# Self-Supervised Learning Feed-Forward Nets

#### CSCI 601 471/671 NLP: Self-Supervised Models

https://self-supervised.cs.jhu.edu/sp2023/



[Slide credit: Andrej Karpathy, Fei Fei Li, Chris Manning and many others]

#### HW<sub>2</sub> is released

- Did you see it?
- Due Tuesday noon, in 120 hours!

# Recap: Contextual Representation of Meaning

- Language is complex, and *context* can completely change the meaning of a word in a sentence.
- Example:
  - I let the kids outside to *play*.
  - He had never acted in a more famous *play* before.
  - It wasn't a *play* the coach would approve of.
- Previous models (e.g., Word2Vec) only have one representation per word
  - They can't capture these ambiguities.
- Need a model which captures the different nuances of the meaning of words given the surrounding text.
- **Approach:** build representations that are **contextual** via **neural networks**.

#### **Chapter Plan**

- 1. Defining neural networks (feed-forward nets)
- 2. Neural nets: brief history
- 3. Word2Vec as a simple neural network
- 4. Training neural networks: analytical back-propagation
- 5. Backprop in practice

**Chapter goal:** Get really comfortable with thinking, designing and building neural networks — very powerful modeling tools.

#### Neural Network

- Neural Networks are functions!
  - Function class for approximating real-valued, discrete-valued and vector valued target functions.
  - NN:  $X \to Y$  where  $X = [0,1]^n$ , or  $\mathbb{R}^n$  and  $Y = [0,1]^d$ ,  $\{0,1\}^d$
- Example: A **2-layer** neural network
  - The input, hidden and output variables are represented by nodes
  - The links are the weight parameters
  - Arrows denote direction of information flow through the network

 $f(\mathbf{x}) = W_2 g(W_1 \mathbf{x}) \qquad \mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^d$ 

 $g(\mathbf{z}) = [\sigma(z_1), \dots, \sigma(z_h)]$  (nonlinearity)  $\sigma(z_i) = \frac{1}{1+e^{-x}}$  (sigmoid function)

•  $W_1 \in \mathbb{R}^{h \times n}$  and  $W_2 \in \mathbb{R}^{d \times h}$  are the parameters that need to be learned.



#### Neural Network: Making it bigger

Add more layers, or wider layers!

A **2-layer** neural network



#### A **3-layer** neural network



- This is actually a particular class called "feed-forward" networks.
  - Cascade neurons together
  - Output from one layer is the input to the next
  - Each layer has its own sets of weights



• Inputs multiplied by initial set of weights



• Intermediate "predictions" computed at first hidden layer



- Intermediate predictions multiplied by second layer of weights
- Predictions are fed forward through the network



• Compute second set of intermediate predictions



• Multiply by final set of weights



- Aggregate all the computations in the output
  - e.g. probability of a particular class



• All the intermediate parameters are ought to be learned.



### Why Add Non-linearity?

• Without non-linearity, the overall model amounts to a linear model.

$$f(\mathbf{x}) = W_2 g(W_1 \mathbf{x})$$
  

$$\tilde{f}(\mathbf{x}) = W_2 W_1 \mathbf{x} = W_3 \mathbf{x} \text{ (a linear function)}$$
  

$$drop g$$

- A linear function cannot approximate complex tasks.
- Non-linearity adds capacity to the model to approximate any continuous function to arbitrary accuracy given sufficiently many hidden units.
  - See <u>"universal approximation theorem"</u>



Cannot separate red and blue points with linear classifier

#### Activation/Nonlinearity Functions



Leaky ReLU  $\max(0.1x, x)$ 

 $\begin{array}{l} \mathsf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$ 



10

#### Demo time!

• Link: <u>https://playground.tensorflow.org/</u>

#### Quiz Time

- Given a neural network and an input, how do you compute its predictions?
  - You start from the input. Calculate the output of each layer (starting from the first layer), until you get to the output.
  - 2. You start from the end and make your way to the first layer.

