

---

---

---

---

---



# Ethical Algo Design in the Generative Era



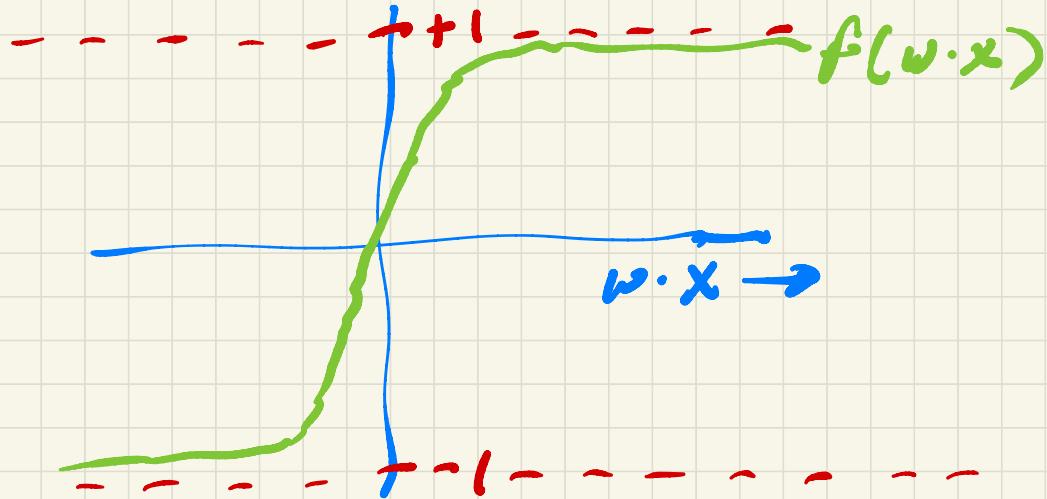
# Generative AI

- What is it ?
- E.g. LLMs
- Next-word prediction
- Training, annotation
- History / precursors
- Embeddings
- Autoencoders / decoders

# Neural Net Background

- single "neuron": on input  $x \in \mathbb{R}^d$ ,  
output =  $f(w \cdot x)$   
 $\approx f(v_1 x_1 + v_2 x_2 + \dots + v_d x_d)$

- typical choice for  $f$ : sigmoid



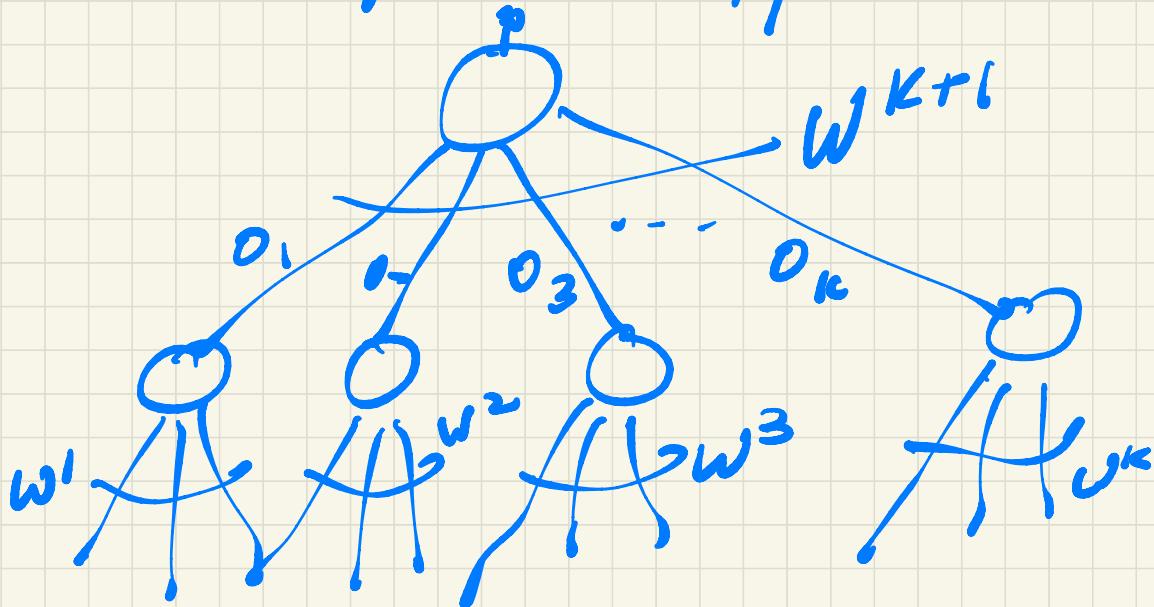
- training: weights  $w$  are parameters, loss fn is

