MLPs, Representations, Ngram Language Models

Yoav Goldberg

Binary:

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$

$$f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$

Binary:

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$
 $\operatorname{sign}(\mathbf{w} \cdot \mathbf{x} + b)$

$$f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$



$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$

$$sign(\mathbf{w} \cdot \mathbf{x} + b) \\ \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

$$f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$

Binary: $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ $sign(\mathbf{w} \cdot \mathbf{x} + b)$ $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

$$f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$
$$\arg \max_{i} (\mathbf{W} \cdot \mathbf{x} + \mathbf{b})_{[i]}$$

Binary: $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ $sign(\mathbf{w} \cdot \mathbf{x} + b)$ $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

$$f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$

arg max_i $(\mathbf{W} \cdot \mathbf{x} + \mathbf{b})_{[i]}$
arg max softmax $(\mathbf{W} \cdot \mathbf{x} + \cdot b)_{[i]}$

Binary: $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ $sign(\mathbf{w} \cdot \mathbf{x} + b)$ $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

Multi class:

$$f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$

arg max_i $(\mathbf{W} \cdot \mathbf{x} + \mathbf{b})_{[i]}$
ded? arg max softmax $(\mathbf{W} \cdot \mathbf{x} + \mathbf{b})_{[i]}$

is this really needed?

xor(0, 0) = 0xor(1, 0) = 1xor(0, 1) = 1xor(1, 1) = 0

 $(0, 0) \cdot \mathbf{w} + b < 0$ $(0, 1) \cdot \mathbf{w} + b \ge 0$ $(1, 0) \cdot \mathbf{w} + b \ge 0$ $(1, 1) \cdot \mathbf{w} + b < 0$

$$\mathbf{W} = ?$$



 $\phi(x_1, x_2) = [x_1 \times x_2, x_1 + x_2]_{:}$

 $\phi(x_1, x_2) = [x_1 \times x_2, x_1 + x_2]_{:}$



 $\phi(x_1, x_2) = [x_1 \times x_2, x_1 + x_2]_{:}$



 $\phi(x_1, x_2) = [x_1 \times x_2, x_1 + x_2]_{:}$



 $f(\mathbf{x}) = \mathbf{w} \cdot \phi(\mathbf{x}) + b$

 $\phi(x_1, x_2) = [x_1 \times x_2, x_1 + x_2]_{:}$

• Can make the transfer function **trainable**:

$$\phi(\mathbf{x}) = g(\mathbf{W}' \cdot \mathbf{x} + \mathbf{b}')$$
$$g(x) = \max(x, 0)$$

Multilayer Networks

Binary: $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ $sign(\mathbf{w} \cdot \mathbf{x} + b)$ $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

Multi class:

$$f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$$

 $\arg\max_{i} \operatorname{softmax}(\mathbf{W} \cdot \mathbf{x} + \mathbf{b})_{[i]}$

Non-Linear Classifier

 $f_{\theta}(\mathbf{x}) = \mathbf{w}g(\mathbf{W}' \cdot \mathbf{x} + \mathbf{b}') + b$

Non-Linear Classifier

Multi-layer Perceptron (MLP):

$$f_{\theta}(\mathbf{x}) = \mathrm{NN}_{\mathrm{MLP2}}(\mathbf{x}) = \mathbf{y}$$
$$\mathbf{h}^{1} = g^{1}(\mathbf{x}\mathbf{W}^{1} + \mathbf{b}^{1})$$
$$\mathbf{h}^{2} = g^{2}(\mathbf{h}^{1}\mathbf{W}^{2} + \mathbf{b}^{2})$$
$$\mathbf{y} = \mathbf{h}^{2}\mathbf{W}^{3}$$

Non-Linear Classifier

Multi-layer Perceptron (MLP):

$$f_{\theta}(\mathbf{x}) = \mathrm{NN}_{\mathrm{MLP2}}(\mathbf{x}) = \mathbf{y}$$
$$\mathbf{h}^{1} = g^{1}(\mathbf{x}\mathbf{W}^{1} + \mathbf{b}^{1})$$
$$\mathbf{h}^{2} = g^{2}(\mathbf{h}^{1}\mathbf{W}^{2} + \mathbf{b}^{2})$$
$$\mathbf{y} = \mathbf{h}^{2}\mathbf{W}^{3}$$

