

Cast time:
NB
80w LR

Neural Networks (INLP ch. 3)

CS 685, Spring 2021

Advanced Topics in Natural Language Processing

<http://brenocon.com/cs685>

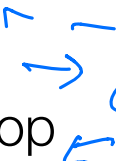

https://people.cs.umass.edu/~brenocon/cs685_s21/

Brendan O'Connor

College of Information and Computer Sciences

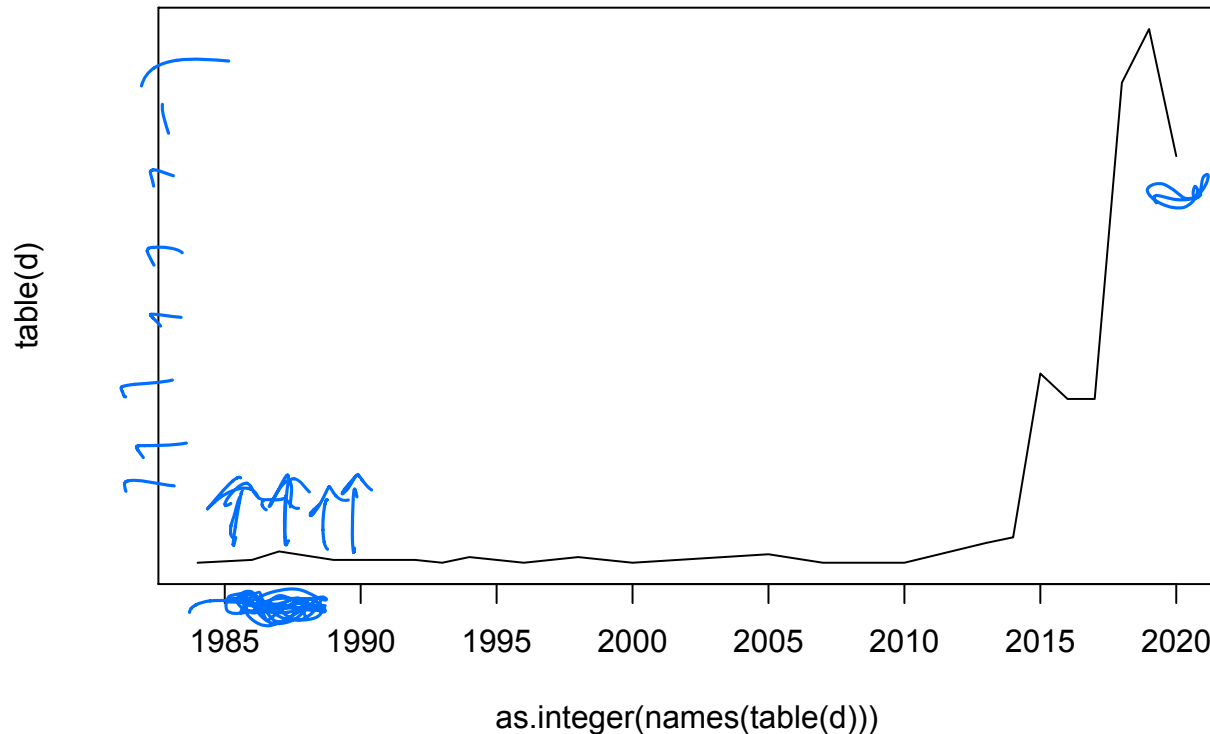
University of Massachusetts Amherst

Neural Networks in NLP

- Motivations:
 - Word sparsity => denser word representations
 - Nonlinearity ~~_____~~
 - Models
 - BoE / Deep Averaging
 - Learning
 - Backprop 
 - Dropout  Regularization
- Handwritten blue notes:* "Gradients" is written next to the "Backprop" item, and "Regularization" is written next to the "Dropout" item.

The Second Wave: NNs in NLP

- % of ACL paper titles/venues with “connectionist/connectionism”, “parallel distributed”, “neural network”, or “deep learning”
- <https://www.aclweb.org/anthology/>



NN Text Classification

- Goals:

- Avoid feature engineering
- Generalize beyond individual words

- General model architectures that work well for many different datasets (and tasks!)

- For medium-to-large labeled training datasets, deep learning methods generally outperform feature-based LogReg

Δ normalization
use n-grams
use lexicons

(100s ??)

Word sparsity



Row vector
 $\rightarrow \langle 0, 2, 0, 0, 1 \rangle$
 ↑ "cat" ↑ "man"

- Alternate view of Bag-of-Words classifiers: every word has a “one-hot” representation.

- Represent each word as a vector of zeros with a single 1 identifying the index of the word length ✓

- Doc BOW \mathbf{x} = average of all words' vectors

vocabulary

i
hate
love
the
movie
film

movie = $\langle 0, 0, 0, 0, 1, 0 \rangle$

film = $\langle 0, 0, 0, 0, 0, 1 \rangle$

vec add.

what are the issues of representing a word this way?

- Scaling
 (if dense repr)

- OOVs

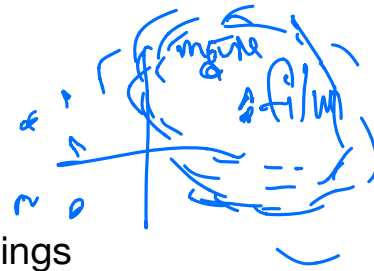
- Compo/ctx



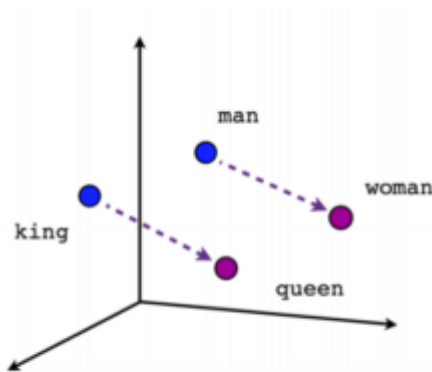
all word vcs
are orthogonal

Learned!

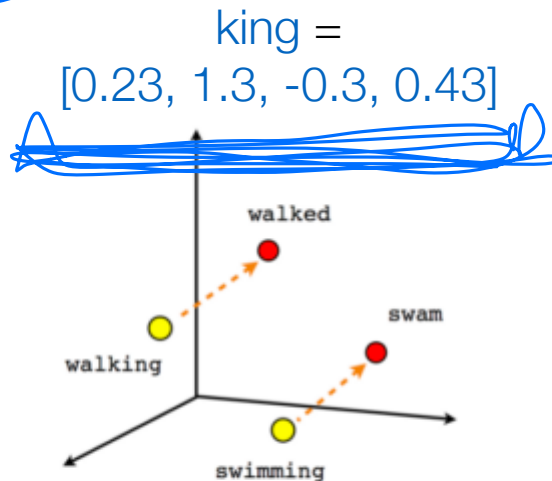
Word embeddings



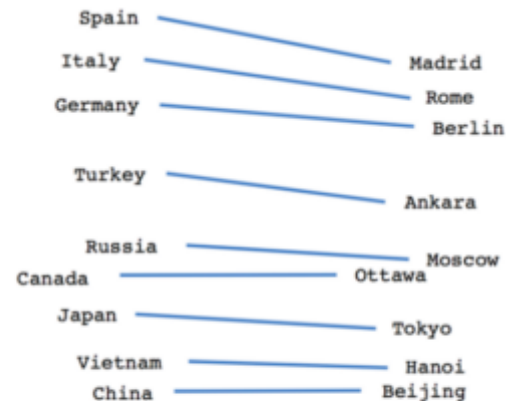
- Represent words with low(ish)-dimensional vectors called **embeddings**
- Today: word embeddings are the first “lookup” layer in an NN. Every word in vocabulary has a vector — these are model parameters.
- Ideally: semantically similar words get similar vectors. Or other semantic properties??



Male-Female



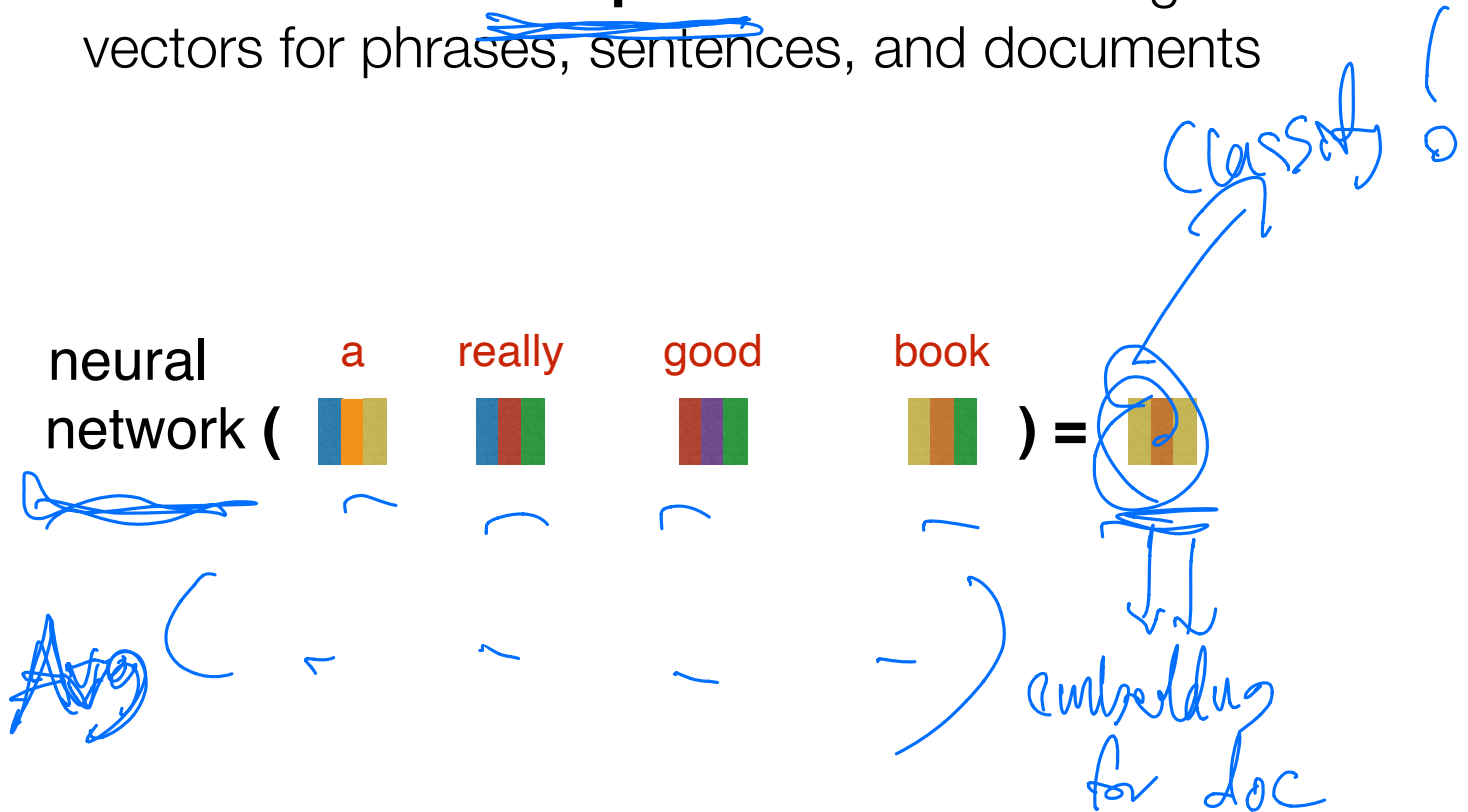
Verb tense



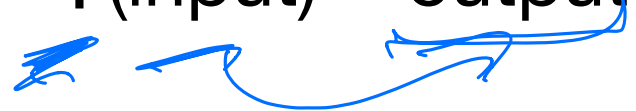
Country-Capital

composing embeddings

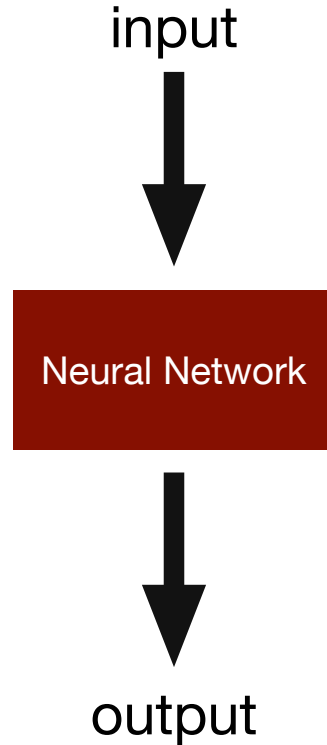
- neural networks **compose** word embeddings into vectors for phrases, sentences, and documents



