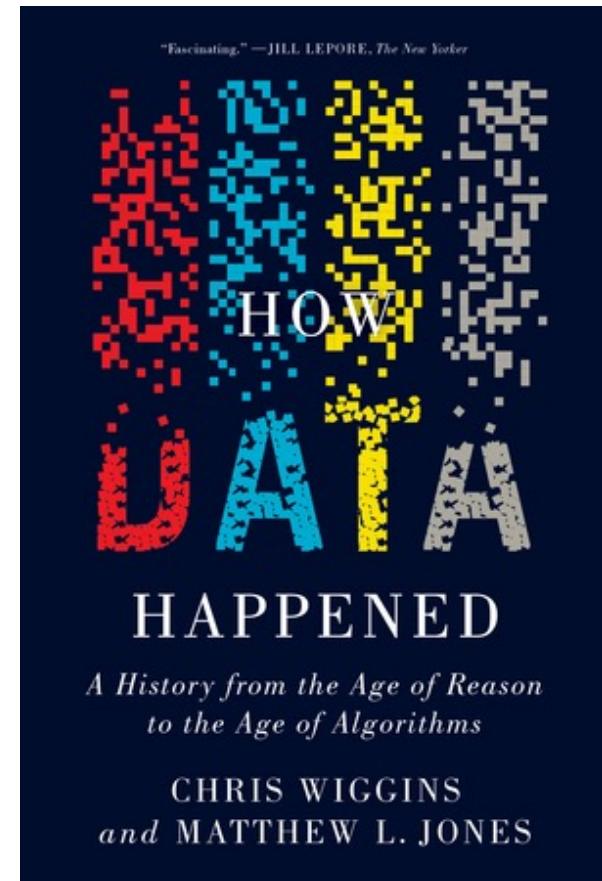


What we'll cover today

- History of machine learning
- Epistemology of machine learning
- How can we use it for social scientific research?
- Under the hood of important ML techniques

History of machine learning

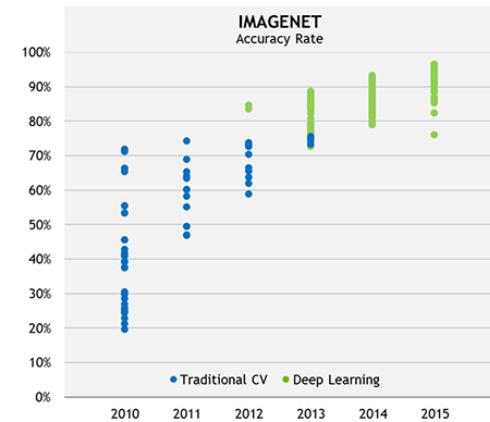
- Traces its origins back to the mid-20th century, early pioneers like Alan Turing.
- NSA pursued Signal Intelligence (SIGINT) Analysis: had a lot of data, and less interested in “it works” than “doing it right”
- The field gained momentum in the 1950s and 1960s with the development of simple neural networks and the exploration of algorithms like the perceptron.
- "AI winters" of the 1970s and 1980s due to limited computational power and high expectations.
- The revival of machine learning came in the late 1980s and early 1990s, fueled by the advent of more powerful computers.



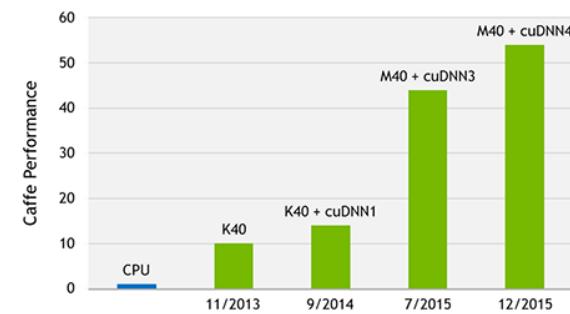
“Big Bang of Machine Learning”

- GPU processing + Massive data from platforms + Deep learning.
- Powerful deep Artificial Neural Networks.
- 2022 ChatGPT revolution: transformers and all text on the internet.

2015: A MILESTONE YEAR IN COMPUTER SCIENCE



50X BOOST IN DEEP LEARNING IN 3 YEARS



Machine learning and Surveillance Capitalism

- Economic dimension of the “Big Bang”
- Machine learning become the foundation for platforms:
 - Online targeted advertising.
 - Behavioral modification.
 - Maximizing engagement.
- These are based on prediction logic – they don’t care why it works.
- “Will this individual buy the product if we show them the advertisement?”

Machine learning is what makes your personal data valuable. This is what makes machine learning valuable.



What is machine learning?

- Boundary between statistics and machine learning is blurry
- Some methods are both, such as linear regression
- Less about techniques than about the culture and aims

- ML focuses on $\mathbf{Y} = f(\mathbf{X}) = \mathbf{X}^T \boldsymbol{\beta}$, statistics focuses on .



