# Tuesday, March 10

## Midterm 1

Minimum: 46
Median: 74
Mean: 74
Maximum: 99
Standard deviation: 16

Grades in canvas!
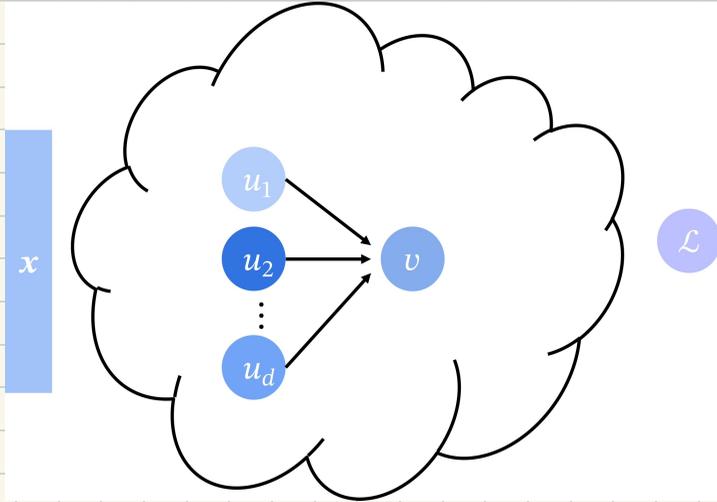
Long tail, will curve e.g. 50% to D+/C-

## Plan

- Backpropagation
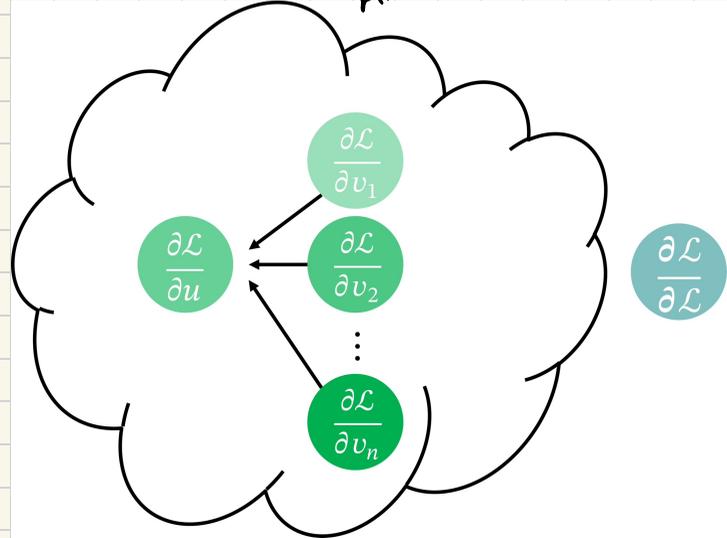- Convolutional Neural Networks

# Forward

Already know $u_1, \ldots, u_d$



$$v(u_1, \ldots, u_d)$$

# Backward

Already know $\frac{\partial \mathcal{L}}{\partial v_1}, \ldots, \frac{\partial \mathcal{L}}{\partial v_n}$



$$\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial v_1} \frac{\partial v_1}{\partial u} + \ldots + \frac{\partial \mathcal{L}}{\partial v_n} \frac{\partial v_n}{\partial u}$$

# Linear Algebraic View of Backpropagation

Neural nets known since 1950's

Backprop known since 1960's

Why now?

GPUs used starting 2010's

really good at matrix multiplication, so training and inference cast as matrix multiplication

**Forward**

$$v = Wu \qquad\qquad v_i = \sum_{j=1}^{d} [W]_{i,j} \, u_j \qquad\qquad \frac{\partial v_i}{\partial u_j} = [W]_{i,j}$$

**Backward**

$$\nabla_u \mathcal{L} = W^T \nabla_v \mathcal{L} \qquad\qquad \frac{\partial \mathcal{L}}{\partial u_j} = \sum_{i=1}^{n} \frac{\partial \mathcal{L}}{\partial v_i} \frac{\partial v_i}{\partial u_j} = \sum_{i=1}^{n} [W^T]_{j,i} \frac{\partial \mathcal{L}}{\partial v_i}$$

Note: Actually want gradient of parameters

$$\frac{\partial \mathcal{L}}{\partial [W]_{i,j}} = \frac{\partial \mathcal{L}}{\partial v_i} \frac{\partial v_i}{\partial [W]_{i,j}} = \frac{\partial \mathcal{L}}{\partial v_i} u_j$$

$$\nabla_W \mathcal{L} = (\nabla_v \mathcal{L}) u^T$$

Neural networks are a powerful tool for learning complex functions, but notoriously difficult to understand because of their many parameters and nonlinearities. The big surprise of deep learning is that, even though their loss landscapes are not convex functions, gradient descent works remarkably well to train neural networks. We don't have good intuition for what happens in high-dimensional space (the loss function depends on the number of parameters, which can be very large), but it *seems* like there's almost always a parameter update that at least slightly improves the loss function so we rarely get stuck in a bad local minimum.