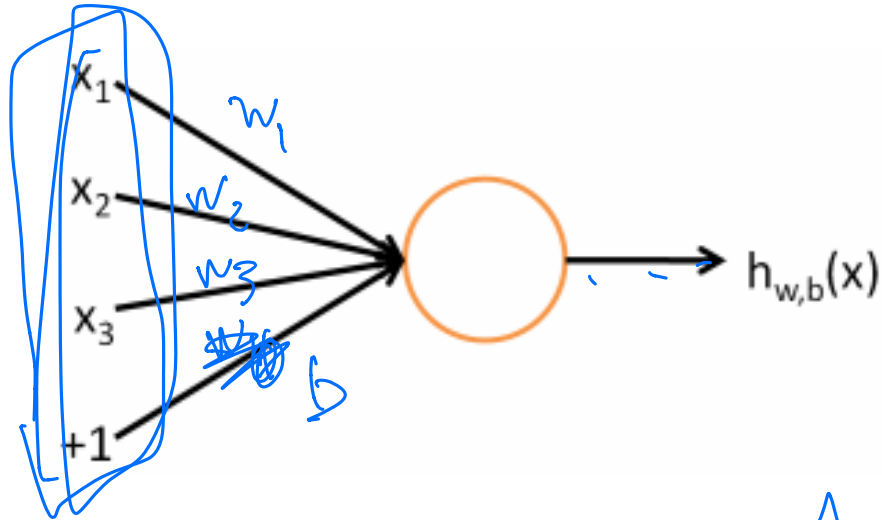
what is deep learning? *ANNs*

$$f(\text{input}) = \text{output}$$


what is deep learning?



Logistic Regression by Another Name: Map inputs to output



Input

$$x_1 \dots x_d = \vec{x}$$

Output

$$f\left(\underbrace{\sum_i x_i w_i + b}_{\in \mathbb{R}}\right)$$

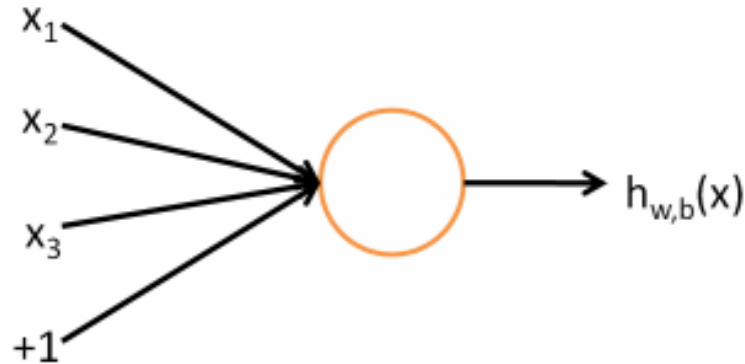
Activ

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

logrs.
sigmoid
 $\in (0, 1)$

→ (0.1)

Logistic Regression by Another Name: Map inputs to output



Input

Vector $x_1 \dots x_d$

Output

$$f\left(\sum_i w_i x_i + b\right)$$

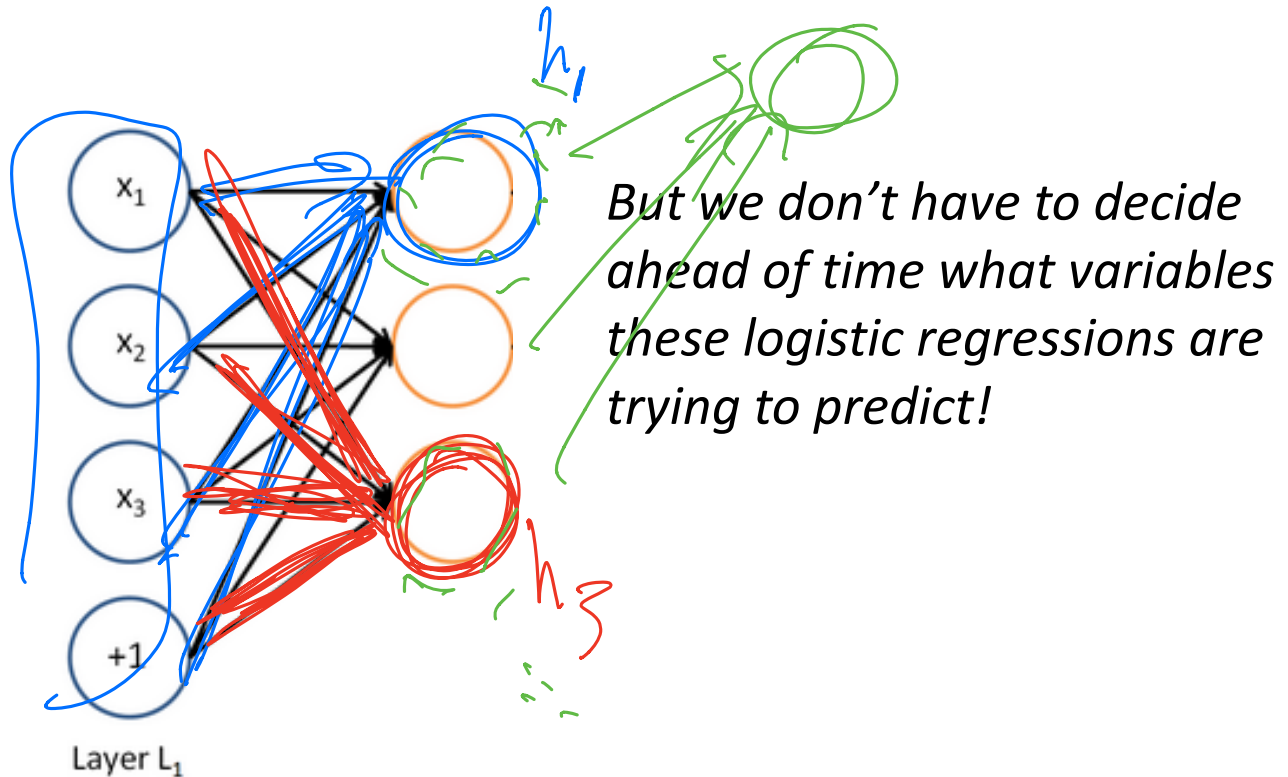
Activation

$$f(z) \equiv \frac{1}{1 + \exp(-z)}$$

pass through
nonlinear sigmoid

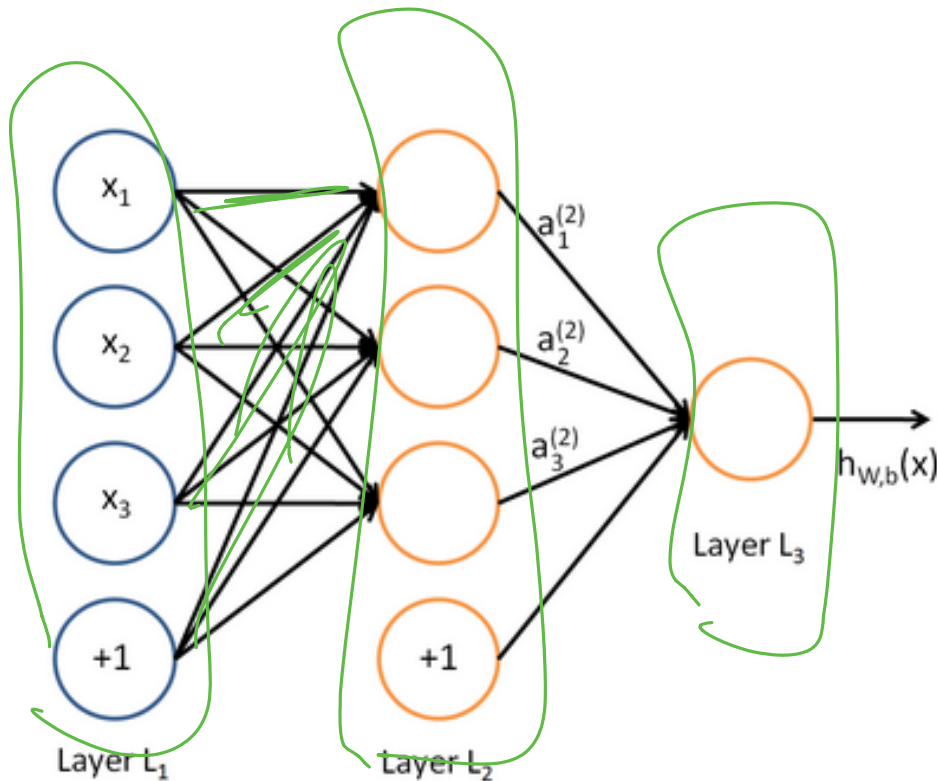
NN: kind of like several intermediate logregs

If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs ...



NN: kind of like several intermediate logregs

... which we can feed into another logistic regression function

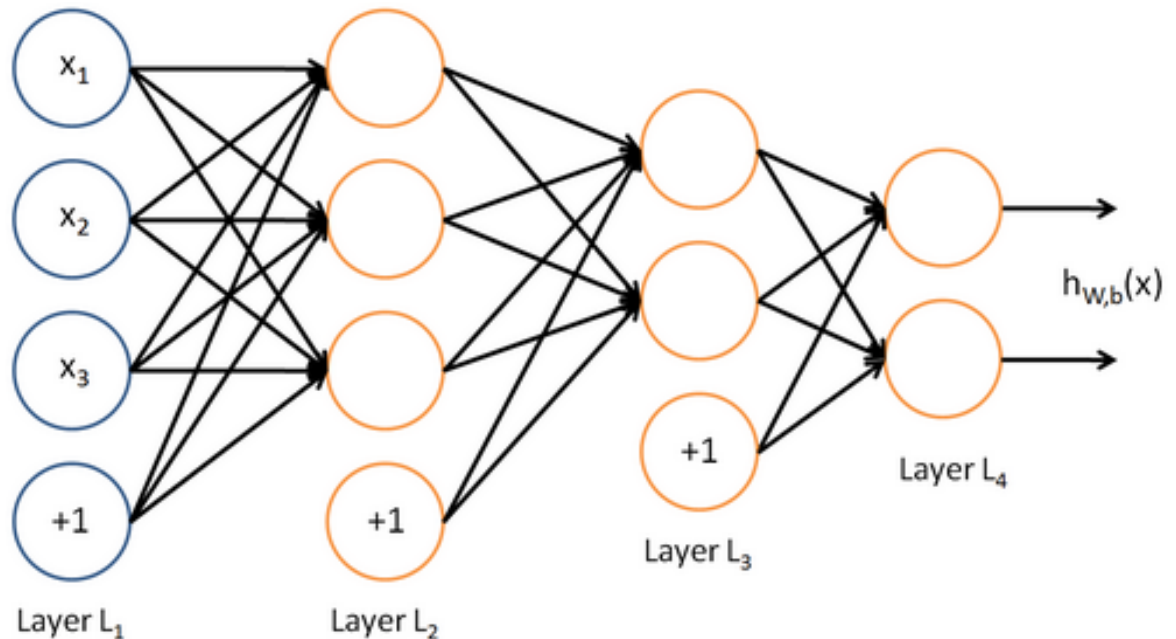


It is the loss function that will direct what the intermediate hidden variables should be, so as to do a good job at predicting the targets for the next layer, etc.