Sigmoid

$$\sigma(x) = 1/(1+e^{-x})$$



tanh

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



hard-tanh

$$hardtanh(x) = \begin{cases} -1 & x < -1 \\ 1 & x > 1 \\ x & otherwise \end{cases}$$



ReLU (rectifier, rectified linear unit)

$$\operatorname{ReLU}(x) = \max(0, x) = \begin{cases} 0 & x < 0\\ x & \text{otherwise} \end{cases}$$



How many layers to use? And how wide should they be?

- No hard and fast rules.
- In vision, we see that "deeper is better".
- Not always the case in text / sequences (though with transformers, we may be starting to see this).
- Can think of each layer as **transforming** the previous layer (remember the xor example).
- Narrower layers "compress" the information in the previous layer. Wider layers introduce redundancies.

Dropout

- At each iteration, select a random subset of "neurons" and "drop" them.
- Like training 2ⁿ different networks.
- Prevents co-adaptation of neurons (prevents neurons from depending on each other).

Dropout

$$NN_{MLP2}(\mathbf{x}) = \mathbf{y}$$

$$\mathbf{h^1} = g^1(\mathbf{x}\mathbf{W^1} + \mathbf{b^1})$$

$$\mathbf{h^2} = g^2(\mathbf{h^1}\mathbf{W^2} + \mathbf{b^2})$$

$$\mathbf{y} = \mathbf{h^2}\mathbf{W^3}$$

$$\begin{split} \mathrm{NN}_{\mathrm{MLP2}}(\mathbf{x}) =& \mathbf{y} \\ & \mathbf{h}^{1} =& g^{1}(\mathbf{x}\mathbf{W}^{1} + \mathbf{b}^{1}) \\ & \mathbf{m}^{1} \sim \mathrm{Bernouli}(r^{1}) \\ & \tilde{\mathbf{h}}^{1} =& \mathbf{m}^{1} \odot \mathbf{h}^{1} \\ & \mathbf{h}^{2} =& g^{2}(\tilde{\mathbf{h}}^{1}\mathbf{W}^{2} + \mathbf{b}^{2}) \\ & \mathbf{m}^{2} \sim \mathrm{Bernouli}(r^{2}) \\ & \tilde{\mathbf{h}}^{2} =& \mathbf{m}^{2} \odot \mathbf{h}^{2} \\ & \mathbf{y} =& \tilde{\mathbf{h}}^{2}\mathbf{W}^{3} \end{split}$$

Initialization

- With a (log) linear model, initialization doesn't matter much.
- With MLPs or more complex networks, initialization is **crucial** for achieving good performance.

Initialization

- With a (log) linear model, initialization doesn't matter much.
- With MLPs or more complex networks, initialization is **crucial** for achieving good performance.

$$\mathbb{R}^{n \times m} \sim \text{Uniform}[-\epsilon, +\epsilon]$$

$$\epsilon = \frac{\sqrt{6}}{\sqrt{m+n}}$$

Xavier Glorot et al's suggestion:

Language Identification for 6 languages based on letter bigram counts.



$$\hat{\mathbf{y}} = f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

 $\hat{y} = \operatorname*{argmax}_{i} \hat{\mathbf{y}}_{[i]}$

assume 28 letters (including space).

the vector **x** is 784 dimensional vector.

each entry is the count for a particular letter pair.

$$\hat{\mathbf{y}} = f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

 $\hat{y} = \operatorname*{argmax}_{i} \hat{\mathbf{y}}_{[i]}$

consider the values $\, \hat{\mathbf{y}} \,$

$$\hat{\mathbf{y}} = f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

 $\hat{y} = \operatorname*{argmax}_{i} \hat{\mathbf{y}}_{[i]}$

consider the 6 columns of ${\bf W}$

$$\hat{\mathbf{y}} = f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

 $\hat{y} = \operatorname*{argmax}_{i} \hat{\mathbf{y}}_{[i]}$

consider the 784 rows of W

$$\hat{\mathbf{y}} = f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

 $\hat{y} = \operatorname*{argmax}_{i} \hat{\mathbf{y}}_{[i]}$

think of **x** as a sum of one-hot vectors.