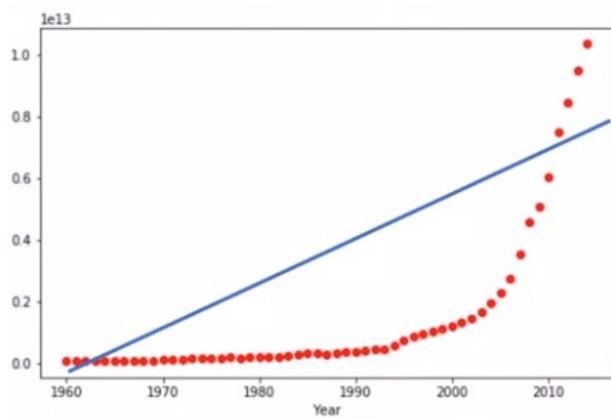
The Two Cultures

- Statistics:
 - *Understand* how an outcome is related to inputs.
 - Test theory
- Machine learning:
 - *Predict* new unseen data
 - Models don't need to be understandable
 - Inductive

Machine learning	Statistics
network, graphs	model
weights	parameters
learning	fitting
generalization	test set performance
supervised learning	regression/classification
unsupervised learning	density estimation, clustering
large grant = \$1,000,000	large grant = \$50,000
nice place to have a meeting: Snowbird, Utah, French Alps	nice place to have a meeting: Las Vegas in August

Advantages with Machine Learning

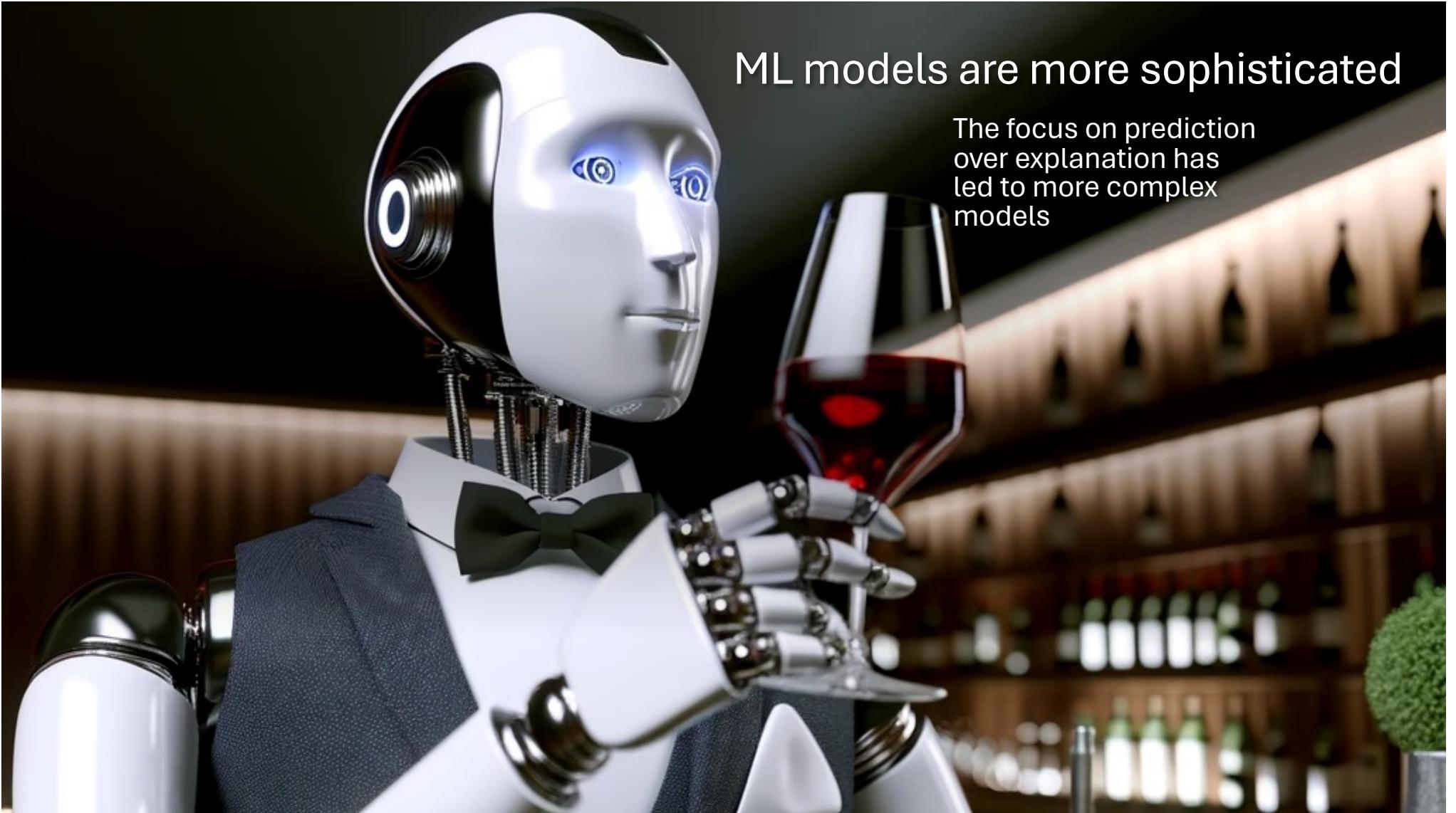
Reality is not linear!



