maximize probability of data =  $\prod_{i=1}^n p(y^{(i)} | x^{(i)}; \Theta)$   
w.r.t parameters  $\Theta$

$\Leftrightarrow$  minimizing negative loglikelihood  
=  $-\sum_{i=1}^n \log p(y^{(i)} | x^{(i)}; \Theta)$

*↑ all params of model*

How to minimize loss/choose parameters

Gradient Descent OR Normal Equations

Gradient Descent

# Modeling Strategies

	Linear	Kernelized	Neural
Predictions	$\underbrace{w^T}_{\text{param}} \underbrace{x}_{\text{input}}$ <p>or</p> $w^T \underbrace{\phi(x)}_{\text{add more features}}$	$\sum_{i=1}^n \alpha_i \underbrace{k(x^{(i)}, x)}_{\substack{\text{training} \\ \text{example}}}$ <p>Rose prediction on similarity between <math>x</math> &amp; <math>x^{(i)}</math>'s</p>	$\underbrace{w^T}_{\text{param}} \underbrace{f(x)}_{\substack{\text{learned} \\ \text{non-linear function} \\ \text{of } x}}$
How to make complex prediction?	Add features manually "feature engineering"	Choose polynomial or RBF kernel, instead of dot product	Add layers + non-linearities Train layers to <u>learn</u> new features only one where we learn features from data ⇒ Needs most data
Efficiency	Linear in size of $\phi(x)$ & #data	Kernel trick! but quadratic in #data	Linear in # data & #layers SGD to take more steps on 1 loop over data
Bias/Variance	High bias many functions are not linear	low bias (universal approximators) high variance easy to overfit	
Hyper-parameters	Choice of features, $L_1/L_2$ regularization	Choice of kernel $L_2$ reg.	#layers, #neurons / layer which non-linearity sigmoid / tanh / relu Dropout, CNN: kernel size RNN: word vector size