

Week 6

19/30/25

↳ Quiz

↳ Last class: 3 4 8 1 0 2.4 avg

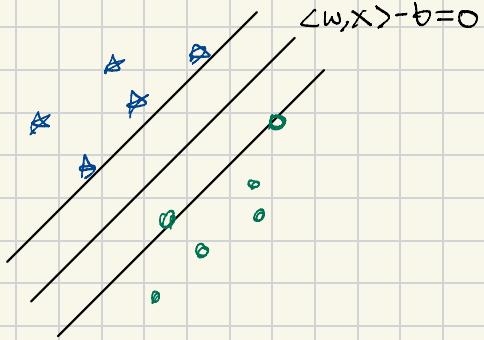
- most mathematical topic
- working on problem helped?

↳ Pset 5 due tonight

↳ Flipped class Thursday

↳ OH Friday 10-11:30, 12:30-1:00

Review



Today: Why is $\langle x^{(i)}, x^{(j)} \rangle$ so useful?

- a new (simple) algorithm
- kernel trick (why dual sum)
- extend trick to other algorithms

Primal

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \quad \text{s.t. } y^{(i)} [\langle w, x^{(i)} \rangle - b] - 1 \geq -\xi_i, \quad \xi_i \geq 0$$

variables = $d+1+n$ // $w \in \mathbb{R}^d$, $b \in \mathbb{R}$, $\xi \in \mathbb{R}_+^n$

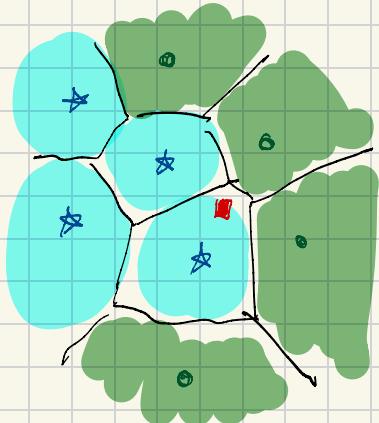
$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \boxed{\langle x^{(i)}, x^{(j)} \rangle} \quad \text{s.t. } \alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y^{(i)} = 0$$

variables = n // $\alpha \in \mathbb{R}_+^n$

K-Nearest Neighbors

- On new point x ,
- find k points "closest" to x from training points
- aggregate labels
 - ↳ mode for classification
 - ↳ average for regression

Voronoi

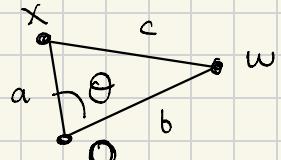


"close" by cosine similarity

$$\cos(\theta) = \frac{\langle x, w \rangle}{\|w\|_2 \|x\|_2}$$

↑ angle between x, w

Proof:



$$a = \|x\|_2$$

$$\cos \theta = \frac{a^2 + b^2 - c^2}{2ab} \quad \text{by Law of Cosines}$$

$$\cos(\theta) = \frac{\langle x, w \rangle}{\|x\|_2 \|w\|_2}$$

Equivalent:

- large angle θ
- small cosine sim $\cos \theta$
- large distance $\|x - w\|_2$

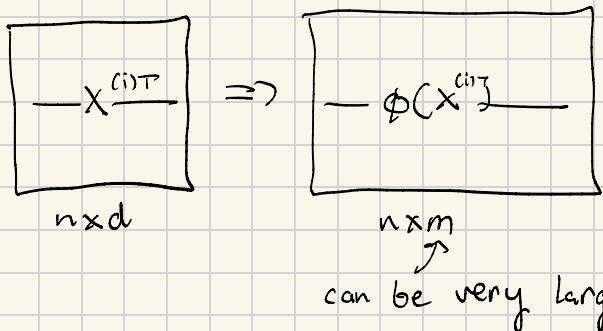
MNIST: 0 1 2 3 4 5 6 7 8 9

x_{new}	0	1	2	3	4	5	6	7	8	9
2	2	2	2	2	2					
7	7	7	7	7	7					
9	9	9	9	4	9					

Non-parametric (like Naive Bayes), how can we improve?