$$l(w) = \sum_{(x,y) \in S} (f(w \cdot x) - y)^2$$

differentiable

# Multiple Layers

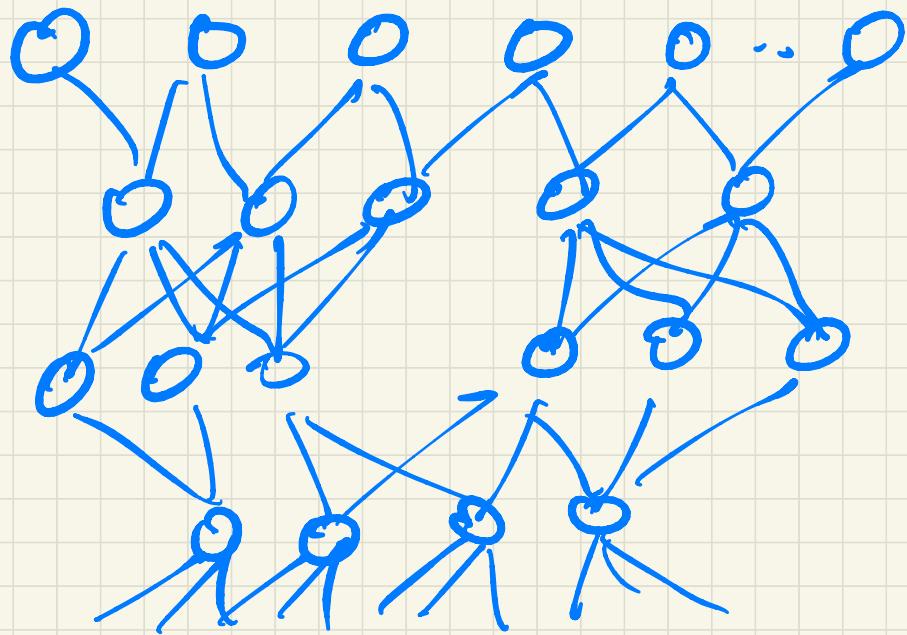
Final output



- params now  $w^1, w^2, \dots, w^{k+1}$
- still differentiable  
(chain rule  $\rightarrow$  backprop)