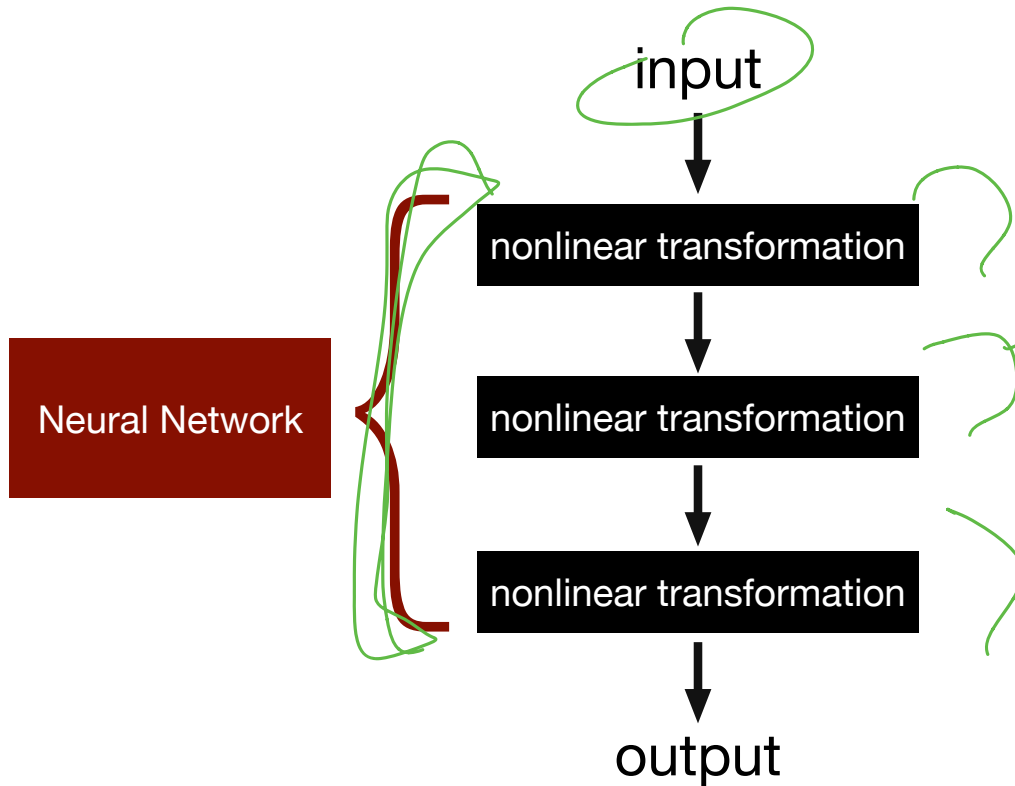
NN: kind of like several intermediate logregs

Before we know it, we have a multilayer neural network....

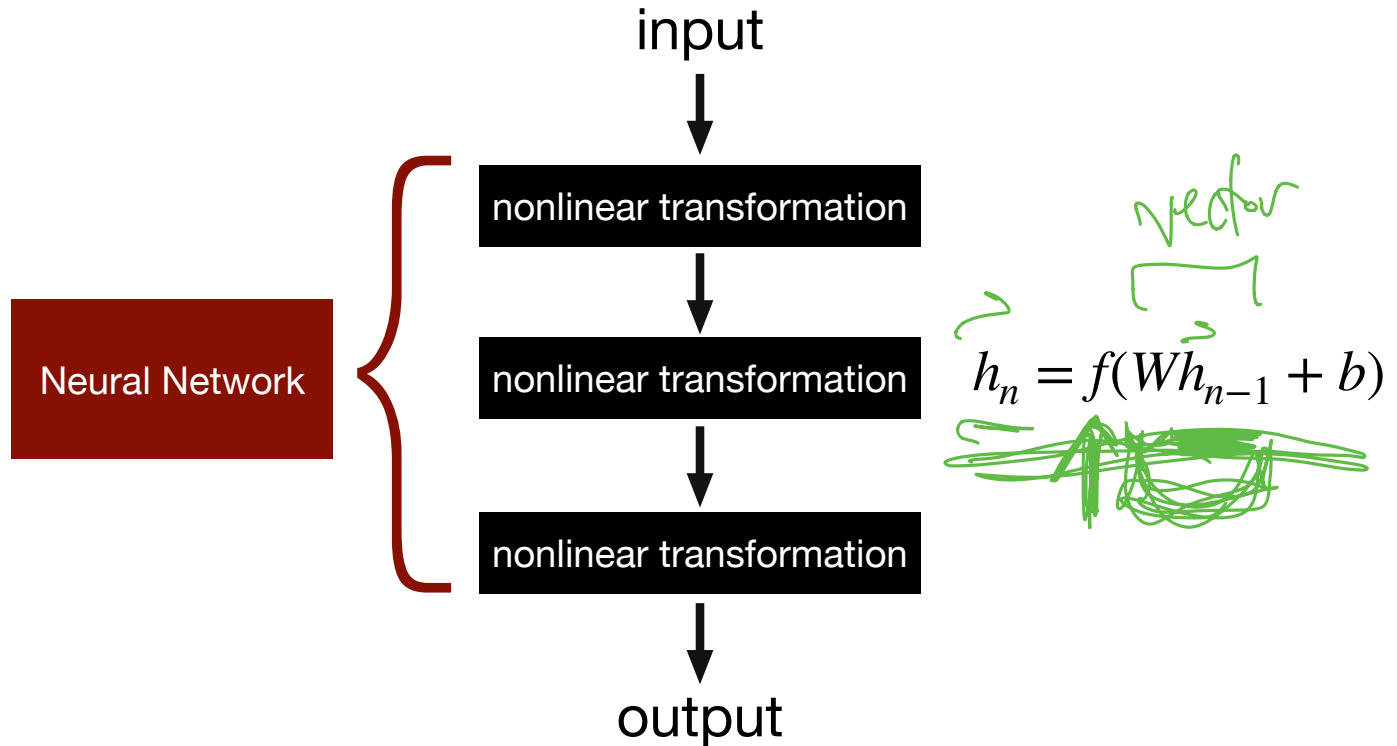
a.k.a. **feedforward network** (see INLP on terminology)



what is deep learning?



what is deep learning?



Nonlinear activations

- “Squash functions”!

- Logistic / Sigmoid

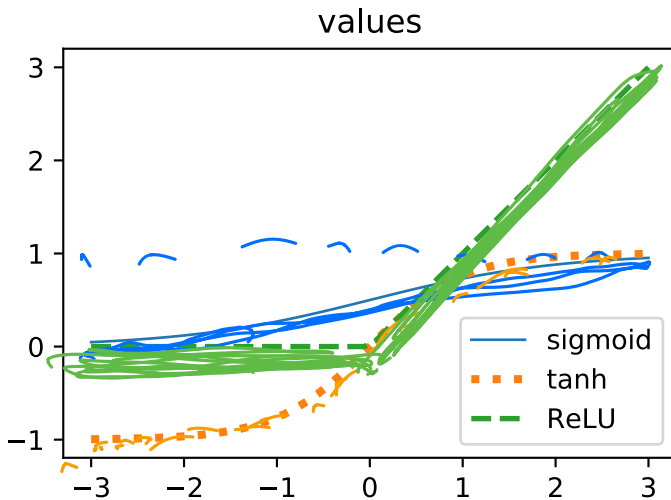
$$f(x) = \frac{1}{1 + e^{-x}} \quad (1)$$

- tanh

$$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1 \quad (2)$$

- ReLU

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases} \quad (3)$$



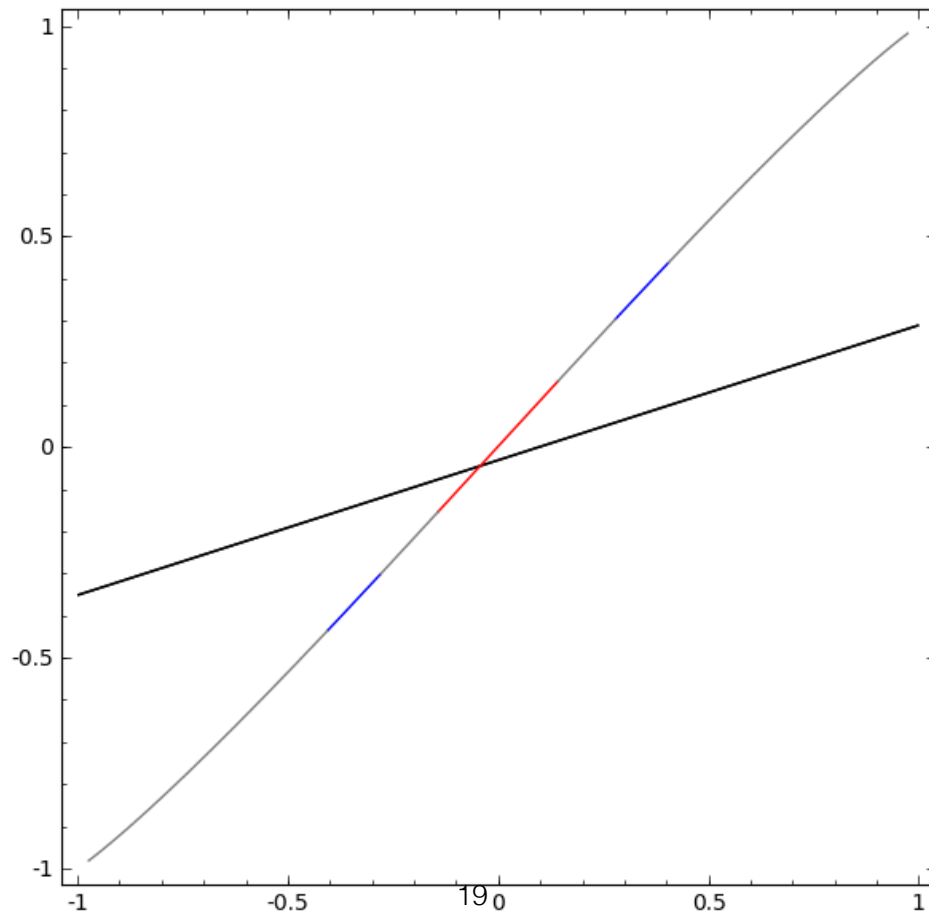
“Softplus”

is a multi-layer neural network with no nonlinearities
(i.e., f is the identity $f(x) = x$)
more powerful than a one-layer network?