what is **xW** ?

$$\hat{\mathbf{y}} = f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

 $\hat{y} = \operatorname*{argmax}_{i} \hat{\mathbf{y}}_{[i]}$
Representations

what happens if we add layers?

$$\mathbf{\hat{y}} = g(\mathbf{x}\mathbf{W})\mathbf{U}$$

 $\mathbf{W} \in \mathbb{R}^{784 \times 30} \ \mathbf{U} \in \mathbb{R}^{30 \times 6}$

Language Modeling

Let's talk about sequences

• Predicting how a sequence will continue.

_תל א

ראיתי כלב ח_

מה אתה _

Language Model

 $p(x_i|x_1,...,x_{i-1})$

Language Model: Markov Assumption

 $p(x_i|x_1, \dots, x_{i-1}) \approx p(x_i|x_{i-4}, x_{i-3}, x_{i-2}, x_{i-1})$

Language Model: Markov Assumption

$$p(x_i|x_1, ..., x_{i-1}) \approx p(x_i|x_{i-4}, x_{i-3}, x_{i-2}, x_{i-1})$$

(condition only on last n items)

this is called n-gram language model

Language Model

LM can also be used to assign a probability to a sequence.

$$p(x_1, ..., x_n) = p_{LM}(x_1 | *S*, *S*) \\ \times p_{LM}(x_2 | *S*, x_1) \\ \times p_{LM}(x_3 | x_1, x_2) \\ \times p_{LM}(x_4 | x_2, x_3)$$

. . .

 $\times p_{LM}(x_n | x_{n-2}, x_{n-1})$

Language Model

- Very useful (used in Speech Recognition, Machine Translation.. and many others).
- Does not have to be over natural language.
- Huge research topic. We'll see a neural LM.

Neural LM

$$p(x_k | x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1}) = \text{softmax}(\text{MLP}(\mathbf{x}))$$
$$\mathbf{x} = encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$$

Neural LM

$$p(x_k | x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1}) = \text{softmax}(\text{MLP}(\mathbf{x}))$$
$$\mathbf{x} = encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$$

$\operatorname{softmax}(g(g(\mathbf{x}\mathbf{W^1} + \mathbf{b^1})\mathbf{W^2} + \mathbf{b^2})\mathbf{W^3} + \mathbf{b^3})$



 $\mathbb{R}^{d_{out}}$ $\operatorname{softmax}(\Box)$ $\mathbb{R}^{d_{out}}$ \uparrow $\Box W^3 + b^3$ \uparrow \mathbb{R}^{d_2} $g(\Box W^2 + b^2)$ \uparrow \mathbb{R}^{d_1} $g(\Box W^1 + b^1)$ \uparrow $\mathbb{R}^{d_{in}}$ \mathbf{X} \uparrow $encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$



 $\operatorname{encode}(x_1, x_2, x_3, x_4)$

We have k elements in a vocabulary of size |V|

 $\operatorname{encode}(x_1, x_2, x_3, x_4)$

We have k elements in a vocabulary of size |V| 4 10

 $V = \{A, B, C, D, E, F, G, H, I, J\}$

 $\operatorname{encode}(D, A, G, C)$

A = [1,0,0,0,0,0,0,0,0,0]B = [0, 1, 0, 0, 0, 0, 0, 0, 0, 0]C = [0,0,1,0,0,0,0,0,0]D = [0,0,0,1,0,0,0,0,0]E = [0,0,0,0,1,0,0,0,0]F = [0,0,0,0,0,1,0,0,0]G = [0,0,0,0,0,0,1,0,0,0]H = [0,0,0,0,0,0,0,1,0,0]I = [0,0,0,0,0,0,0,0,1,0]J = [0,0,0,0,0,0,0,0,0,1]