Feature Transform

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$$



$$\text{Kernel} \quad K(x, x') = \langle \phi(x), \phi(x') \rangle$$

$$K \in \mathbb{R}^{n \times n} = \begin{matrix} \phi(x) & \phi(x)^T \\ n \times m & m \times n \end{matrix}$$

Big Q: How is this faster?

Polynomial Kernel

$$K(x, x') = (\langle x, x' \rangle + 1)^p \quad \text{for integer } p$$

$\uparrow O(d)$ time

$$\phi(x) = \begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_3 \\ x_1^2 \\ x_2^2 \\ x_3^2 \\ \sqrt{2}x_1 x_3 \\ \sqrt{2}x_1 x_2 \\ \sqrt{2}x_2 x_3 \end{bmatrix} \quad \begin{matrix} d=3 \\ p=2 \end{matrix}$$

$\leftarrow O(d^p)$

Gaussian Kernel

$$K(x, x') = e^{-\frac{\|x - x'\|_2^2}{2\sigma^2}} \quad \leftarrow O(d) \text{ time}$$

$\phi(x)$ is infinite series

Reparameterization Trick

What if $\mathbf{x}\mathbf{x}^T$ doesn't naturally appear?

Replace \mathbf{x}_w with $\mathbf{x}\mathbf{x}^T \mathbf{v}$

Why is this okay?

$$\text{WLOG, } w = \sum_{i=1}^n v_i x^{(i)} = \mathbf{x}^T \mathbf{v}$$

Proof: Suppose for contradiction not

i.e., $w \neq v + u$

$$u^T x^{(i)} = 0 \quad \forall i$$

$$\mathbf{x}_w = \mathbf{x}_v + \mathbf{x}_u = \mathbf{x}_v$$

Logistic Regression

$$\underset{w}{\operatorname{argmin}} \mathcal{L}_{\text{CE}}(\sigma(\mathbf{x}_w), y)$$

vs

$$\underset{v}{\operatorname{argmin}} \mathcal{L}_{\text{CE}}(\sigma(\mathbf{x}\mathbf{x}^T \mathbf{v}), y)$$

Linear Regression

$$\underset{w}{\operatorname{argmin}} \mathcal{L}_{\text{MSE}}(\mathbf{x}_w, y)$$

vs

$$\underset{v}{\operatorname{argmin}} \mathcal{L}_{\text{MSE}}(\mathbf{x}\mathbf{x}^T \mathbf{v}, y)$$