key idea: intermediate  
neurons can learn  
new/better features

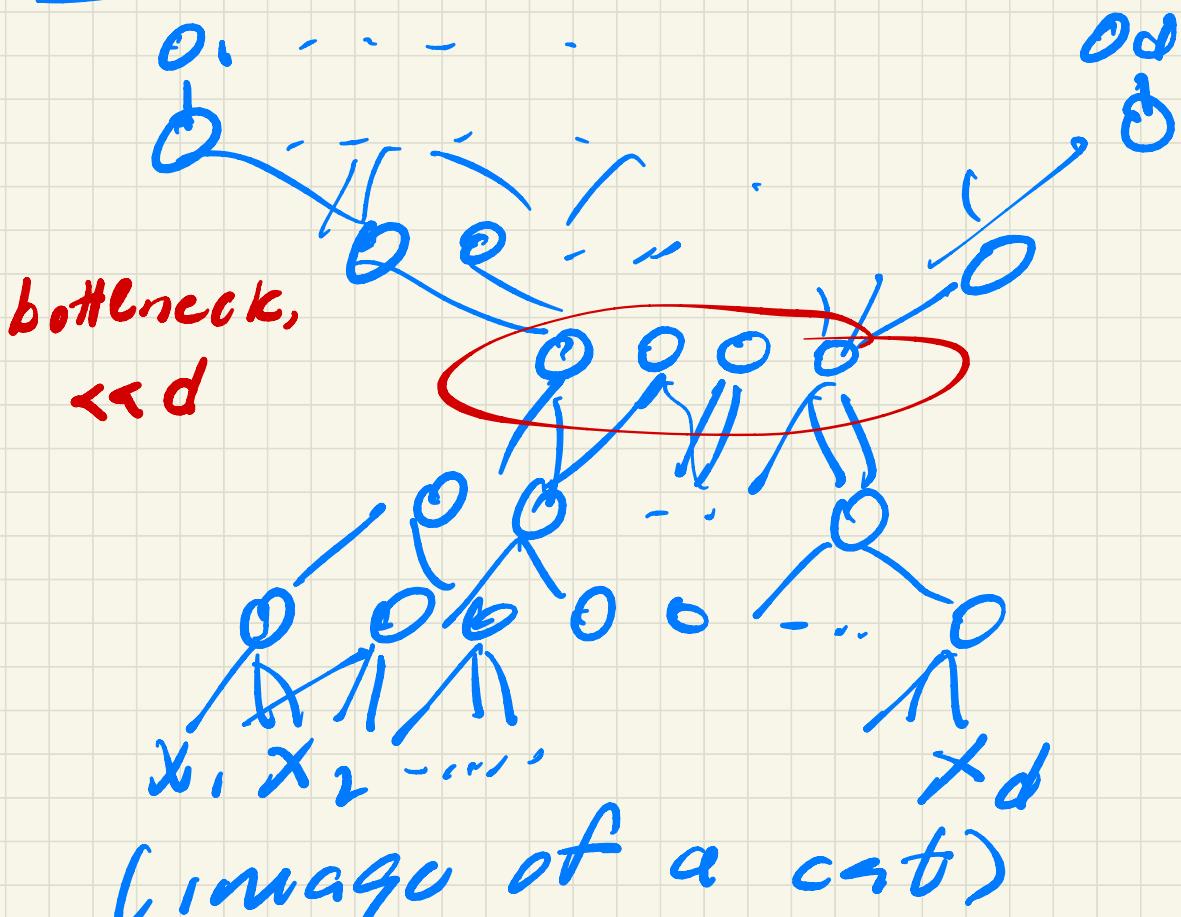
# Multiple Outputs



$x_1, x_2, \dots, x_d$

- now outputs are high dimensional, still differentiable

# Autoencoders



goal: minimize

$$\sum_{i=1}^d (o_i - x_i)^2$$

# Fairness in GAI

- Compare e.g. group fairness in consumer lending
- Fairness in pronouns?
- In tone?
- Biases vs. correlations, use cases

# Toxicity

- Subjectivity  
  & context
- Quotations  
  & censorship
- Offensive  
  opinions
- Nuance &  
  indirection

# Toxicity & Fairness

- Training data curation
- Must automate
  - human annotation (unpleasant)
- Train "guardrail" models
- Apply to training, prompts & output

# Privacy in GAI

- Regurgitation
- Virtual copying  
(e.g. code)
- Stylistic mimicry  
(more later)
- Could do DP model training, but impractical
- Robustness to membership inference attacks

# Copyright & IP Concerns

- Stylistic mimicry
- Threat to creative livelihoods?

# Copyright & IP

- Tech, policy & legal mechanisms
- Content attribution & compensation
- Machine Unlearning & model disengagement
  - DP
  - Sharding

# Adversarial Attacks

- E.g. image manipulation,  
prompt injection,  
data poisonings,  
spam filter attacks, ...
- Robust adversarial training:
  - instead of minimizing e.g.  
$$\sum_{(x_i, y_i)} (h(x_i) - y_i)^2$$
  - minimize something like  
$$\max_{\delta: \|\delta\|_1 \leq d} \sum (h(x_i + \delta_i) - y_i)^2$$

# Plagiarism & Cheating

- College essays,  
writing samples, etc.
- Tool vs. cheat ?
- Author verification

# Author Verification

- Train model to detect human/LLM (arms race)
- Redlist/greenlist (watermarking)

# Hallucinations

- Plausible but verifiably wrong
- E.g. citations

# Hallucinations

- User education
- Independent verification
- External sources  
(e.g. citation DBs)
- Disclaimers
- Content attribution

A bit more on hallucinations...  
(Kalai & Vempala, MM&MK, HW2)

- stylized facts, e.g.  
IMDB records - unambiguous  
true/false (also citations)
- schema or tuples:  
 $\langle \text{author}, \dots, \text{title}, \text{year}, \text{journal} \rangle$
- huge space of "factoids",  
only a fraction are facts
- training data contains  
only facts
- safest: only generate  
facts from training
- But in ML we want  
to generalize...

# Missing Mass Estimation

- How many facts are missing in data?

- Good-Turing: let

$$\underline{\mathcal{L}} = \# \text{facts occurring only once in training}$$

$n$

- Then

| prob. next sample -  $\mathcal{L}/n$  /  
| is a new fact

- Kalai-Vempala: for any model

hallucination rate  $\geq \mathcal{L} - \text{model miscalibration}$   
(one measure of overfitting)

# Disruption to Work

- Passing bar exams,  
MBA courses, etc.
- IP concerns again

# Job disruption/creation

- Prompt engineers
  - “English is the new programming language”

Strongest  
defense:

Use case

specialization