$$f(W_3 f(W_2 f(W_1 x)))$$

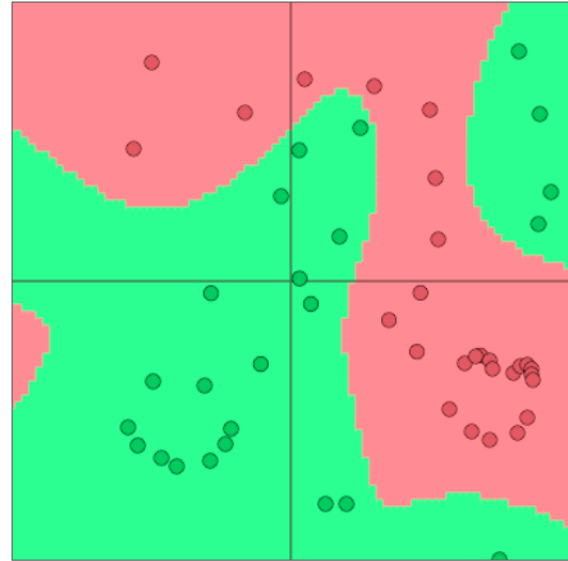
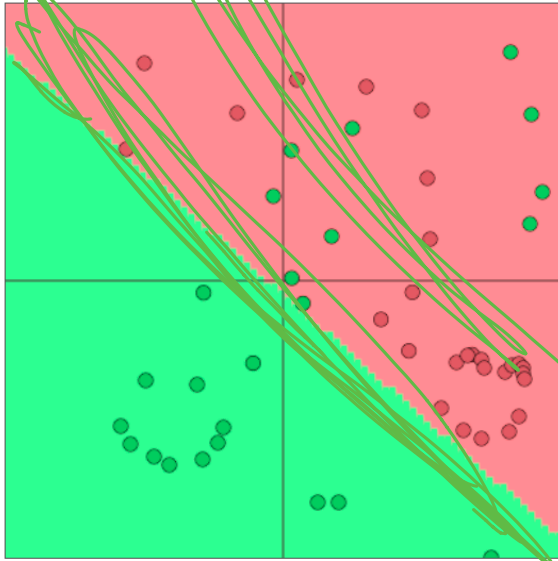
$$\equiv \underbrace{W_3 W_2 W_1}_W x$$

why nonlinearities?

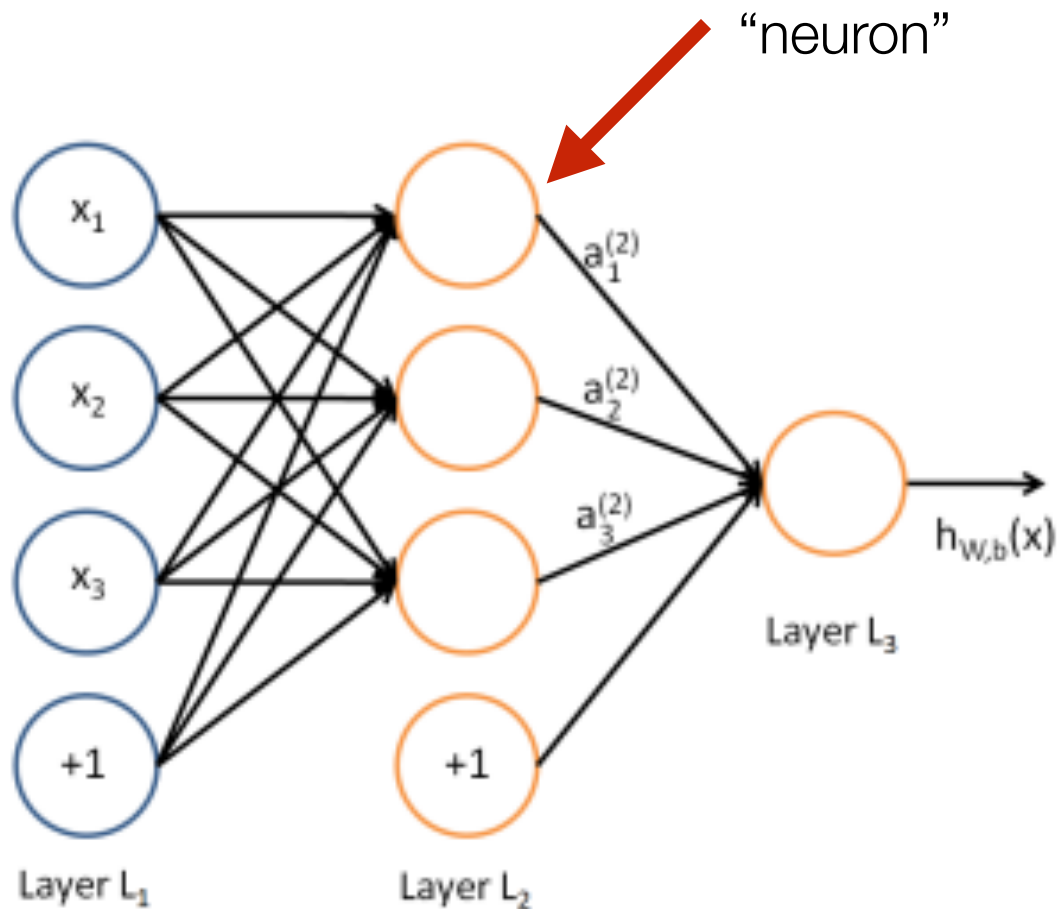


credit for figure:
Christopher Olah

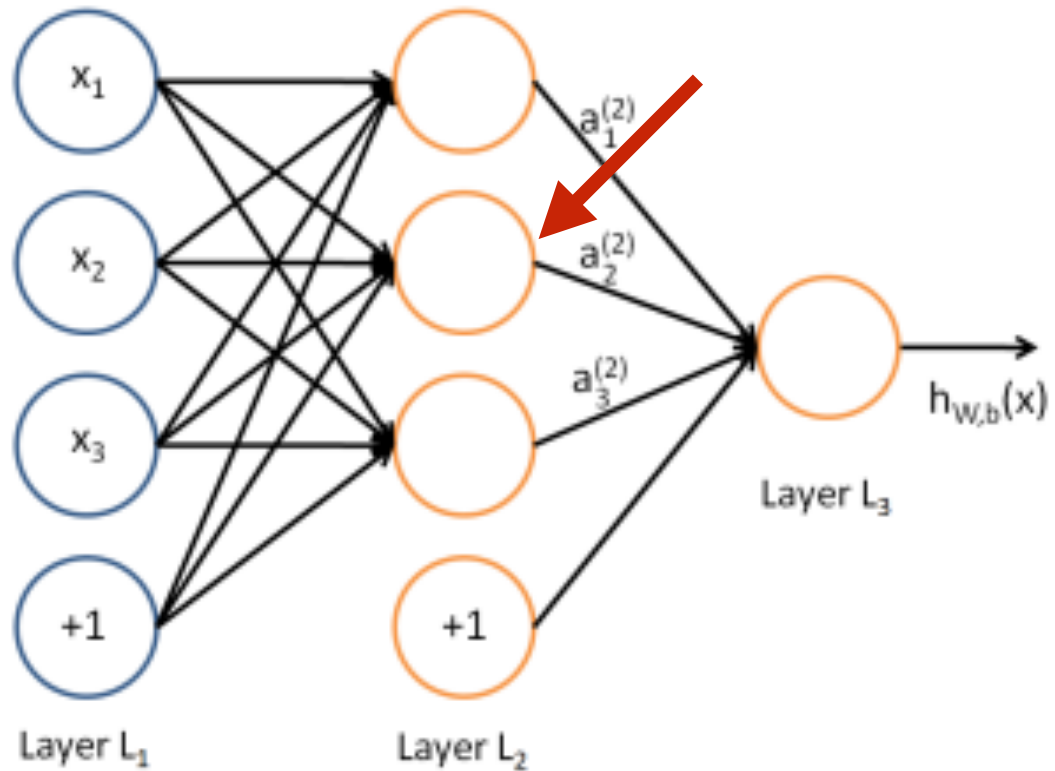
why nonlinearities?



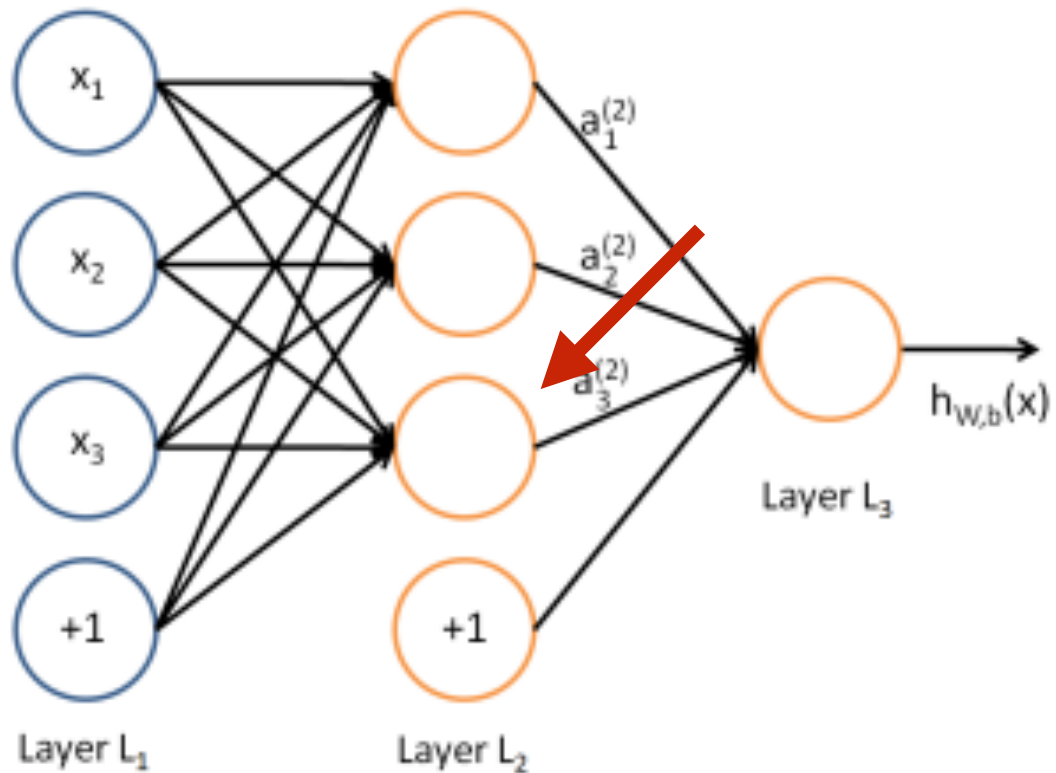
Tensorflow Playground



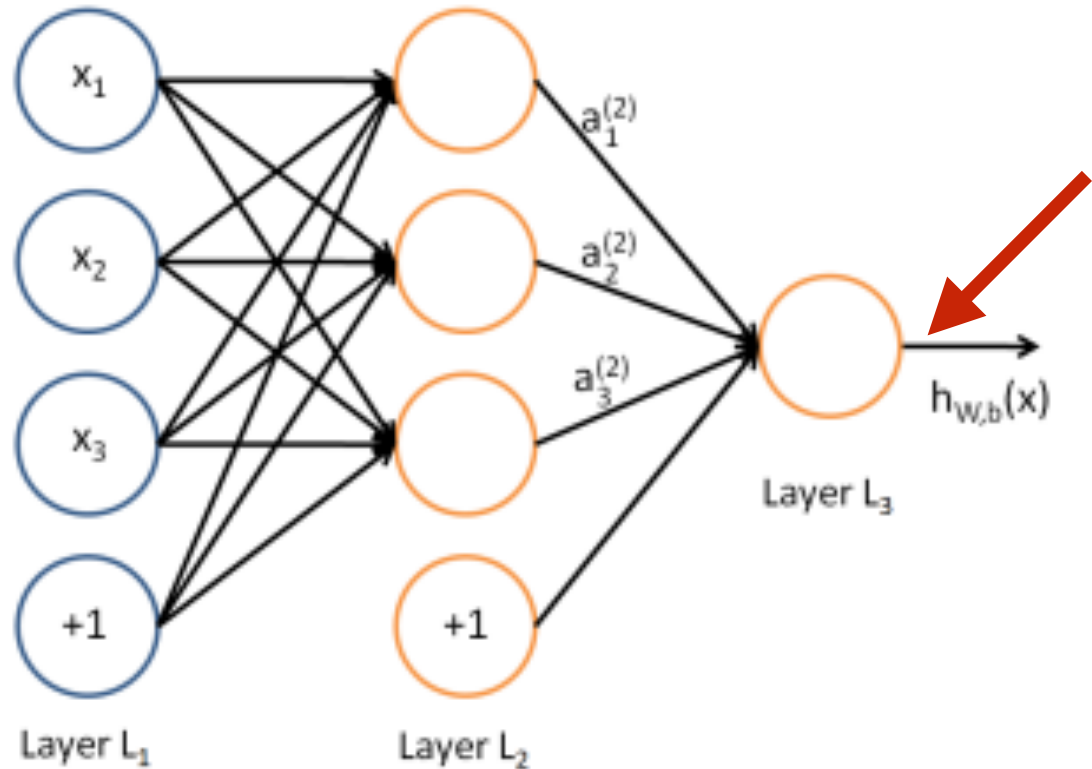
$$a_1^{(2)} = f\left(W_{11}^{(1)}x_1 + W_{12}^{(1)}x_2 + W_{13}^{(1)}x_3 + b_1^{(1)}\right)$$



$$a_2^{(2)} = f\left(W_{21}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{23}^{(1)}x_3 + b_2^{(1)}\right)$$