Population heterogeneity and causal complexity

- A cause doesn't always have the same effect
- Social scientific theory usually contains “if's” and “then's”, but they cannot be captured by quantitative models.
- *On average* is often a terribly destructive way of dealing with data



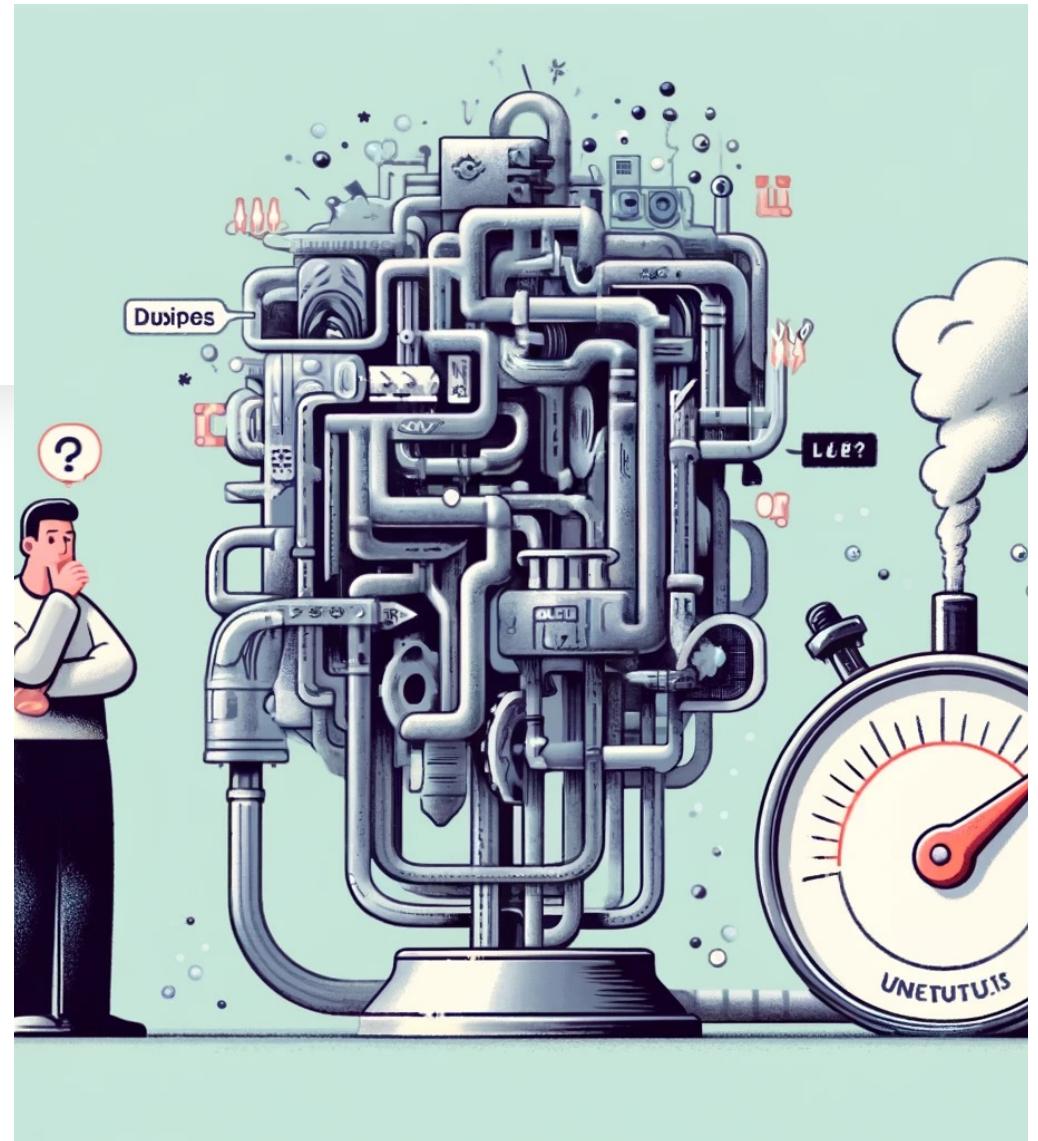


ML models are more sophisticated

The focus on prediction over explanation has led to more complex models

Downsides of ML

- The predictions can be an opaque function of the inputs: hard to answer the *why*?
- The models may not easily provide meaningful estimates of uncertainty
- Estimation often more computationally intensive