A = [1,0,0,0,0,0,0,0,0,0]

[0,0,0,1,0,0,0,0,0,0] + [1,0,0,0,0,0,0,0,0,0] + [0,0,0,0,0,0,0,1,0,0,0] + [0,0,1,0,0,0,0,0,0,0,0] =

 $\mathbf{v}_D + \mathbf{v}_A + \mathbf{v}_G + \mathbf{v}_C$

encode(D, A, G, C)

Encoding k elements

B = [0, 1, 0, 0, 0, 0, 0, 0, 0, 0]C = [0,0,1,0,0,0,0,0,0]D = [0,0,0,1,0,0,0,0,0,0]E = [0,0,0,0,1,0,0,0,0]F = [0,0,0,0,0,1,0,0,0]G = [0,0,0,0,0,0,1,0,0,0]H = [0,0,0,0,0,0,0,1,0,0]I = [0,0,0,0,0,0,0,0,1,0]J = [0,0,0,0,0,0,0,0,0,1]

A = [1,0,0,0,0,0,0,0,0,0]

 $\mathbf{v}_D + \mathbf{v}_A + \mathbf{v}_G + \mathbf{v}_C$

[0,0,0,1,0,0,0,0,0,0]

[1,0,0,0,0,0,0,0,0,0]

[0,0,0,0,0,0,1,0,0,0]

[0,0,1,0,0,0,0,0,0,0]

[1,0,1,1,0,0,1,0,0,0]

what does this miss?

Encoding k elements

encode(D, A, G, C)

A = [1,0,0,0,0,0,0,0,0,0]B = [0, 1, 0, 0, 0, 0, 0, 0, 0, 0]C = [0,0,1,0,0,0,0,0,0]D = [0,0,0,1,0,0,0,0,0] $\mathsf{E} = [0, 0, 0, 0, 1, 0, 0, 0, 0, 0]$ F = [0,0,0,0,0,1,0,0,0]G = [0,0,0,0,0,0,1,0,0,0]H = [0,0,0,0,0,0,0,1,0,0]I = [0,0,0,0,0,0,0,0,1,0]J = [0,0,0,0,0,0,0,0,0,1]

 $\mathbf{v}_D \circ \mathbf{v}_A \circ \mathbf{v}_G \circ \mathbf{v}_C$

encode(D, A, G, C)

 $\mathbf{v}_D \circ \mathbf{v}_A \circ \mathbf{v}_G \circ \mathbf{v}_C$



+ [1,0,0,0,0,0,0,0,0,0,0] + [0,0,0,0,0,0,0,0,0,0] + [0,0,1,0,0,0,0,0,0,0] =

[0,0,0,1,0,0,0,0,0,0]

- F= [-0.28, -0.26, -0.24, 0.31]G= [-0.32, -0.42, -0.21, 0.18]H= [0.00, 0.01, 0.06, 0.14]
- D = [-0.15, -0.31, 0.34, 0.00]E = [-0.10, -0.37, 0.01, 0.40]
- B= [0.29, 0.02, -0.46, -0.39]C= [-0.46, 0.24, -0.16, 0.08]
- $A = \begin{bmatrix} -0.32, 0.09, 0.33, -0.44 \end{bmatrix}$

W

$(\mathbf{v}_D + \mathbf{v}_A + \mathbf{v}_G + \mathbf{v}_C)\mathbf{W}$

+ [1,0,0,0,0,0,0,0,0,0,0] + [0,0,0,0,0,0,0,1,0,0,0] + [0,0,1,0,0,0,0,0,0,0,0] = [1,0,0,1,0,0,0,1,0,0,0]

[0,0,0,1,0,0,0,0,0,0]

- $E = \begin{bmatrix} -0.13, -0.31, 0.34, 0.00 \end{bmatrix}$ $E = \begin{bmatrix} -0.10, -0.37, 0.01, 0.40 \end{bmatrix}$ $F = \begin{bmatrix} -0.28, -0.26, -0.24, 0.31 \end{bmatrix}$
- C = [-0.46, 0.24, -0.16, 0.08]D = [-0.15, -0.31, 0.34, 0.00]
- A= [-0.32, 0.09, 0.33,-0.44] B= [0.29, 0.02,-0.46,-0.39]