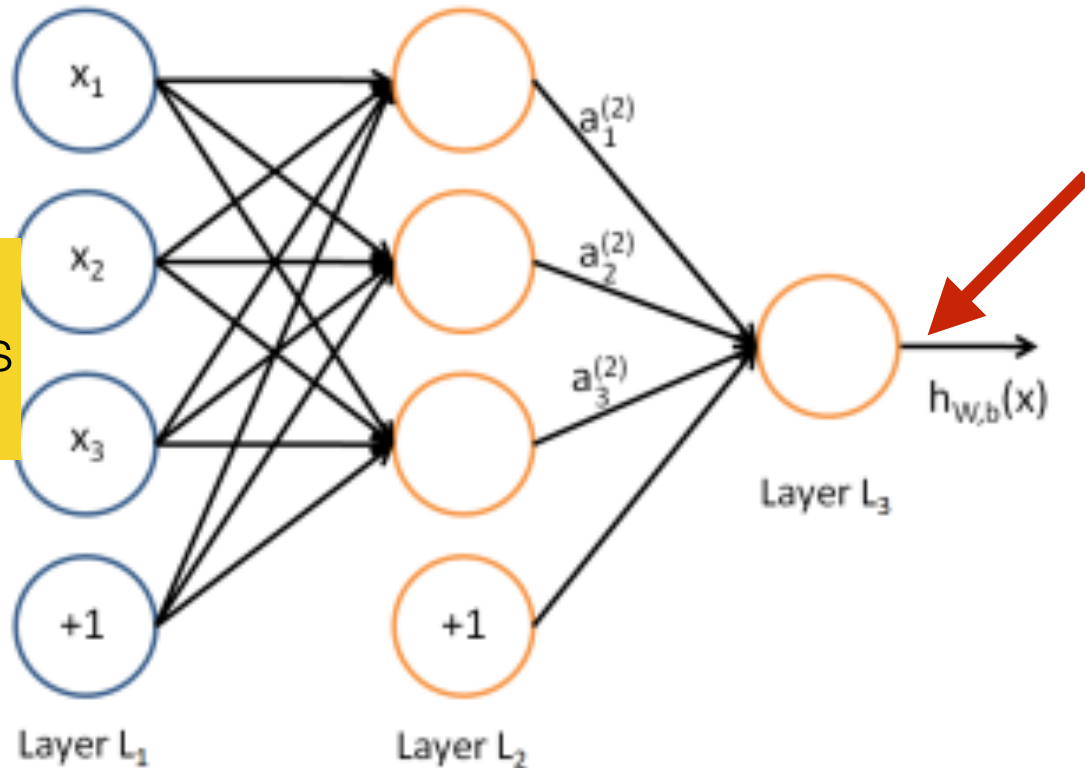


$$a_3^{(2)} = f\left(W_{31}^{(1)}x_1 + W_{32}^{(1)}x_2 + W_{33}^{(1)}x_3 + b_3^{(1)}\right)$$



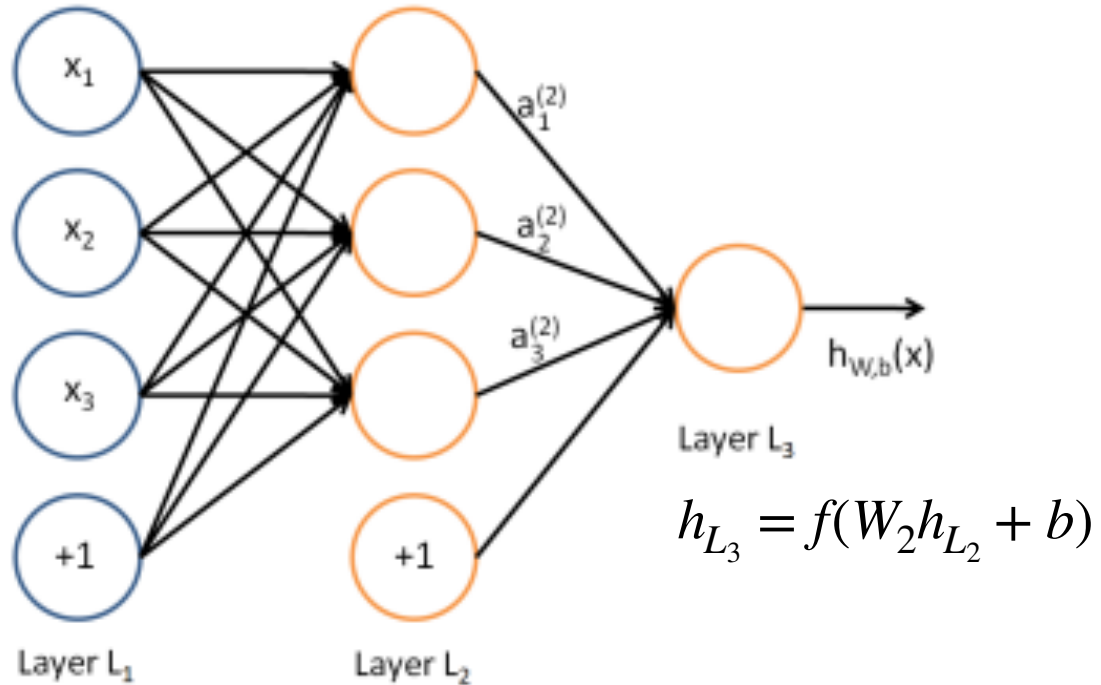
$$h_{W,b}(x) = a_1^{(3)} = f\left(W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + W_{13}^{(2)} a_3^{(2)} + b_1^{(2)}\right)$$

we will be
learning the x 's
and the W 's!



$$h_{W,b}(x) = a_1^{(3)} = f\left(W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + W_{13}^{(2)} a_3^{(2)} + b_1^{(2)}\right)$$

in matrix-vector notation...



$$h_{L_2} = f(W_1 x + b)$$

Dracula is a really good book!



neural
network



Positive

softmax function

- let's say I have 3 classes (e.g., positive, neutral, negative)
- use multiclass logreg with “cross product” features between input vector \mathbf{x} and 3 output classes. for every class c , i have an associated weight vector β_c , then

$$P(y = c | \mathbf{x}) = \frac{e^{\beta_c \mathbf{x}}}{\sum_{k=1}^3 e^{\beta_k \mathbf{x}}}$$