# Motivating Convolutional Layers

Goal: Classification on an image

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}$$
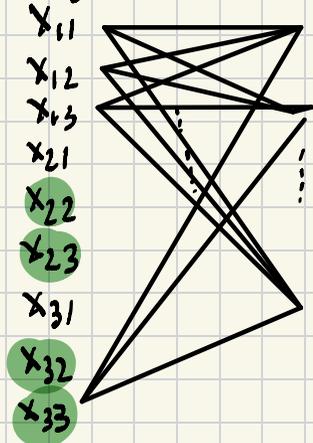
3x3

← Each value is color of a pixel

We would like:

1. Fewer parameters in layer

2. Preserve context

3. Shift invariance

Using linear layer...

$x_{11}$
$x_{12}$
$x_{13}$
$x_{21}$
$x_{22}$
$x_{23}$
$x_{31}$
$x_{32}$
$x_{33}$

Q: If 400 × 400 pixels, Linear layer where input has same size as output, how many parameters?

# Convolutional Neural Network

"kernel"/weights

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} * \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} w_{11}x_{11} + w_{12}\cdot x_{12} + x_{21}w_{21} & w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23} \\ & \\ \ldots & \ldots \end{bmatrix}$$
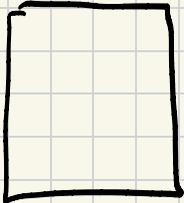
Q: $\begin{bmatrix} 1 & 1 & 1 \\ -1 & -1 & -1 \\ 1 & 1 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}$

$* \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}$

Q: If 400 × 400 pixels, Convolutional layer with d×d kernel, how many parameters?

# CNNs continued:

- Padding : add variables around edges
- Strides : distance between kernel applications
- Pooling : aggregate values
- Channels and 3d convolution :

$$\left[\left[\quad\right]\right] \star \left[\begin{smallmatrix}1 & 1 \\ -1 & -1\end{smallmatrix}\right] = \boxed{\phantom{xxx}}$$

- Transposed Convolution : <u>Increase</u> dimension

Input                    kernel

$$\begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix} \star^{T} \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix} + \begin{bmatrix} 2 & 2 \\ -2 & -2 \end{bmatrix} + \begin{bmatrix} 3 & 3 \\ -3 & -3 \end{bmatrix}$$

Q: How do we train ?

## Thursday, March 12

- Happy spring break!
- No OH Monday 3/23

## Plan

Self-attention

We process

↳ "rich" features with linear layers

↳ "local" features (image/audio/video) with convolutional layers

Q: How do we process sequential data such as text?

↑
size of text
changes eg
sentence, paragraph, page

Goals:

↳ Decompose text into tokens

↳ Represent tokens as vectors

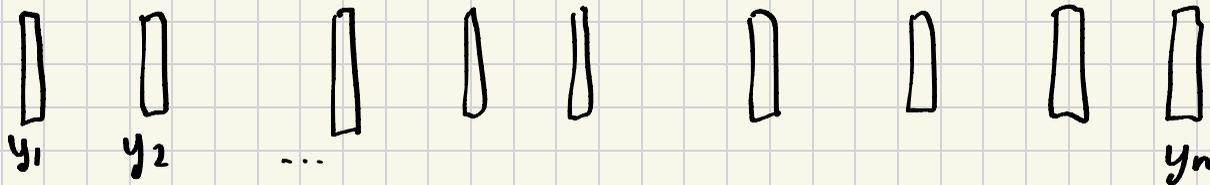↳ Process "similar" tokens together

↳ Encode order of tokens

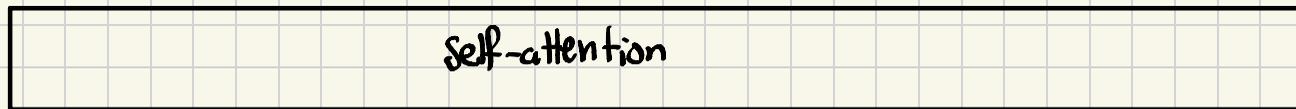Input: The weather today is warm and sunny.