Quiz Time (2)

- What is needed to fully specify a neural network?
  - <sup>1.</sup> Architecture (which input goes through what function etc.)
  - <sup>2.</sup> Parameters of the function (the weights)
  - 3. Both

## Quiz Time (3)

- What makes neural networks expressive functions?
  - 1. Activations (non-linearities)
  - 2. Depth (number of hidden layers)
  - <sup>3.</sup> Width (number of variables in each hidden layer)
  - 4. All the above

#### **Chapter Plan**

- 1. Feed-forward networks
- 2. Neural nets: origins and brief history
- 3. Word2Vec as a simple neural network
- 4. Training neural networks: analytical back-propagation
- 5. Backprop in practice

## Artificial Neurons: An Inspiration from Nature

- A single node in your neural network
  - Accept information from multiple inputs
  - Transmit information to other neurons
- A neuron's function is inspired by its biological counterpart:
  - Apply some function on inputs signals
  - If output of function over threshold, neuron "fires"





#### Artificial Neurons: Not Quite Analogous to Nature

Biological neurons: complex connectivity



Neurons in an artificial neural network: organized based on a highly regular structure for computational efficiency



#### Very Brief History of Neural Networks

- 1. Single-layer neural networks (1943-1969)
- 2. Symbolic AI & knowledge engineering (1970-1985)
- 3. Multi-layer NNs and symbolic learning (1985-1995)
- 4. chirp chirp .... Statistical learning/probabilistic models (1995-2010)
- 5. Deep networks and self-supervised learning (2010-?)

## A Neuron as a Mathematical Model of Computation



• McCulloch and Pitts (1943) showed how linear threshold units can be used to compute logical functions



An alternative model of computation (comparable to "Turing Machine")

# Perceptron: Imitating Nature's Learning Process

- Rosenblatt (1959) developed the Perceptron Algorithm, an iterative, hillclimbing algorithm for learning the weights of a linear threshold unit.
  - A single neuron with a fixed input, it can incrementally change weights and learn to produce a fixed output using the Perceptron learning rule.
- Update weights by:

$$w_i = w_i + \eta(t - o)x_i$$

- where  $\eta$  is the "learning rate," *t* is the teacher output, and *o* is the network output.
- If output is correct do nothing.
- If output is higher than *t*, lower weights on active inputs
- If output is lower than t, increase weights on active inputs

#### Perceptron: Demise

- *Perceptons* (1969) by Minsky and Papert illuminated the limitations of the perceptron.
- It showed that:
  - Shallow (2-layer) networks are unable to learn or represent many classification functions (e.g. XOR)
  - Only the linearly separable functions are learnable.
- Also, there was an understanding that deeper networks were infeasible to train.
- Result: work on neural-networks dissipated during the 70's and early 80's!



#### Neural Net Resurgence (1986)

- Interest in NNs revived in the mid 1980's due to the rise of "connectionism."
- Backpropagation algorithm was [re-]introduced for training three-layer NN's.
  - Generalized the iterative "hill climbing" method to approximate networks with multiple layers, but no convergence guarantees.



[Learning representations by back-propagating errors, Rumelhart, Hinton & Williams 1986; for a broader context, see: <u>http://people.idsia.ch/~juergen/who-invented-backpropagation.html</u>]

#### Second NN Demise (1995-2010)

- Generic Back-Propagation did **not** generalize that well to training **deeper** networks.
  - Overfitting / underfitting remained an issue.
  - Computers were still quite slow
- Little theoretical justification for underlying methods.
- Machine learning research moved to graphical/probabilistic models and kernel methods.

#### Deep Learning Revolution (2010...)

- Improved methods developed for training deep neural works.
- Particular successes with:
  - Convolutional neural nets (CNNs) for vision (2012 AlexNet showed 16% error reduction on ImageNet).
  - Recurrent neural nets (RNNs) for machine translation and speech recognition.
  - Deep reinforcement learning for game playing (e.g., AlphaGo).
- Massive data and specialized hardware:
  - Large collections of [mostly] supervised (crowdsourced) training data has been critical.
  - Efficient processing of this big data using specialized hardware (Graphics Processing Units, GPUs) has been critical.

#### Deep Learning Revolution (2010...)

#### How it start

#### How it's going



#### **Chapter Plan**

- 1. Feed-forward networks
- 2. Neural nets: brief history
- 3. Word2Vec as a simple neural network
- 4. Training neural networks: analytical back-propagation
- 5. Backprop in practice

#### **Chapter Plan**

- 1. Feed-forward networks
- 2. Neural nets: brief history
- 3. Word2Vec as a simple neural network
- 4. Training neural networks: analytical back-propagation
- 5. Backprop in practice

#### Word2Vec as Neural Net

• Word2Vec is actually a very simple neural net!