softmax function

output is vector $\in \mathbb{R}^d$ but non-neg
sum to 1

$$\text{softmax}(x) = \frac{e^x}{\sum_j e^{x_j}}$$

x is a vector

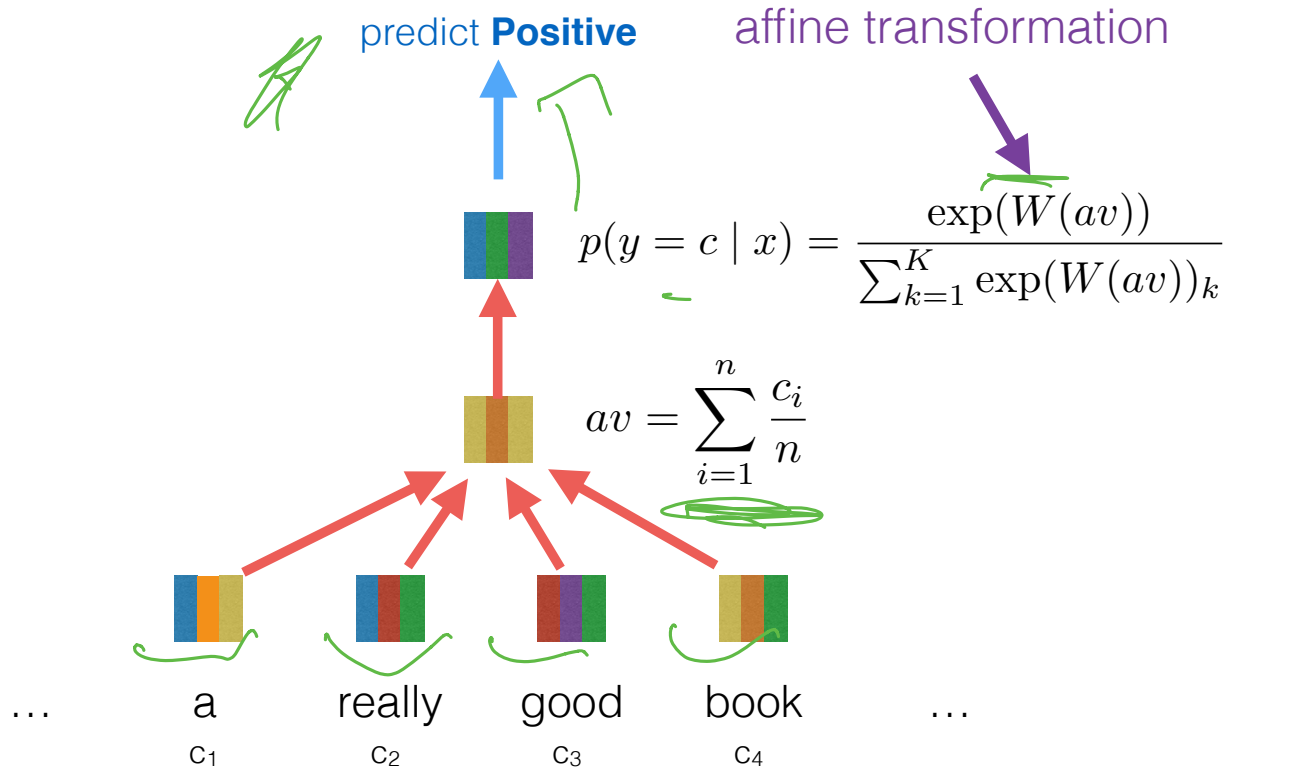
$$x \in \mathbb{R}^d$$

x_j is dimension j of x

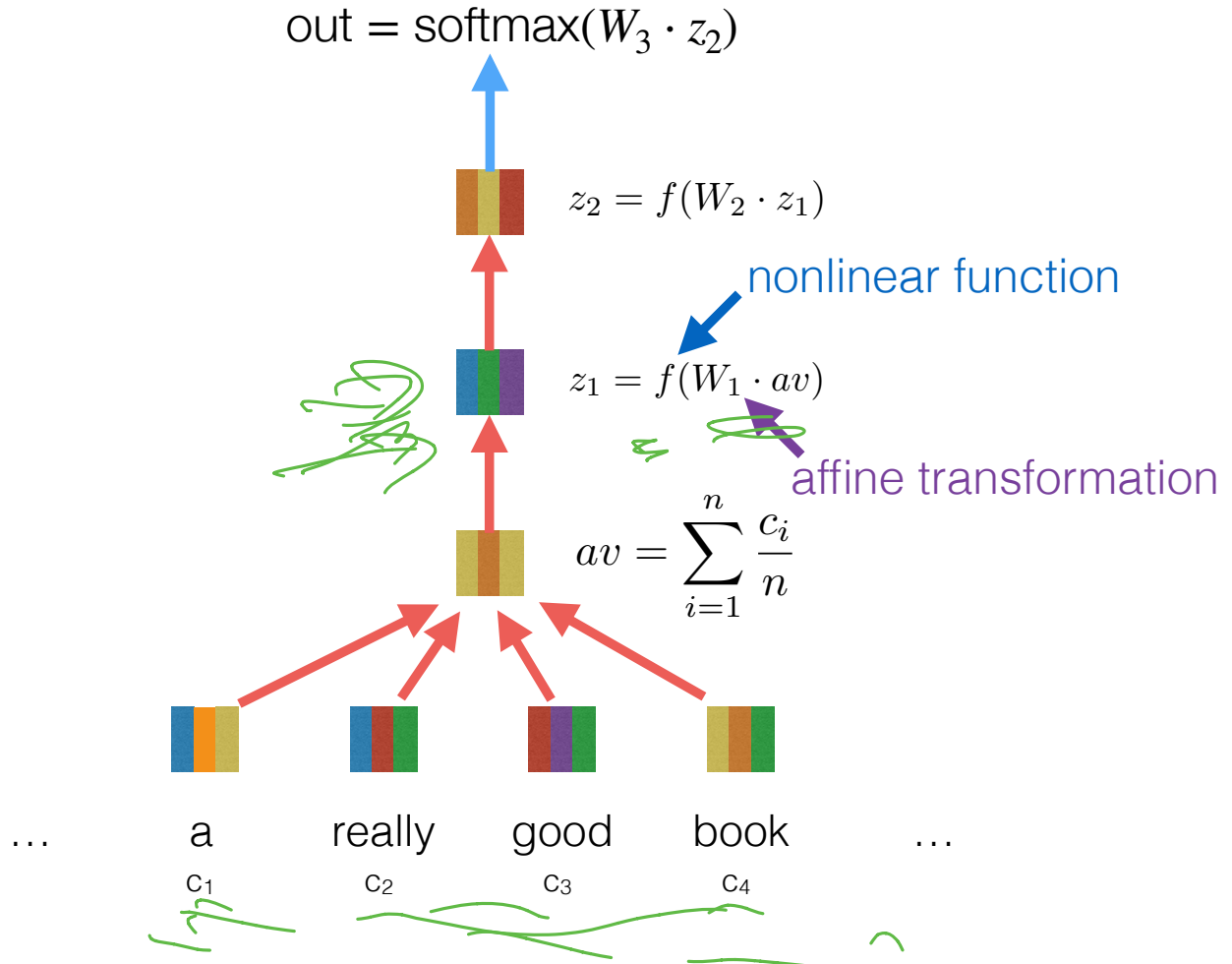
each dimension j of the softmaxed output
represents the probability of class j

$$[\text{softmax}(x)]_j = P(y=j/x)$$

“bag of embeddings”