W

$= \mathbf{v}_D \cdot \mathbf{W} + \mathbf{v}_A \cdot \mathbf{W} + \mathbf{v}_G \cdot \mathbf{W} + \mathbf{v}_C \cdot \mathbf{W}$

$(\mathbf{v}_D + \mathbf{v}_A + \mathbf{v}_G + \mathbf{v}_C)\mathbf{W}$

$(\mathbf{v}_D + \mathbf{v}_A + \mathbf{v}_G + \mathbf{v}_C)\mathbf{W}$

 $= \mathbf{v}_D \cdot \mathbf{W} + \mathbf{v}_A \cdot \mathbf{W} + \mathbf{v}_G \cdot \mathbf{W} + \mathbf{v}_C \cdot \mathbf{W}$

sum of rows in W

each row corresponds to a certain vocabulary item.

still sum of rows in **W** but **W** has 4x many rows.

	A(-3) = [0.42 - 0.15 0.12 0.02]	
	$B(-3) = \begin{bmatrix} 0.28 & -0.15 & -0.11 & 0.32 \end{bmatrix}$	
	C(-3) = [0.15, -0.24, 0.23, 0.41]	
	D(-3) = [-0.12, -0.24, 0.12, -0.34]	
	E(-3) = [-0.42, -0.21, 0.08, 0.40]	
	F(-3) = [0.20, 0.11, -0.31, 0.33]	
	G(-3) = [0.07, -0.05, 0.16, 0.23]	
	H(-3)= [0.28, 0.03, 0.22, -0.49]	
[0,0,0,1,0,0,0,0,0]	l(-3)= [0.08, 0.39, -0.25, 0.27]	
0	J(-3)= [0.10,-0.42,-0.37, 0.35]	
[1.0.0.0.0.0.0.0]	A(-2)= [-0.00, 0.41, 0.19, 0.49]	
	B(-2)= [0.24, 0.48, 0.34, -0.42]	
	C(-2) = [-0.46, 0.22, 0.24, -0.21]	
[0,0,0,0,0,0,1,0,0,0]	D(-2) = [-0.11, -0.48, 0.18, -0.22]	
0	E(-2) = [-0.32, 0.10, -0.41, -0.43]	
	$F(-2) = \begin{bmatrix} 0.32, 0.02, -0.22, 0.06 \end{bmatrix}$	
[0,0,1,0,0,0,0,0,0]	G(-2) = [-0.31, -0.36, 0.09, 0.39]	
=	H(-2) = [0.01, -0.22, -0.09, -0.15]	
[0.0.0.1.0.0.0.0.0.1.0.0.0.0.0.0.0.0.0.0	I(-2) = [0.01, 0.10, -0.16, -0.21]	
	$J(-2) = \begin{bmatrix} -0.24, 0.40, -0.34, -0.13 \end{bmatrix}$	
	A(-1) = [-0.23, -0.36, 0.02, 0.32] B(-1) = [-0.34, 0.04, 0.18, 0.00]	
	D(-1) = [-0.34, 0.04, -0.10, -0.00] C(-1) = [-0.40, -0.02, 0.10, -0.16]	
	$D(-1) = \begin{bmatrix} 0.40, -0.02, 0.10, -0.10 \end{bmatrix}$	
	$E(-1) = \begin{bmatrix} 0.10, 0.07, 0.10, 0.01 \end{bmatrix}$ $E(-1) = \begin{bmatrix} 0.40, 0.27, 0.33, 0.36 \end{bmatrix}$	
	$F(-1) = \begin{bmatrix} 0.04 & -0.13 & -0.43 & 0.39 \end{bmatrix}$	
	G(-1) = [0.44, 0.38, 0.03, -0.39]	
	H(-1) = [0.41, -0.23, 0.33, -0.08]	
	I(-1) = [-0.50, -0.16, -0.42, -0.27]	
	J(-1) = [-0.15, 0.41, 0.46, -0.16]	
	A(+0)= [-0.11, 0.03, 0.20, 0.50]	
	B(+0) = [0.16, -0.34, 0.20, -0.21]	
	C(+0) = [0.05, -0.13, -0.23, -0.31]	
	D(+0)= [0.13,-0.02, 0.38,-0.09]	
	E(+0)= [0.30, 0.39, 0.10, 0.38]	
	F(+0) = [-0.16, -0.31, -0.02, -0.34]	
	G(+0)= [0.06,-0.04, 0.02,-0.32]	
	$\begin{array}{l} G(+0) = \; [\; 0.06, -0.04, \; 0.02, -0.32] \\ H(+0) = \; [\; 0.25, \; 0.30, \; 0.29, \; 0.24] \end{array}$	
	$\begin{array}{l} G(+0) = \ [\ 0.06, -0.04, \ 0.02, -0.32] \\ H(+0) = \ [\ 0.25, \ 0.30, \ 0.29, \ 0.24] \\ I(+0) = \ [\ 0.40, -0.17, -0.18, -0.19] \end{array}$	