#### Word2Vec (SkipGram): Recap



d

d

words

 $|V_w|$ 

- **Objective:** find embeddings such that the words that co-occur are close.
- Represent each word as a *d* dimensional vector.
- Represent each context as a *d* dimensional vector.



#### Word2Vec (SkipGram): Recap



- **Objective:** find embeddings such that the words that co-occur are close.
- Represent each word as a *d* dimensional vector.
- Represent each context as a *d* dimensional vector.





Probability that if you

randomly pick a word

nearby "ants", that it is "car"



#### Word2Vec as Neural Network

• Word2Vec is basically a very simple feed-forward neural network!



- Any algorithms for neural networks could also be applied to Word2Vec.
  - For example, we can use Backpropagation (forthcoming!) to train Word2Vec. (HW3!!)
- We can extend W<sub>2</sub>V by modifying the architecture.
  - Contextual self-supervised representations: next chapter!

#### **Chapter Plan**

- 1. Feed-forward networks
- 2. Neural nets: brief history
- 3. Word2Vec as a simple neural network
- 4. Training neural networks: back-propagation
- 5. Backprop in practice



#### Derivatives

- First let's get the notation right:
- The arrow shows functional dependence of z on y,
   i.e. given y, we can calculate z.
  - For example:  $z(y) = 2y^2$





#### Quiz time!

• If  $z(x, y) = y^4 x^5$  what is the following derivative  $\frac{\partial z}{\partial y}$ ?

1. 
$$\frac{\partial z}{\partial y} = 4y^3 x^5$$
  
2.  $\frac{\partial z}{\partial y} = 5y^4 x^4$   
 $\frac{\partial z}{\partial y} = 2x^3 x^4$ 

3. 
$$\frac{\partial z}{\partial y} = 20y^3x^4$$

4. None of the above



#### Gradient

• Given a function with 1 output and *n* inputs

$$f(\mathbf{x}) = f(x_1, x_2, \dots, x_n) \in \mathbb{R}$$



• Its gradient is a vector of partial derivatives with respect to each input

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n$$

(always assume vectors are column vectors, i.e., they're in  $\mathbb{R}^{n \times 1}$ )

#### Quiz time!

- If  $z(x, y) = y^4 x^5$  what is the following gradient  $\nabla z$ ?
  - 1.  $\nabla z(x,y) = 4y^3x^5$
  - 2.  $\nabla z(x, y) = (5y^4x^4, 20y^3x^4)$
  - 3.  $\nabla z(x, y) = (5y^4x^4, 4y^3x^5)$
  - 4. None of the above



#### Jacobian Matrix: Generalization of the Gradient

- Given a function with *m* outputs and *n* inputs
  - $\mathbf{f}(\mathbf{x}) = [f_1(x_1, x_2, \dots, x_n), \dots, f_m(x_1, x_2, \dots, x_n)] \in \mathbb{R}^m$
- It's Jacobian is an **m** x **n** matrix of partial derivatives:  $(\mathbf{J}_{\mathbf{f}}(\mathbf{x}))_{ij} = \frac{\partial f_i}{\partial x_i}$

$$\mathbf{J}_{\mathbf{f}}(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n} \text{ or } \left( \mathbf{J}_{\mathbf{f}}(\mathbf{x}) \right)_{ij} = \frac{\partial f_i}{\partial x_j}$$

- When m=1 (scalar-valued function), Jacobian reduces to  $\nabla^T \mathbf{f}(\mathbf{x})$  (gradient transpose).
- When m=n=1 (single-variable function), Jacobian reduces to the derivative of **f**.



#### Jacobian for Matrix Inputs

• Given a function with *m* outputs and *n*×*p* inputs

$$\mathbf{f}(\mathbf{X}) = [f_1(\mathbf{X}), \dots, f_m(\mathbf{X})] \in \mathbb{R}^m, \text{ where } \mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix} \in \mathbb{R}^{n \times p}$$

• Jacobian is a  $m \times n \times p$  **tensor** (i.e., matrix of matrices) of partial derivatives:

$$\left(\mathbf{J}_{\mathbf{f}}(\mathbf{X})\right)_{ijk} = \frac{\partial f_i}{\partial x_{jk}}$$

• The Jacobian math holds if you keep adding more dimensions to the input or output.