Tokens: "the" "weather" "today" "is" "warm" "and" "sun" "ny" "."

Vectors:

↰ using
autoencoders $x_1$ $x_2$ ... $x_n$

| self-attention |

$y_1$ $y_2$ ... $y_n$

Notice: same number of inputs and outputs

# Intuition

Write $y_i$ as linear combination of $x_j$ weighted by similarity.

$$y_i = \sum_{j=1}^{n} \langle x_i, x_j \rangle \, x_j = \sum_{j=1}^{n} \langle q_i, k_j \rangle \, v_j$$

similarity between $x_i$ and $x_j$

input vectors

Notice each $x_k$ appears in three different places:

↳ Key

↳ query

↳ value

---

where are the parameters??

$$q_i = W_q \, x_i$$

$$k_j = W_k \, x_j$$

$$v_j = V_j \, x_j$$

$$X \in \mathbb{R}^{n \times d}$$

$$Q = X W_q, \quad K = X W_k, \quad V = X W_v$$
$$\phantom{Q = X W_q,} {\scriptstyle d \times d} \quad\quad {\scriptstyle d \times d} \quad\quad {\scriptstyle d \times d}$$

$$S = Q K^T = X W_q W_k^T X^T$$
$$\phantom{S = Q K^T =} {\scriptstyle n \times n}$$

$$[S]_{i,j} = \langle q_i, k_j \rangle$$

$$S = \begin{bmatrix} \langle q_1, k_1 \rangle & \langle q_1, k_2 \rangle \dots \\ & \vdots \end{bmatrix} \frac{1}{\sqrt{d}}$$
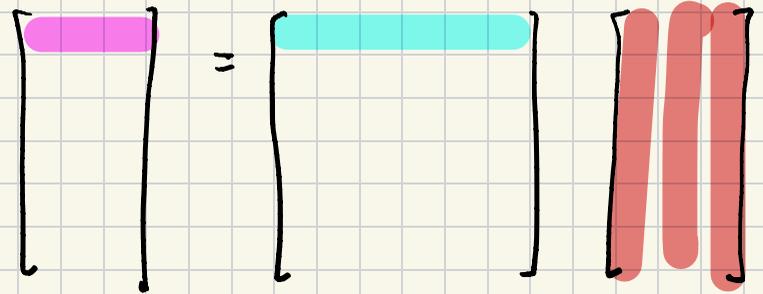"scores"

$$A = \begin{bmatrix} \text{softmax}( \langle q_1, k_1 \rangle & \langle q_1, k_2 \rangle \dots \\ \text{softmax}( & ) \\ & \vdots \end{bmatrix}$$
"attention"

$$Y = A V$$
$$\phantom{Y =} {\scriptstyle n \times n} \; {\scriptstyle n \times d}$$



$$Y_i = \sum_{j=1}^{n} \text{sim}(i,j) \, V_j$$

Suppose we have three tokens (eg "I love climbing") represented in $\mathbb{R}^2$

$$X = \begin{bmatrix} 1 & 1 \\ 0 & -1 \\ 1 & 0 \end{bmatrix} \qquad W_q = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \qquad W_k = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \qquad W_v = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$K =$ $\qquad$ $V =$ $\qquad$ $Q =$

$A =$
$\uparrow$ no need to simplify

## Attention Mask

Goal: We want tokens to only "look back"

$$
A = \begin{array}{c} \\ \\ \\ \\ i \end{array}
\begin{bmatrix}
1 & & & & j & \\
.3 & .7 & & & & \\
.1 & .4 & .5 & & & \\
0 & 0 & .1 & .9 & & \\
0 & .2 & 0 & .6 & .2 & \\
.9 & 0 & 0 & 0 & 0 & .1
\end{bmatrix}
$$

$[A]_{i,j}$ is how much token $i$ attends to $j$
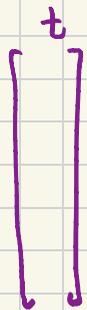
Q: Side Benefit with auto regression?

# Positional Encoding

"Jack gave water to Jill"

vs.

"Jill gave water to Jack"

$$\begin{bmatrix} \\ \\ t \\ \\ \\ \end{bmatrix}$$

$\sin(2^0 t)$

$\sin(2^1 t)$

$\vdots$

$\sin(2^n t)$

$t \quad t+1$