still sum of rows in **W** but **W** has 4x many rows.

alternatively: = $\mathbf{v}_D \cdot \mathbf{W}' + \mathbf{v}_A \cdot \mathbf{W}'' + \mathbf{v}_G \cdot \mathbf{W}''' + \mathbf{v}_C \cdot \mathbf{W}''''$

still sum of rows in **W** but **W** has 4x many rows.

alternatively: = $\mathbf{v}_D \cdot \mathbf{W}' + \mathbf{v}_A \cdot \mathbf{W}'' + \mathbf{v}_G \cdot \mathbf{W}''' + \mathbf{v}_C \cdot \mathbf{W}''''$

 $\mathbf{W} = \mathbf{W}' \circ \mathbf{W}'' \circ \mathbf{W}''' \circ \mathbf{W}''''$

	$\begin{array}{llllllllllllllllllllllllllllllllllll$	
[0,0,0,1,0,0,0,0,0,0]	$\begin{array}{llllllllllllllllllllllllllllllllllll$	W'
Ο	I(-3)= [0.08, 0.39, -0.25, 0.27] $J(-3)= [0.10, -0.42, -0.37, 0.35]$ $A(-2)= [-0.00, 0.41, 0.19, 0.49]$ $B(-2)= [0.24, 0.48, 0.34, -0.42]$ $C(-2)= [-0.46, 0.22, 0.24, -0.21]$	
[1,0,0,0,0,0,0,0,0,0]	D(-2) = [-0.11, -0.48, 0.18, -0.22] E(-2) = [-0.32, 0.10, -0.41, -0.43] F(-2) = [0.32, 0.02, -0.22, 0.06] G(-2) = [-0.31, -0.36, 0.09, 0.39] H(-2) = [0.01, -0.22, -0.09, -0.15]	W''
Ο	I(-2) = [0.01, 0.10, -0.16, -0.21] $J(-2) = [-0.24, 0.40, -0.34, -0.13]$ $A(-1) = [-0.23, -0.38, 0.02, 0.32]$ $B(-1) = [-0.34, 0.04, -0.18, -0.00]$ $C(-1) = [0.40, -0.02, 0.10, -0.16]$	
[0,0,0,0,0,0,1,0,0,0]	D(-1)= [0.13,-0.07,-0.19,-0.01] $E(-1)= [0.40, 0.27,-0.33, 0.36]$ $F(-1)= [0.04,-0.13,-0.43, 0.39]$ $G(-1)= [0.44, 0.38, 0.03,-0.39]$ $H(-1)= [0.41,-0.23, 0.33,-0.08]$ $I(-1)= [-0.50,-0.16,-0.42,-0.27]$	W'''
Ο	$\begin{array}{r} J(-1) = & [-0.15, 0.41, 0.46, -0.16] \\ \hline A(+0) = & [-0.11, 0.03, 0.20, 0.50] \\ B(+0) = & [0.16, -0.34, 0.20, -0.21] \\ C(+0) = & [0.05, -0.13, -0.23, -0.31] \\ D(+0) = & [0.13, -0.02, 0.38, -0.09] \end{array}$	
[0,0,1,0,0,0,0,0,0,0]	E(+0) = [0.30, 0.39, 0.10, 0.38] $F(+0) = [-0.16, -0.31, -0.02, -0.34]$ $G(+0) = [0.06, -0.04, 0.02, -0.32]$ $H(+0) = [0.25, 0.30, 0.29, 0.24]$ $I(+0) = [0.40, -0.17, -0.18, -0.19]$ $J(+0) = [0.27, 0.33, -0.42, -0.07]$	W''''