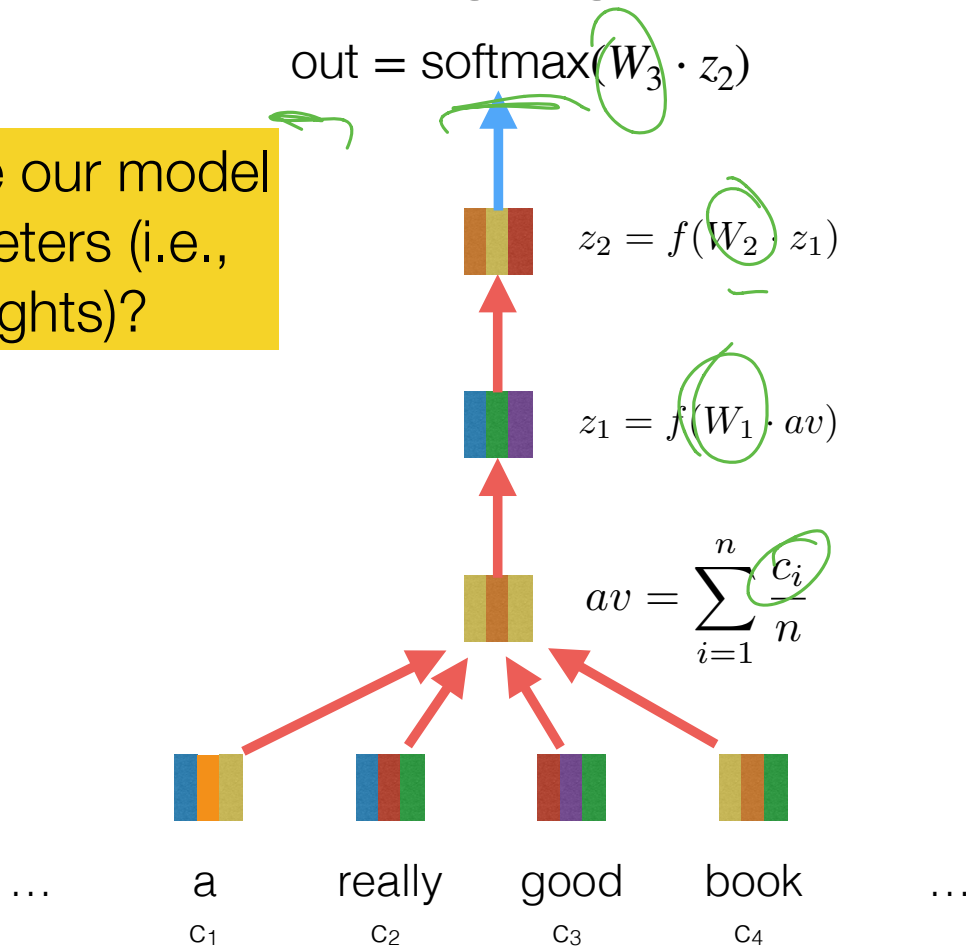


deep averaging networks

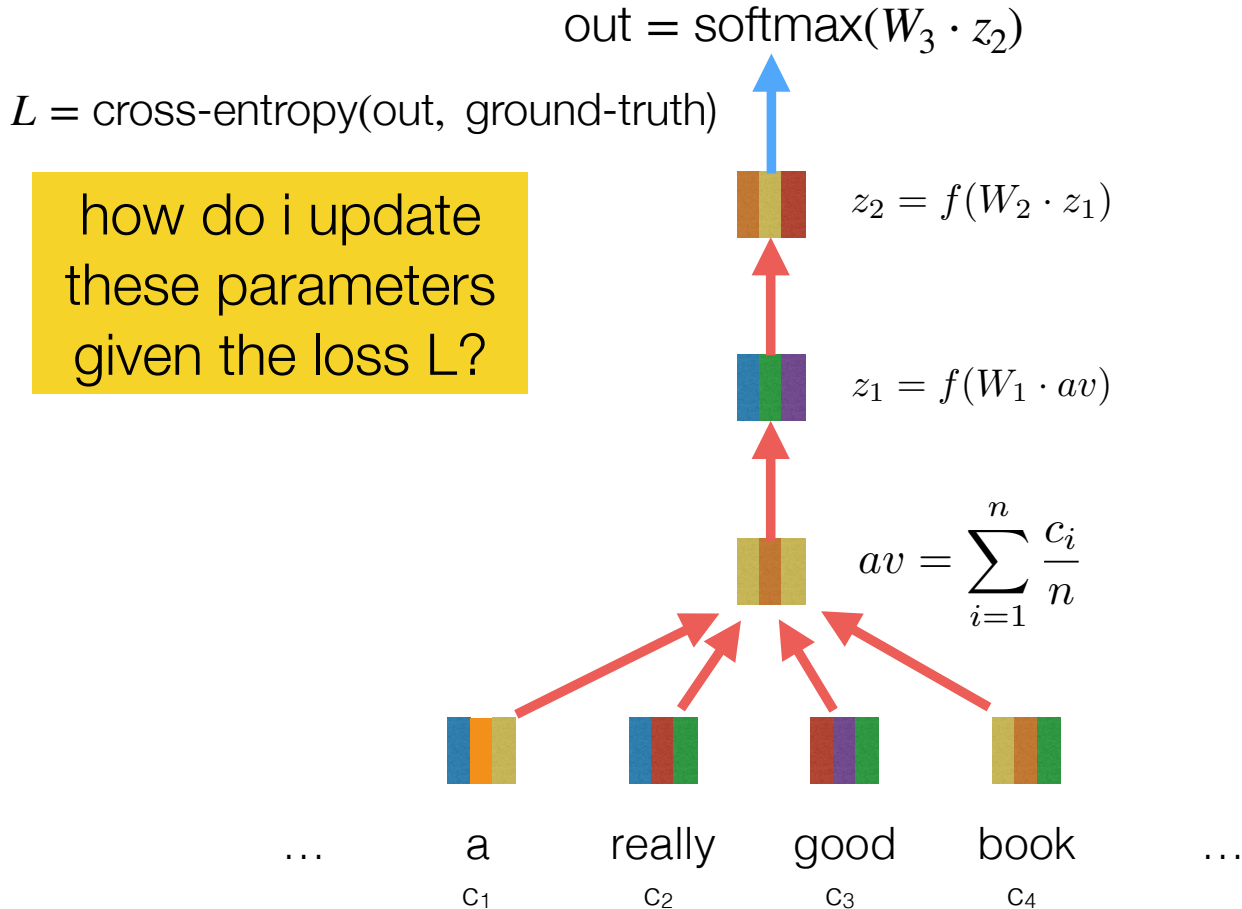


deep averaging networks

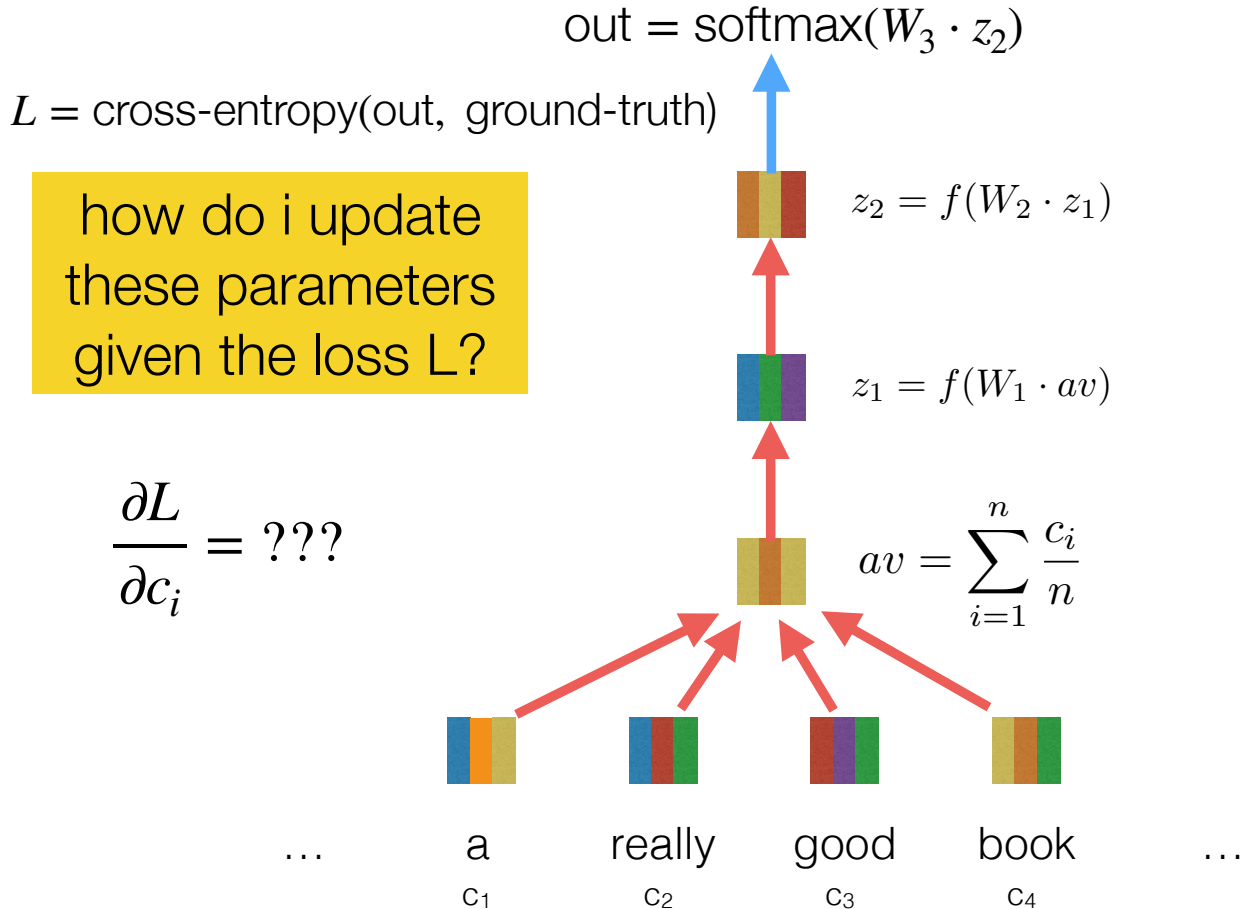
what are our model parameters (i.e., weights)?



deep averaging networks



deep averaging networks



deep averaging networks

chain rule!!!

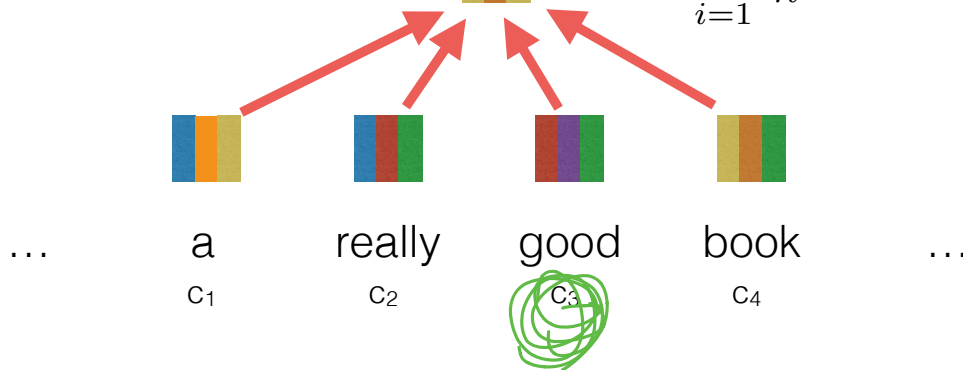
$$\frac{\partial L}{\partial c_i} = \frac{\partial L}{\partial \text{out}} \frac{\partial \text{out}}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial av} \frac{\partial av}{\partial c_i}$$

$$\text{out} = \text{softmax}(W_3 \cdot z_2)$$

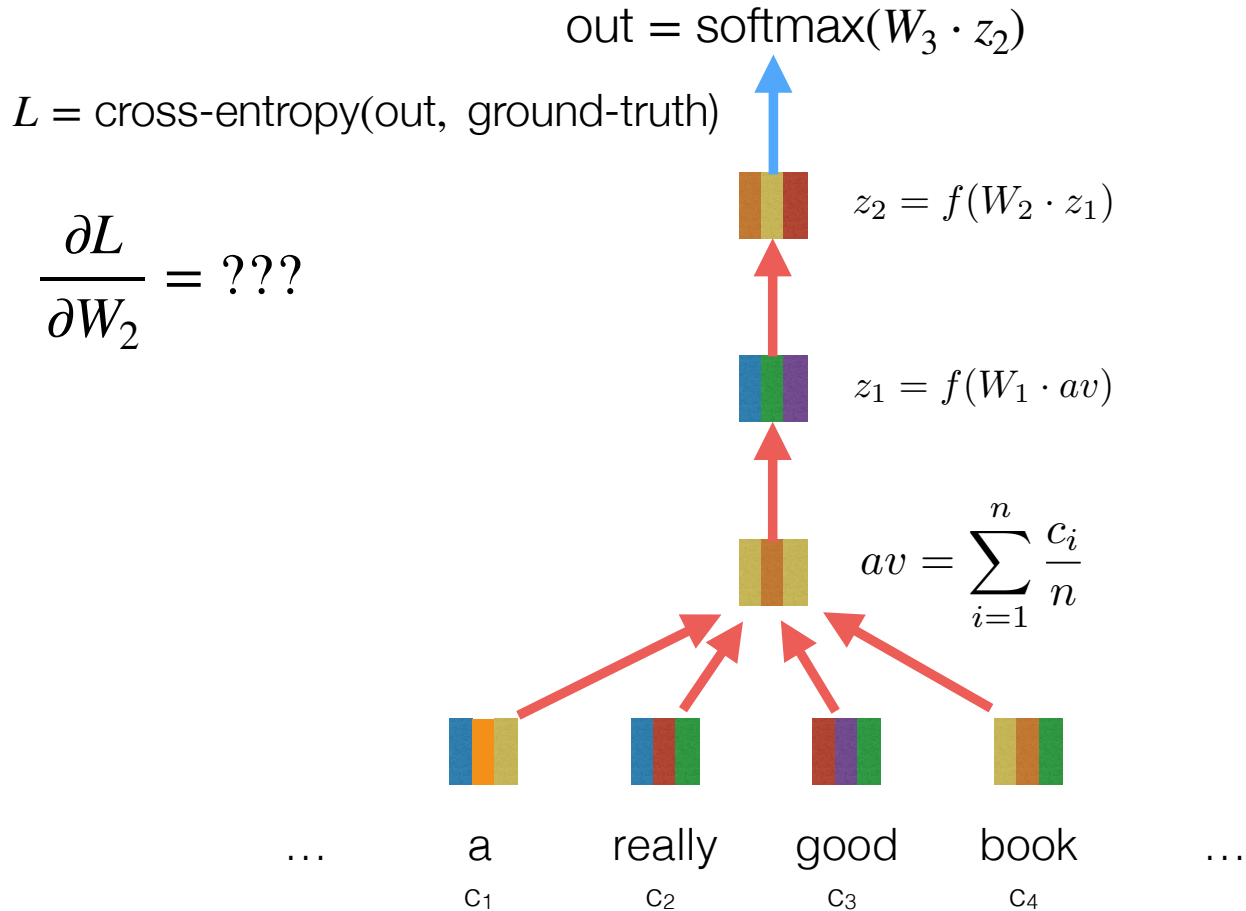
$$z_2 = f(W_2 \cdot z_1)$$

$$z_1 = f(W_1 \cdot av)$$

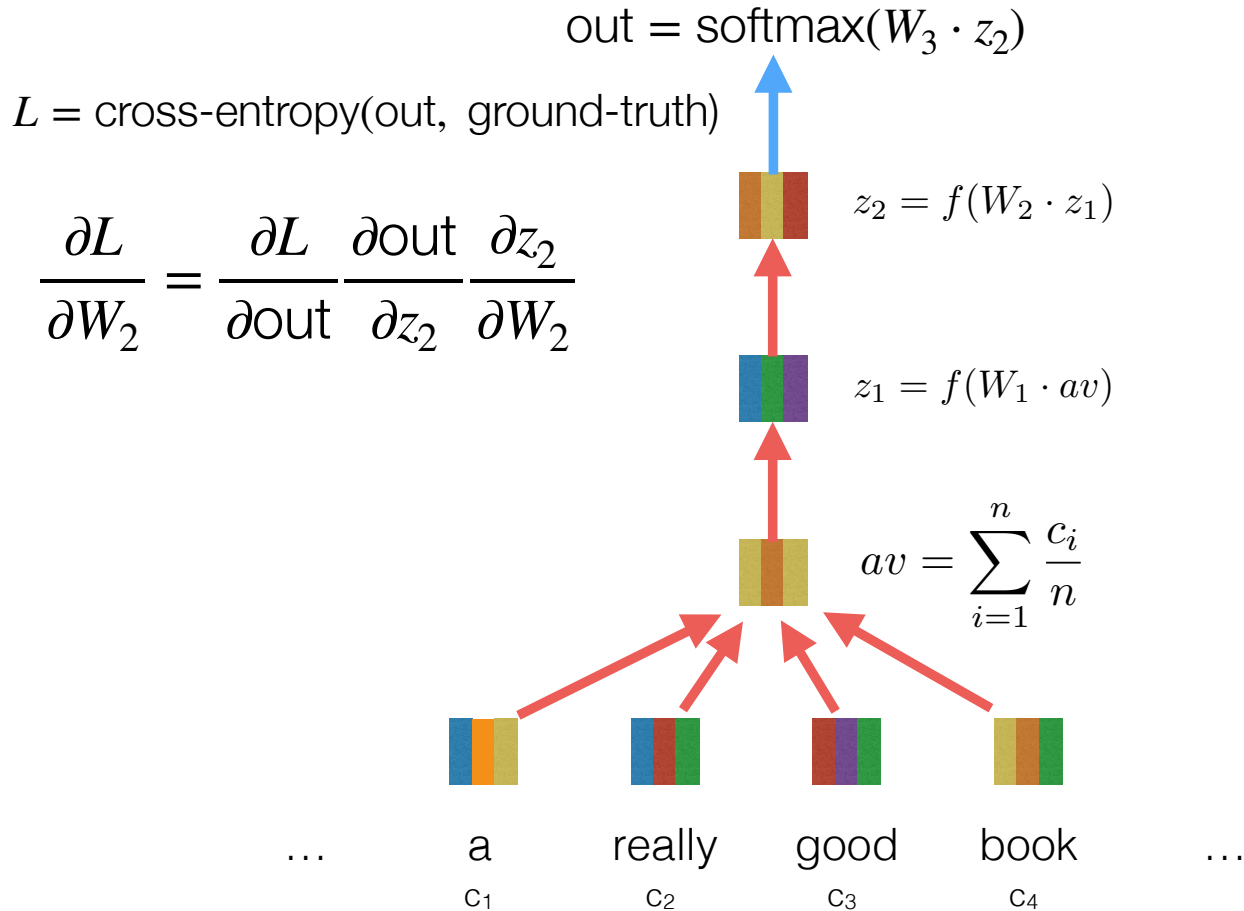
$$av = \sum_{i=1}^n \frac{c_i}{n}$$



deep averaging networks



deep averaging networks



backpropagation


- use the chain rule to compute partial derivatives w/ respect to each parameter
- trick: re-use derivatives computed for higher layers to compute derivatives for lower layers!


$$\frac{\partial L}{\partial c_i} = \frac{\partial L}{\partial \text{out}} \frac{\partial \text{out}}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial a v} \frac{\partial a v}{\partial c_i}$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial \text{out}} \frac{\partial \text{out}}{\partial z_2} \frac{\partial z_2}{\partial W_2}$$

deep learning frameworks make building NNs super easy!

$$\text{out} = \text{softmax}(W_2 \cdot z_2)$$


$$z_1 = f(W_1 \cdot av)$$


$$av = \sum_{i=1}^n \frac{c_i}{n}$$



a really good book

set up the network

```
def __init__(self, n_classes, vocab_size, emb_dim=300,
               n_hidden_units=300):
    super(DanModel, self).__init__()
    self.n_classes = n_classes
    self.vocab_size = vocab_size
    self.emb_dim = emb_dim
    self.n_hidden_units = n_hidden_units
    self.embeddings = nn.Embedding(self.vocab_size,
                                    self.emb_dim)
    self.classifier = nn.Sequential(
        nn.Linear(self.n_hidden_units,
                  self.n_hidden_units),
        nn.ReLU(),
        nn.Linear(self.n_hidden_units,
                  self.n_classes))
    self._softmax = nn.Softmax()
```

deep learning frameworks make building NNs super easy!