Two classes of machine learning models

Unsupervised models

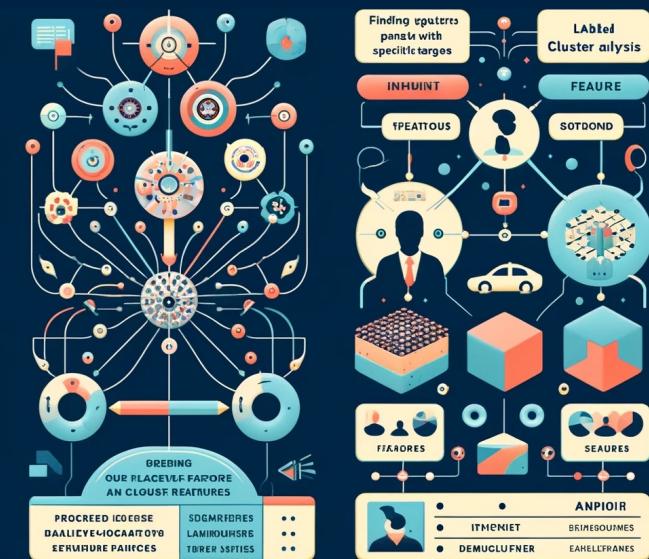
Identify pattern in data, or reduce dimensions.

- No target output to predict, no teacher showing the algorithm what it should aim for, and no immediate measure of success
- E.g. topic modeling, cluster analysis.

Supervised models

Mapping between features and outcomes.

- Uses labelled training data.
- E.g. identifying race from profile picture.



Reinforcement learning: learn-by-doing

- Train AI for skills we don't fully understand: we know how to walk, but it's awfully hard to explain.
- Trial-and-error
- The better the AI walks, the more reward they get.
- The learning can be implemented in many different ways



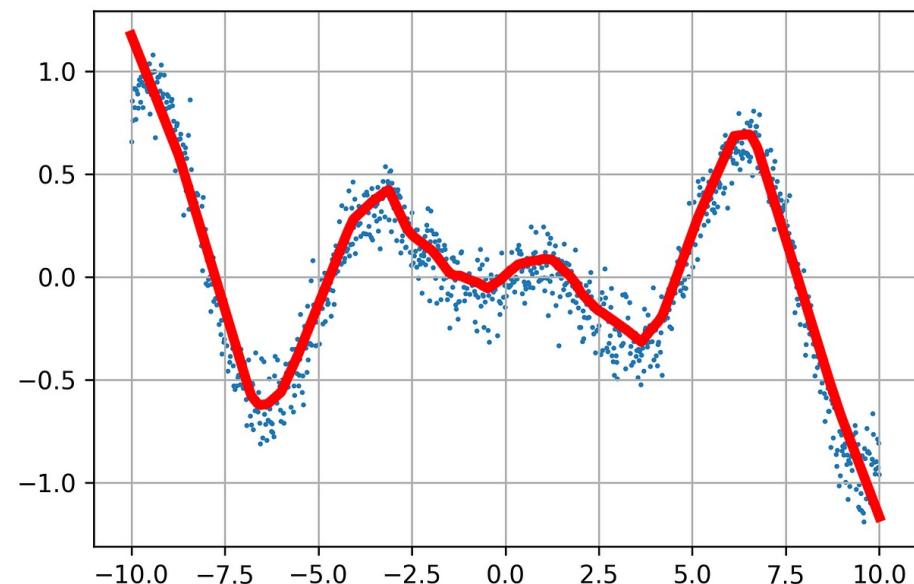
Fitting supervised models

Optimization of parameters with respect to an objective function.

Often uses the good-old OLS model (Ordinary Least Squares)

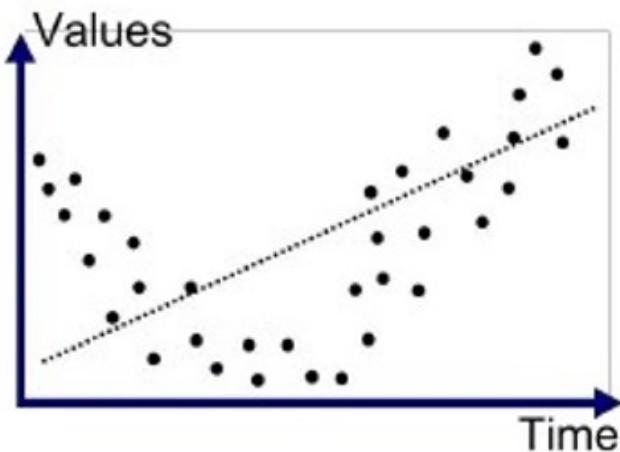
$$\hat{\beta} = \arg \min_{\beta} \sum_i (y_i - \beta' x_i)^2.$$

What distinguishes machine learning methods is that they can fit substantially more flexible functions.