- 1-hot times matrix: row selection
- sum of 1-hot times matrix: row selection + sum
- concat of 1-hot: like using 1-hot from larger vocab

"Embedding Layer"

- Very common in neural network land:
 - associate each vocabulary item with a row in matrix E of dense vectors (dim of row << |V|)
 - concat or sum rows of **E** for input.

"Embedding Layer"

encode(D, A, G, C)

$= \mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]}$

Ε

[-0.15,-0.31, 0.34, 0.00,-0.32, 0.09, 0.33,-0.44,-0.32,-0.42,-0.21, 0.18,-0.46, 0.24,-0.16, 0.08]
$\mathbb{R}^{d_{out}}$ $\operatorname{softmax}(\Box)$ $\mathbb{R}^{d_{out}}$ $\Box W^3 + b^3$ \mathbb{R}^{d_2} $g(\Box W^2 + b^2)$ \mathbb{R}^{d_1} \uparrow $g(\Box \mathbf{W^1} + \mathbf{b^1})$ \uparrow $\mathbb{R}^{d_{in}}$ $\mathbf{E}_{[x_{k-4}]} \circ \mathbf{E}_{[x_{k-3}]} \circ \mathbf{E}_{[x_{k-2}]} \circ \mathbf{E}_{[x_{k-1}]}$ $encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$

"Embedding Layer"

encode(D, A, G, C)

$= \mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]}$

Ε

[-0.15,-0.31, 0.34, 0.00,-0.32, 0.09, 0.33,-0.44,-0.32,-0.42,-0.21, 0.18,-0.46, 0.24,-0.16, 0.08]

how does this relate to what we had before?

what is $(\mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]})\mathbf{W}$

"Embedding Layer"

 $(\mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]})\mathbf{W}$

same row in **E** regardless of position in the sequence.

but **W** will transform this row differently for each position.

 $\mathbb{R}^{d_{out}}$ $\operatorname{softmax}(\Box)$ $\mathbb{R}^{d_{out}}$ $\Box W^3 + b^3$ \mathbb{R}^{d_2} $g(\Box W^2 + b^2)$ \mathbb{R}^{d_1} \uparrow $g(\Box \mathbf{W^1} + \mathbf{b^1})$ \uparrow $\mathbb{R}^{d_{in}}$ $\mathbf{E}_{[x_{k-4}]} \circ \mathbf{E}_{[x_{k-3}]} \circ \mathbf{E}_{[x_{k-2}]} \circ \mathbf{E}_{[x_{k-1}]}$ $encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$

 $\mathbb{R}^{d_{out}}$

what's in **W1**? what's in **W2**? what's in **W3**?

 $\operatorname{softmax}(\Box)$ $\mathbb{R}^{d_{out}}$ $\Box W^3 + b^3$ \mathbb{R}^{d_2} $g(\Box W^2 + b^2)$ \mathbb{R}^{d_1} $g(\Box W^1 + b^1)$ $\mathbb{R}^{d_{in}}$ $\mathbf{E}_{[x_{k-4}]} \circ \mathbf{E}_{[x_{k-3}]} \circ \mathbf{E}_{[x_{k-2}]} \circ \mathbf{E}_{[x_{k-1}]}$ $encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$



A Neural Probabilistic Language Model



• What is the cost of increasing the history length?