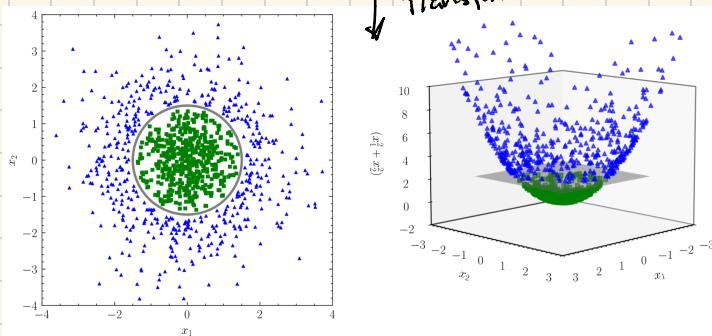
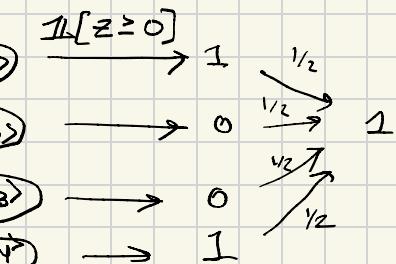
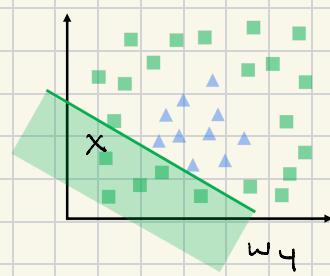
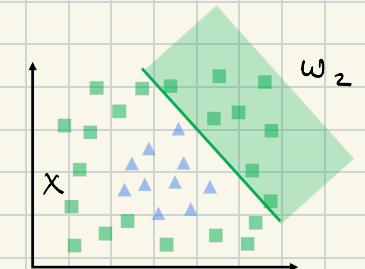
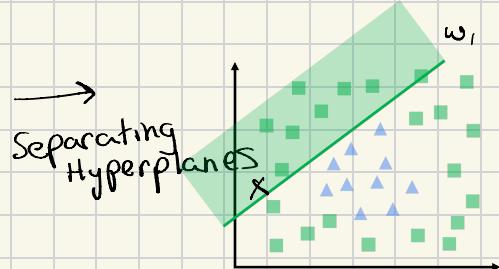
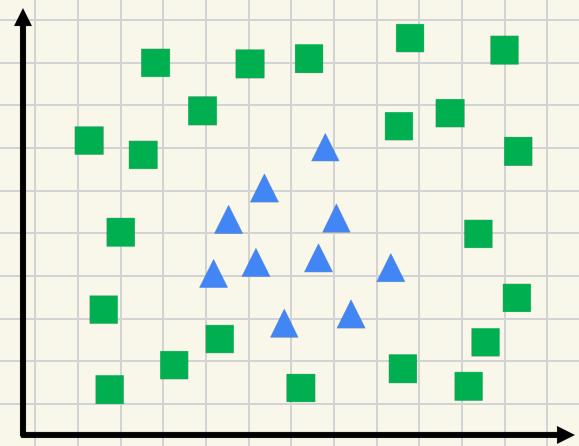
10/2/2025

So far, we have:

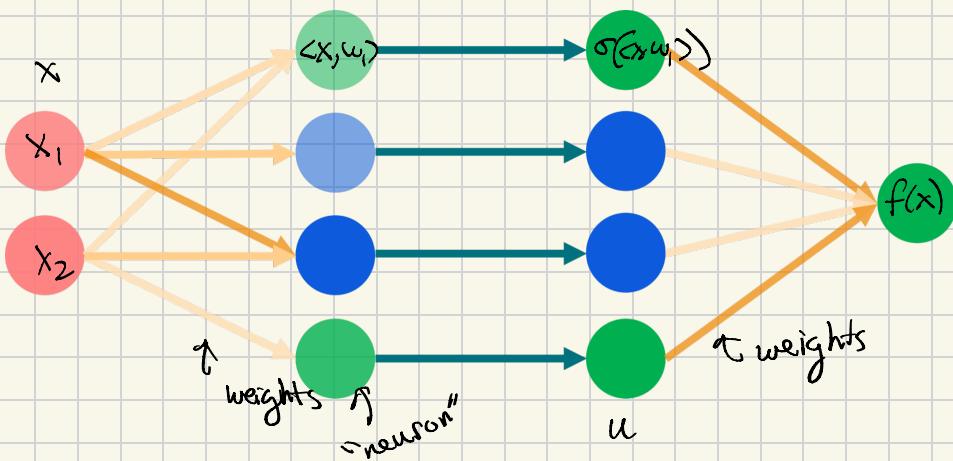
1. Selected features (and transformed them)
2. Trained a model on the features

Neural networks will let us do both simultaneously

Neural Networks



Neural Network Architecture



Choices Include :

- ↳ Activation function, "ReLU"
 - $\sigma(z) = \max(0, z)$ "ReLU"
 - $\sigma(z) = \frac{1}{1+e^{-z}}$ "Sigmoid"
- ↳ Number of layers
- ↳ Number of neurons
- ↳ Loss function
 - MSE
 - Cross Entropy

nn.Sequential(
 nn.Linear(2, 4), nn.ReLU(), nn.Linear(4, 1))

$$w^{(2)} \sigma(w^{(1)}x)$$

Pytorch

`m. Linear(2, 4) nn. relu() nn. Linear(4, 1)`

$$W \in \mathbb{R}^{4 \times 2}$$

$$\sigma(w)$$

$$W \in \mathbb{R}^{1 \times 4}$$

Linear Algebra

Optimization

Big Q: How do we compute $\nabla_{\omega} \mathcal{L}$ (quickly)?