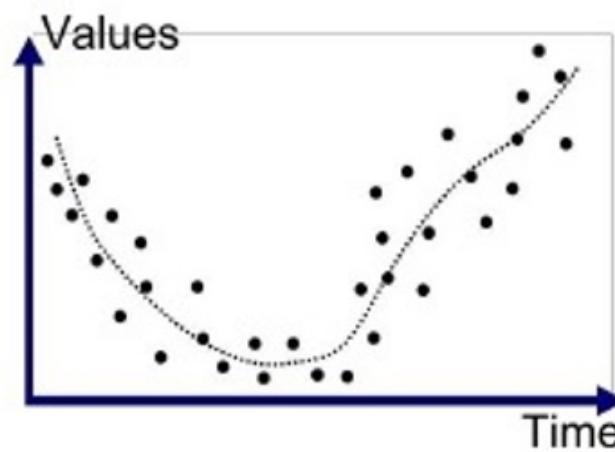


Assessing model performance

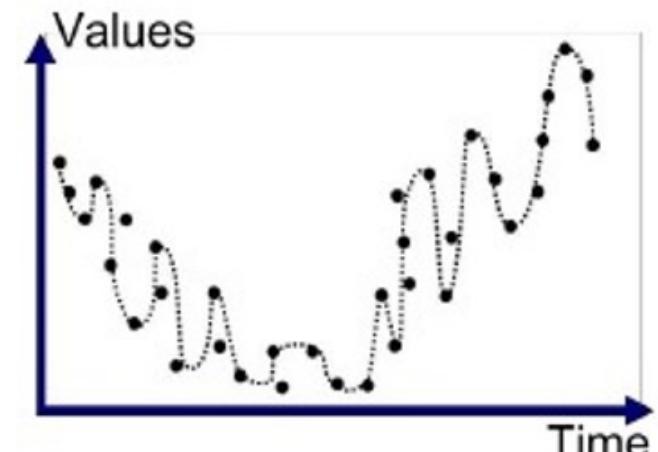
Prediction error within the data set is not a good indicator of prediction in new data



Underfitted



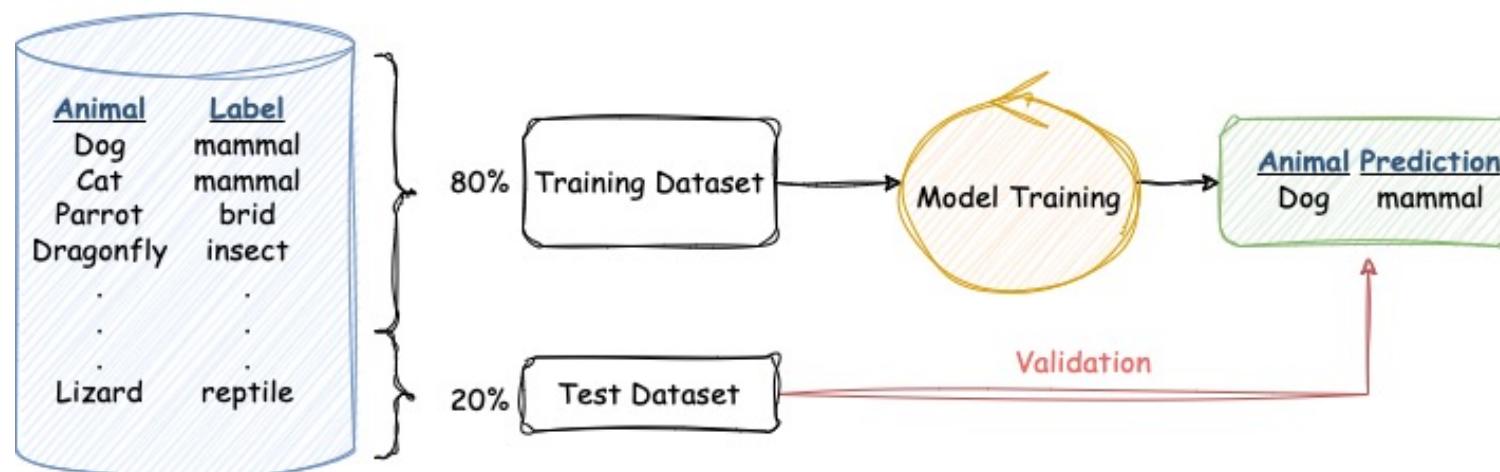
Good Fit/Robust



Overfitted

Sample splitting

Data Splitting - Random State in Machine Learning



Regularization: Keep it simple, stupid model

We can punish the model for being complex, as a way of keeping it as simple as possible.