• We can use a trained LM for scoring a given sentence. (how?)

• We can use a trained LM for comparing two given sentences. (how?)

We can use a trained LM for generating new sentences. (how?)

• We can use K trained LMs for k-class classification. (how?)

• We can use K trained LMs for **k-class** classification. (how?) $p(w_1, w_1, ..., w_n | LM_1)$

$$p(w_1, w_1, ..., w_n | LM_1)$$

$$p(w_1, w_1, ..., w_n | LM_2)$$

$$p(w_1, w_1, ..., w_n | LM_3)$$

$$\hat{y} = \arg\max_k p(w_1, w_1, ..., w_n | LM_k)$$

sequence prediction tasks

- In LM:
 - predict word i based on k previous words.
- But we could also predict label based on k items.
 - which tasks can this be used for?





classification

ACGCGCTCGATCG V

GTGCGCTCGAACG V

Х



tagging

DET ADJ NOUN VERB PREP DET ADJ NOUN The brown fox jumped over the lazy dog

tagging

NOUN

The brown **fox** jumped over

VERB

brown fox jumped over the

PREP

fox jumped **over** the lazy

Window-based approach

back to LM

Training a language model

- Set dimensions of **E**, **W3**, according to vocab size.
- Initial random values for E, W1, W2, W3, b1, b2, b3
- For every n-tuple in some text:
 - try to predict last item based on prev n-1
 - use cross-entropy loss.

What happens after training?

- Consider the columns of **W3**.
- Consider the rows of **E**.

 $\mathbb{R}^{d_{out}}$ $\operatorname{softmax}(\Box)$ $\mathbb{R}^{d_{out}}$ $\Box W^3 + b^3$ what's in **W3**? \mathbb{R}^{d_2} $g(\Box W^2 + b^2)$ columns of W3 \mathbb{R}^{d_1} correspond $g(\Box W^1 + b^1)$ to vocab items! $\mathbb{R}^{d_{in}}$ $\mathbf{E}_{[x_{k-4}]} \circ \mathbf{E}_{[x_{k-3}]} \circ \mathbf{E}_{[x_{k-2}]} \circ \mathbf{E}_{[x_{k-1}]}$ $encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$



toeory as

10 Neural Worder Embeddings

тJ



000 1 15 20 30 10 few half five six two three four several some many other those^{these} all both

chief head chairman director spokesman executive trader

Country and Capital Vectors Projected by PCA 2 China ------≫Beijing 1.5 **Russi**a Japan Moscow 1 Ankara ⇒√Tokvo Turkey -0.5 Poland Germany -0 France ≫Warsaw Berlin 📉 -0.5 Italy -Paris Athens Greece Rome Spain -1 Madrid -1.5 Portugal ->>Lisbon -2 -0.5 -1.5 0.5 1.5 -1 -2 0 2 1

Target Word	BoW5
batman	nightwing
	aquaman
	catwoman
	superman
	manhunter
hogwarts	dumbledore
	hallows
	half-blood
	malfoy
	snape
turing	nondeterministic
	non-deterministic
	computability
	deterministic
	finite-state
florida	gainesville
	fla
	jacksonville
	tampa
	lauderdale
object-oriented	aspect-oriented
	smalltalk
	event-driven
	prolog
	domain-specific
dancing	singing
	dance
	dances
	dancers
	tap-dancing

- We trained a language model.
- We ended up with vector representations of words.
- These representations are useful -- they encode various aspects of word similarity.

tagging + pre-training

we can use the **E** we got from LM training to initialize **E** for the POS tagging task. NOUN

The brown **fox** jumped over

VERB

brown fox jumped over the

(why is that helpful?)

PREP

fox jumped over the lazy

pre-training

- This is a sort of semi-supervised learning or multi-task learning.
- We learn from "unannotated" data.
- We then use the representations on tasks with annotated data.

- We trained a language model.
- We ended up with vector representations of words.
- These representations are useful -- they encode various aspects of word similarity.
- A form of semi-supervised learning.
- More next week.