| ; [`<br>u u?<br>5 w&<br>G u%<br>+ H | `{*Ui`N?><br>?'xJ@8 C<br>\$\dC\$ixX<br>6 oKY:'U<br>i v!F1_i  | [# = R / J<br>4 3 Q d l \$<br>8 g u l [<br>F X 8 % k<br>k i 8 V N   | ) r _ Rt<br>D + . /&<br>Z N K e <<br>L T i K e<br>L T i   |   | 2Y7 Fq{y` V<br>8sr qe\$", G<br>A-x J3/{. r<br>rtN <fuwk 8<br="">r4r &lt;&lt;02* 8</fuwk>  |   |  |
|-------------------------------------|--|---|---|---|---|---|--|
| Wł                                  | ny Use Matrix/T  | ensor F   | orm?  | , , , , , ,<br>A 3%] a 0<br>- EQSM > 0<br>0 ge-; j uW<br>r a>gAw 10   | FW, u1u71<br>FW, l2u71<br>CK  >]T_<br>q G 4]JU<br>W 0 1%]   |   |  |
| Te<br>Jt<br>h)                      | I+0J R G[<br>N8;v *S<br>dpYi V c{  | icq4 e/17<br>{Uta 3Q('<br>oOfv gofY   | - Z 6 L s W<br>T # 7 i z } }<br># R 8 V 1 R I   | * 0 % 5 0 V 0 4<br>X 4 N V ' Z '<br>g { Z " - M {   | * C L @ } x<br>* C L @ } x<br>\$ K U < e 9  | Y 1Y ; e<br>  5 < g +<br>0 nK C   |  |
| <mark>ہِ in</mark> e                | essence, matrix form (   | (multi-varia  | te calculus) i  | is just an exte   | ension of sin   | gle-variable  | /<br>3 K   |
| *7 <b>Ca</b>                        | CUIUS. , } 4 N G . t )<br>H \$ v t % G 2 h A<br>6 : & Y olh. (<br>i E [X D 6 Y D 9 #<br>Z 2 / d # G & #  | <pre><kiii 0="" 7="" <="" th="" w="" z="" {<=""><th>D T K G U O<br/>",/ T #9<br/>p X K 7 #<br/>y X H A 5<br/>g X 1 Z 3 N<br/>N I C N</th><th>, { w = h + g A +<br/>h ^ 6 Z # * G j<br/>G ] % , <b>s</b> f 7 e<br/>&amp; 0 M v &gt; g y<br/>  1 + j i R N<br/>M Z I &gt; t</th><th>Q 0 0 x 7<br/>B J = ^ i<br/>F 0 3 J 9<br/>E h g r i<br/>1 &amp; t w L<br/>&gt; \$ a E c</th><th></th><th></th></kiii></pre> | D T K G U O<br>",/ T #9<br>p X K 7 #<br>y X H A 5<br>g X 1 Z 3 N<br>N I C N   | , { w = h + g A +<br>h ^ 6 Z # * G j<br>G ] % , <b>s</b> f 7 e<br>& 0 M v > g y<br>  1 + j i R N<br>M Z I > t | Q 0 0 x 7<br>B J = ^ i<br>F 0 3 J 9<br>E h g r i<br>1 & t w L<br>> \$ a E c   |   |  |
| <sup>h</sup> Tw                     | (O reasons: @ ₩ & L n<br>= T : & V<br>₩ p n { k  | + q / c<br>y Z Y %<br>, % 8   |   |   |   | U LOY'Z f Ba<br>) W P @ F / - a @<br>9 = # ,   D K Y  | ≠ TU<br>TU^F   |
| 0                                   | Compact derivations  | : with mate   | rix form calc   | ulations we c   | an compute  | a concise 🔄 🗄   | 5 0 % L  |
|                                     | statements.z         T         T         R           i         T         Y         +         -           i         B         Y         I         m           i         B         Y         I         0           i         B         Y         I         0           a         I         c         c         H | 2 y 0<br>1 M z<br>@ ' G =<br>2 o H \ r<br>@ 7 \ H   | 3 d H Q 1 a<br>) Q 9 Z I Y<br>  < L Z I<br>a + u 3 Z<br>9 q f q<br>: d   W  |   |   | f 0 i V u E a \ j<br>f ] 6 0 [ ( L x R F<br># % X ; " 6 o 4 3 0<br>+ v < 1 . V E b 7<br>! } R T Y ! ]   6<br>K a r 1 o 7 z v  | N N<br>N N<br>N N<br>N N<br>N N<br>N N<br>N N<br>N N<br>N N<br>N<br>N N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N<br>N |
| 0                                   | Implementing algorit   | thms in mat   | trix form is n  | nuch faster.  | IQ   (u<br>^U F X c   |   | /" Z,<br>Ri 7{   |
|                                     | <sup>1</sup> O <sup>4</sup> GPUs are optimized   | for VERY FAS  | r matrix/tensor   | operations.   | nz aax<br>{d WV-<br>XB bC\  | K 43 Z & (<br>!8   h & * ! ;<br>w3 g i r + 0  | инк,<br>2%^0<br>х л  |
| * 0 & o 3 ? e +   v<br>B P w ) 6    | ,  | D.<br>R. 9;<br>c. ;<br>h.r.:<br>N.N.r.  | I A w D g<br>P a m v Z w<br>d \ ib + 1<br>Y r r S V o<br>_ > @ h d i *<br>T p f   x i<br>Y P l n B g<br>P 8 h p # C | W U U V V U V V U V U V V U V V V V V V   | IQ 1 3 p<br>IQ 1 3 p<br>IQ 1 3 p<br>JF<br>UI jF<br>UI jF<br>V<br>V<br>V<br>C 1 C<br>I<br>V<br>V<br>C 1 C<br>I<br>V<br>I<br>V<br>V<br>C 1 C<br>I<br>V<br>I<br>V<br>I<br>V<br>I<br>V<br>I<br>V<br>I<br>V<br>I<br>V<br>I<br>V<br>I<br>V<br>I | MA ) e !/ e i<br>MA ) e !/ e i<br>V' Z9" [u R '<br>V' Z9" [u R '<br>V' b" b R '<br>U X 3 }1 n<br>R p 9 G c t g<br>0 (<br>R p 9 G c<br>] 0 (<br>0 (<br>) 0 ( ) 0 (<br>) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( ) 0 ( | с п А А А А А А А А А А А А А А А А А А  |