- This is called regularization.

$$\hat{\beta} = \arg \min_{\beta} \underbrace{\sum_i (y_i - \beta' x_i)^2}_{\text{sum of squared residuals}} + \lambda \overbrace{\sum_{j=1}^p |\beta_j|^q}^{\text{complexity penalty}}$$

- λ : strength of penalty
- q : the type of regularizer ($q=1$ LASSO sets many coefficients to zero; $q=2$ Ridge regression, smooth push to zero)

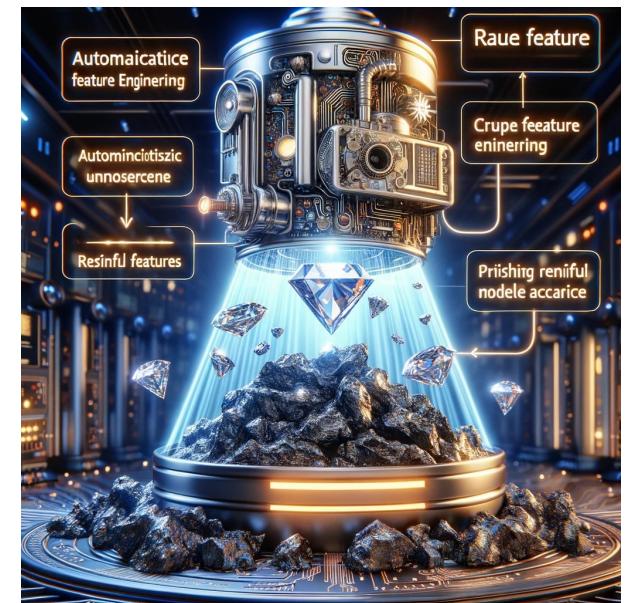
Hyperparameter search

- What λ parameter will lead to the best results?
- We can optimize the parameters themselves!



Automatic feature engineering

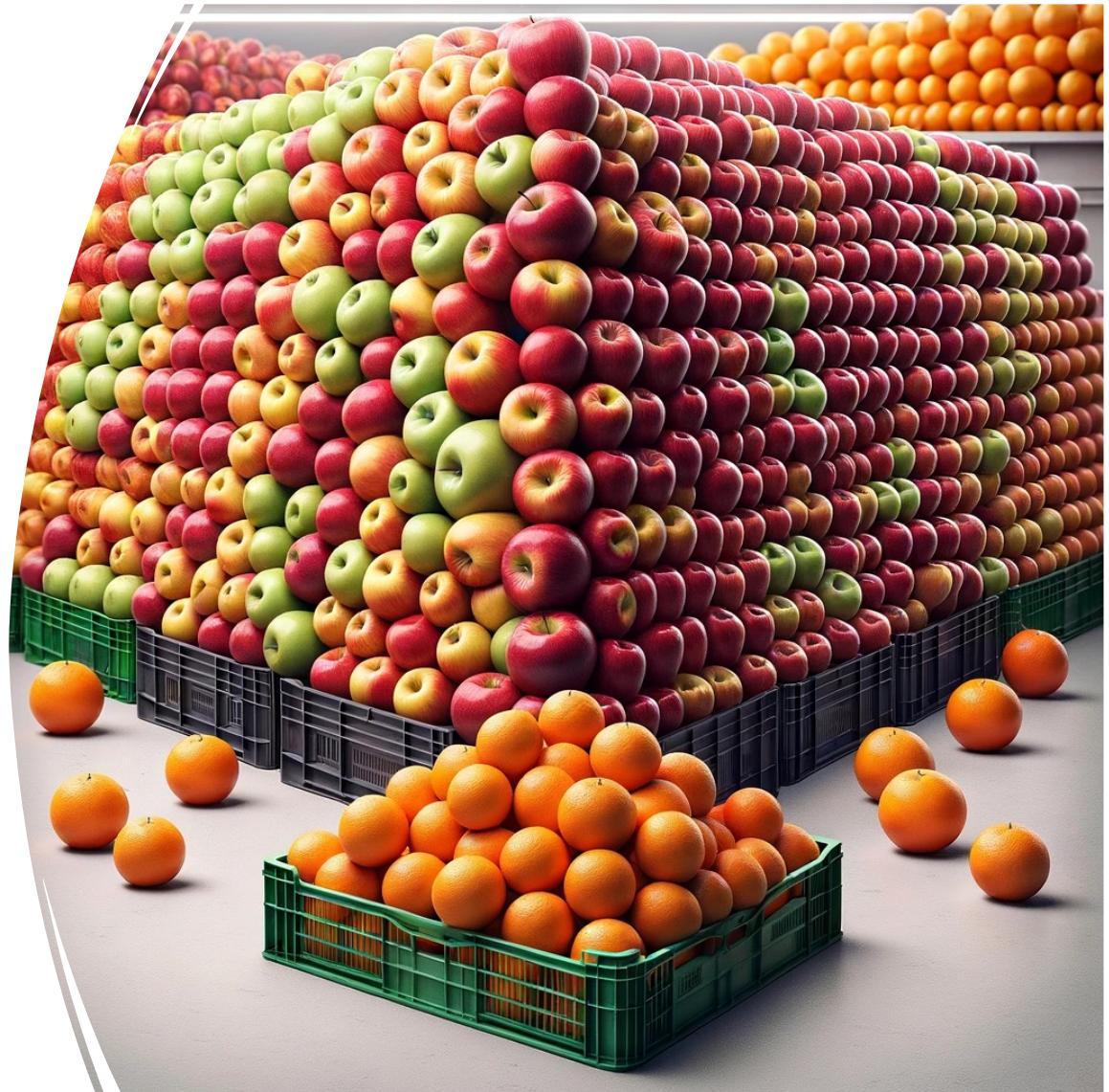
- What variables/features should be included in prediction? How to preprocess?
- How to measure of latent concept like ideology or economic growth?
- Machine learning can figure it out automatically.
- Especially important in e.g. vision or sound, where a pixel says nothing on its own.