Chain Rule Review

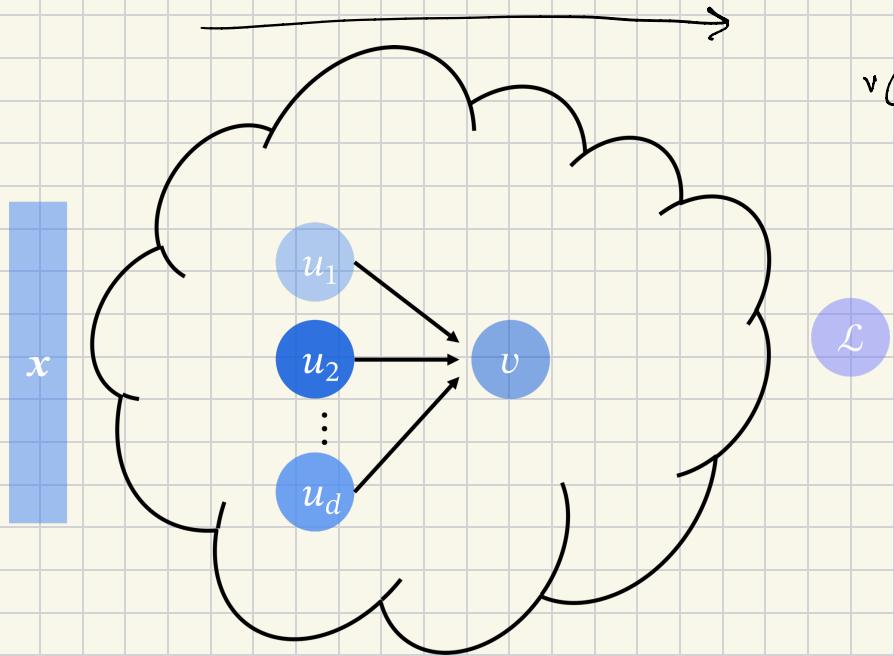
single variable

$$\left. \begin{array}{l} \frac{df}{dx} = \lim_{t \rightarrow 0} \frac{f(x+t) - f(x)}{t} \\ f: \mathbb{R} \rightarrow \mathbb{R} \quad y: \mathbb{R} \rightarrow \mathbb{R} \\ \frac{d}{dx} f(y(x)) = \lim_{t \rightarrow 0} \frac{f(y(x+t)) - f(y(x))}{t} = \lim_{t \rightarrow 0} \frac{f(y(x+t)) - f(y(x))}{y(x+t) - y(x)} \frac{y(x+t) - y(x)}{t} = \frac{df}{dy} \frac{dy}{dx} \end{array} \right.$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R} \quad y_i: \mathbb{R} \rightarrow \mathbb{R}$$

$$\frac{df}{dx}(y_1(x), \dots, y_m(x)) = \left(\frac{df}{dy_1} \frac{dy_1}{dx} + \dots + \frac{df}{y_m} \frac{dy_m}{dx} \right)$$

Forward Pass



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

And so on!