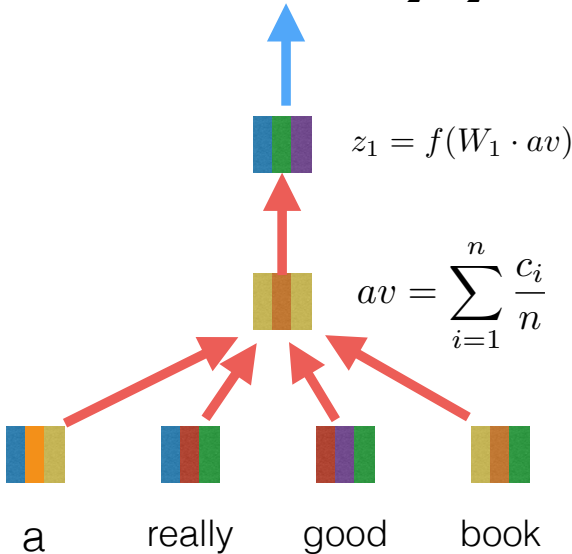
PyTorch

do a forward pass to compute prediction

$$\text{out} = \text{softmax}(W_2 \cdot z_2)$$

$$z_1 = f(W_1 \cdot av)$$

$$av = \sum_{i=1}^n \frac{c_i}{n}$$

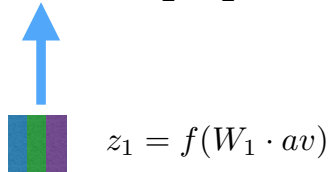


```
def forward(self, batch, probs=False):
    text = batch['text']['tokens']
    length = batch['length']
    text_embed = self._word_embeddings(text)
    # Take the mean embedding. Since padding results
    # in zeros its safe to sum and divide by length
    encoded = text_embed.sum(1)
    encoded /= lengths.view(text_embed.size(0), -1)

    # Compute the network score predictions
    logits = self.classifier(encoded)
    if probs:
        return self._softmax(logits)
    else:
        return logits
```

deep learning frameworks make building NNs super easy!

out = softmax($W_2 \cdot z_2$)



$$av = \sum_{i=1}^n \frac{c_i}{n}$$

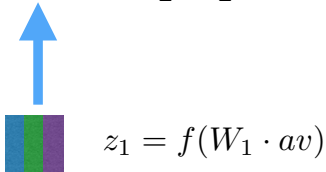


do a backward pass to update weights

```
def _run_epoch(self, batch_iter, train=True):  
    self._model.train()  
    for batch in batch_iter:  
        model.zero_grad()  
        out = model(batches)  
        batch_loss = criterion(out,  
                                batch['label'])  
        batch_loss.backward()  
        self.optimizer.step()
```

deep learning frameworks make building NNs super easy!

out = softmax($W_2 \cdot z_2$)



$$av = \sum_{i=1}^n \frac{c_i}{n}$$



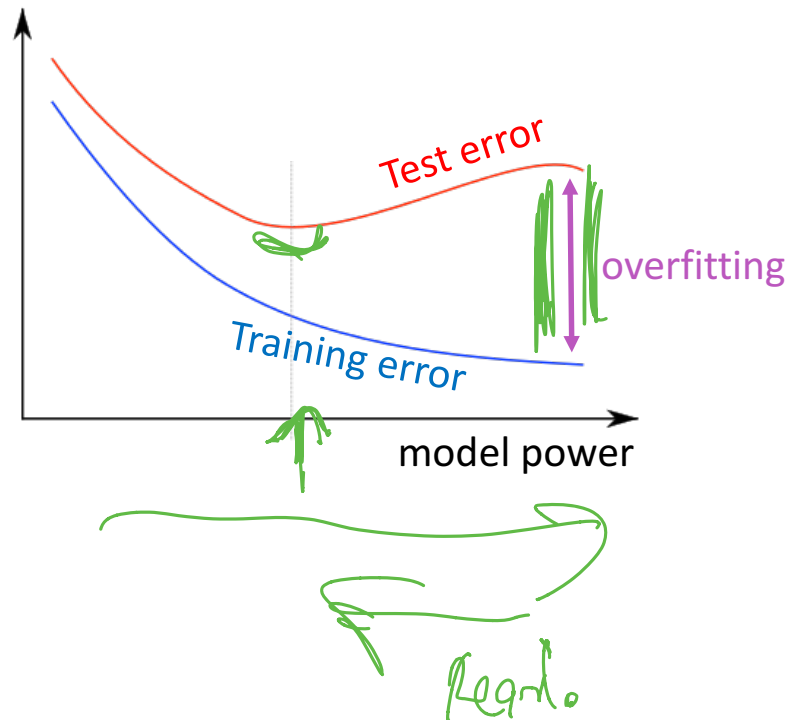
do a backward pass to update weights

```
def _run_epoch(self, batch_iter, train=True):  
    self._model.train()  
    for batch in batch_iter:  
        model.zero_grad()  
        out = model(batches)  
        batch_loss = criterion(out,  
                                batch['label'])  
        batch_loss.backward()  
        self.optimizer.step()
```

that's it! no need to compute gradients by hand!

Regularization

- Regularization prevents **overfitting** when we have a lot of features (or later a very powerful/deep model,++)



L2 regularization

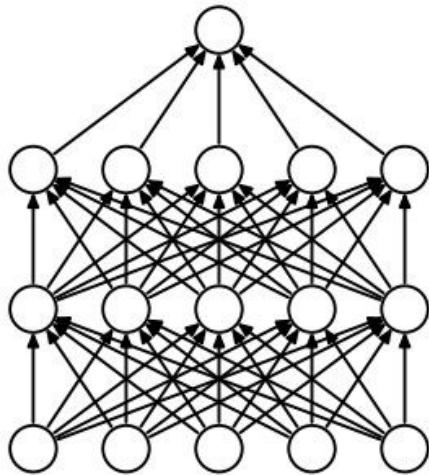
$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}} \right) + \lambda \sum_k \theta_k^2$$

θ represents all of the model's parameters!

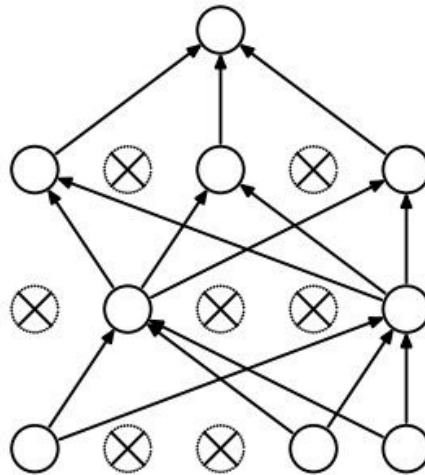
penalizing their norm leads to smaller weights
we are constraining the parameter space
we are putting a prior on our model

Dropout for NNs

randomly set $p\%$ of neurons to 0 in the forward pass



(a) Standard Neural Net

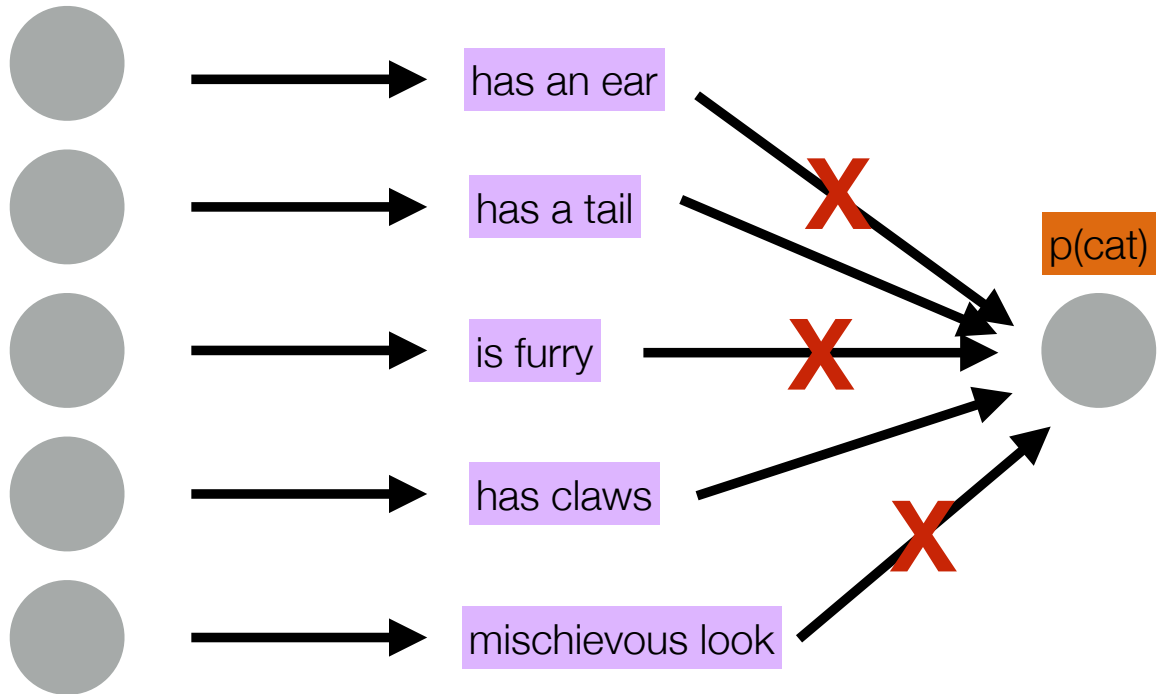


(b) After applying dropout.

[Srivastava et al., 2014]

Why?

randomly set $p\%$ of neurons to 0 in the forward pass



Addressing instability

- Training can be unstable! Therefore some tricks.
 - Initialization — random small but reasonable values can help.
 - Layer normalization (very important for some recent architectures)
- Since performance variance is high, you need to evaluate *multiple runs*
 - whether you're averaging or taking max performance
 - esp for comparisons!

- A few unresolved questions about NNs in NLP
 - Useful architectures?
 - Many: Convolutional, Recurrent, Self/cross-attention
 - Modular systems?
 - Interpretability / explainability?
 - Incorporate prior knowledge?
 - Transferring information across datasets/
languages/etc?
- These are major questions for NLP modeling right now!

o re bake go o

o o re ^{over} baking go o

baked