What does this mean for theory?

Evaluating performance

- You have pictures of 500 apples and 50 oranges.
- You apply a machine learning classification model?
- Your model reaches 92% accuracy.
- Is it a good model?



Performance metrics

Accuracy: Proportion that is correct. Misleading for imbalanced datasets.

Precision: The proportion of *positive* identifications that were actually correct.

Important when the cost of false positives is high (e.g., spam detection).

Recall: The proportion of actual positives that were correctly identified.

Crucial when it's important to capture all positives (e.g., disease screening).

F1-score: The harmonic mean of precision and recall.

Macro: combines the scores from both classes.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

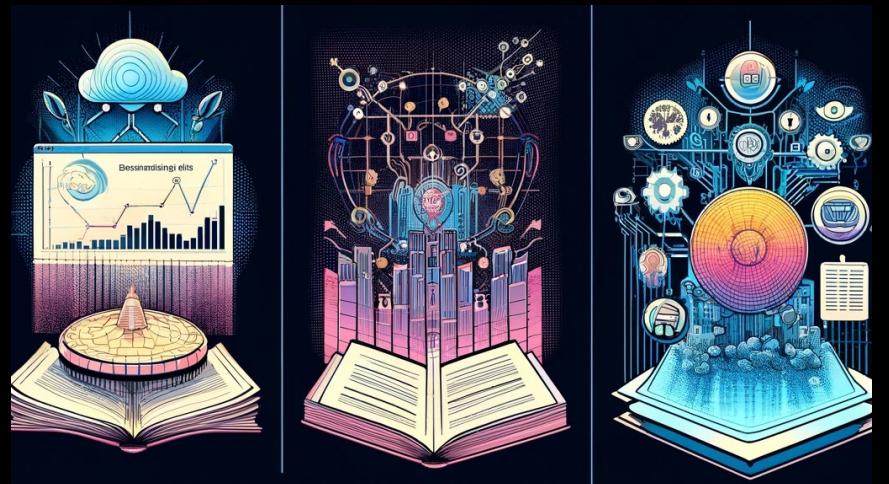
$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

	precision	recall	f1-score	support
0	0.77	0.86	0.81	37584
1	0.84	0.75	0.79	37577
accuracy			0.80	75161
macro avg	0.81	0.80	0.80	75161
weighted avg	0.81	0.80	0.80	75161

What can we actually use ML
for in the social sciences?

The epistemology of machine learning

1. From *variables* to *patterns*



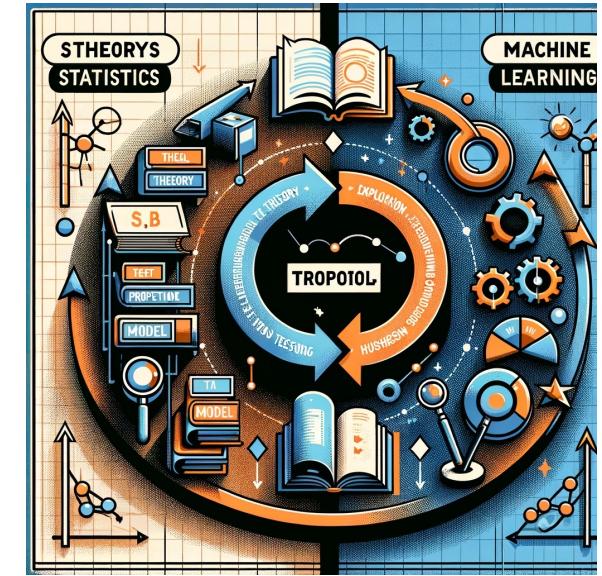
2. From *rules* to *associations*

3. From *surveys* to *sensors*

Machine Learning as a social science approach

Statistics

Deductive. Linear approach: Derive testable propositions from clear theory, test using models.