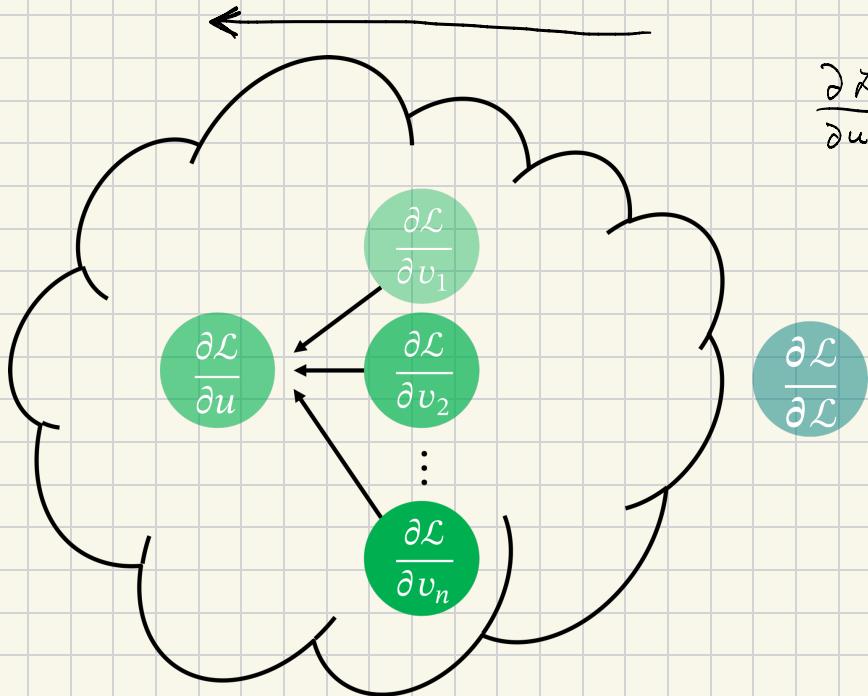
Pytorch:

```
pred = model(x)
```

```
loss = loss_fn(pred, y)
```

Insight: Only need u_1, \dots, u_d to compute v

Backward Pass



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

And so on!

Pytorch:

`loss.backward()`

Insight: Only need $\frac{\partial \mathcal{L}}{\partial v_1}, \dots, \frac{\partial \mathcal{L}}{\partial v_n}$ and $\frac{\partial v_1}{\partial u}, \dots, \frac{\partial v_n}{\partial u}$ to compute $\frac{\partial \mathcal{L}}{\partial u}$

Linear Algebraic View

Neural Nets invented in 1950s

Backprop proposed in 1980s

Q: Why the wait?

A: Only really useful with MANY weights

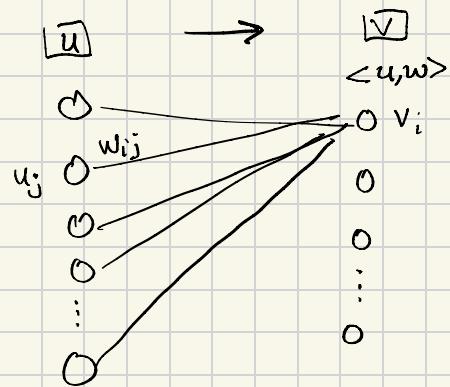
(parameters), which takes storage and compute

In 2010s, researchers realized GPUs

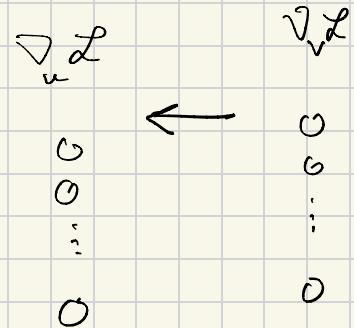
(good at matrix mult for video games)

could be used for neural nets!

Forward



Backward



$$\frac{\partial \mathcal{L}}{\partial u_j} = w_{ij} \quad \nabla_u \mathcal{L} = W^T \nabla_v \mathcal{L}$$

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

Note: We want $\nabla_w \mathcal{L}$, easy from gradient of neurons

Q: Why would Gradient descent work when function highly non-convex?

- Maybe high dimensional loss functions rarely have local minima
- Not well explained by theory

Next week: Two very import architectures

- Convolutions: How can we make networks efficient?
 - ↳ Linear(1000, 1000) has 1m param
- Transformers: How can we process sequential data?
 - ↳ eg text, video