#### Chain Rule

• Function composition:

$$z \circ y(x) = z(y(x)) = z(x)$$

If *z* is a function of *y*, and *y* is a function of *x*, then *z* is a function of *x*, as well.



Then:

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

## Chain Rule for Multivariable Functions

- Let  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{g}: \mathbb{R}^d \to \mathbb{R}^n$ ,  $\mathbf{f}: \mathbb{R}^n \to \mathbb{R}^m$
- Composing them:  $\mathbf{f} \circ \mathbf{g}(\mathbf{x}) = \mathbf{f}(\mathbf{g}(\mathbf{x})): \mathbb{R}^d \to \mathbb{R}^m$

The result looks similar to the single-variable setup:

$$\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{x}) = \mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{x})) \ \mathbf{J}_{\mathbf{g}}(\mathbf{x})$$

Note, the above statement is a **matrix** multiplication! Function  $\mathbf{f} \circ \mathbf{g}$  has *m* outputs and *d* inputs  $\rightarrow m$  by *d* Jacobian



**Quiz Time!** 

Let 
$$x \in \mathbb{R}$$
,  $y: \mathbb{R} \to \mathbb{R}^n$ ,  $z: \mathbb{R}^n \to \mathbb{R}$ 



What is the Jacobean of  $z \circ \mathbf{y}(x) = z(y_1(x), ..., y_n(x))$ ? 1.  $\mathbf{J}_{z \circ \mathbf{y}}(x) = \mathbf{J}_z(\mathbf{y}(x)) \ \mathbf{J}_{\mathbf{y}}(x)$ 2.  $\mathbf{J}_{z \circ \mathbf{y}}(x) = \left[\frac{\partial z}{\partial y_1}, ..., \frac{\partial z}{\partial y_n}\right] \left[\frac{\partial y_1}{\partial x}, ..., \frac{\partial y_n}{\partial x}\right]^{\mathrm{T}}$ 3.  $\mathbf{J}_{z \circ \mathbf{y}}(x) = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \ \frac{\partial y_i}{\partial x}$ 4. All the above!

## We will stop here!

• Today: lots of background about neural networks!

• Next time: training a neural net!