Machine learning

Abductive. Grounded theory. Iterative: explore and develop theory. Generate the research questions, concepts, and hypotheses that can later be rigorously tested with new data.

Abduction

The less famous cousin of induction and deduction.

- Deduction: What would we expect in the data based on what we know?
- Induction: What general rule can be extrapolated from specific data points?
- Abduction: Abduction starts from observing data and then seeks *the best explanation* for those observations. Called “inference to the best explanation”

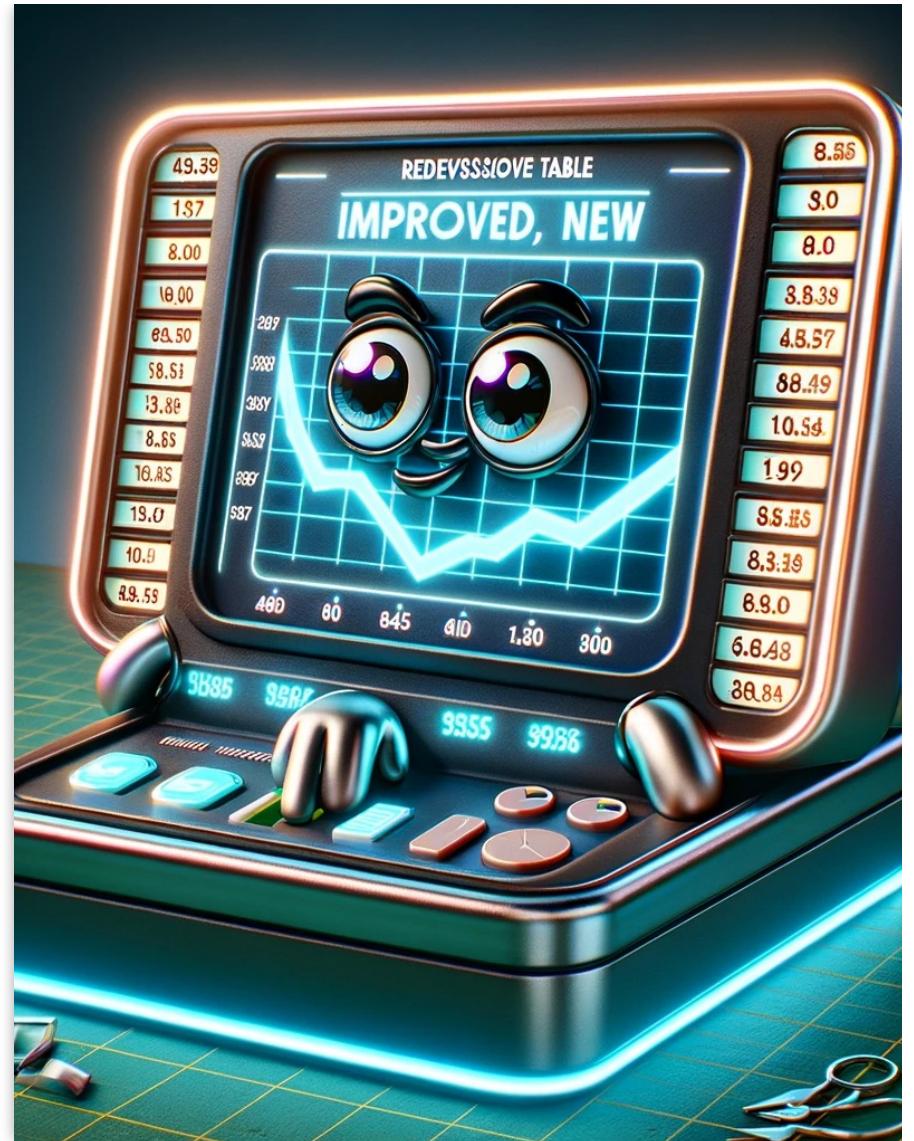
Fits machine learning:

- Machine learning models are good at identifying patterns. Fits the first step of abduction where an anomaly or an interesting pattern is identified.
- These theories are tested through further empirical research.



Improved regression models

- We can use ML as a form of statistics on steroids, with more complex modeling capabilities, handling non-linear relationships, and managing high-dimensional data
- **Interpretable models:** Some machine learning models, such as linear regression with regularization (like Lasso and Ridge), still allow for interpretation. (Yay, regression tables!)
- There are also ways to make sophisticated models interpretable. Bootstrapping can be used to estimate confidence intervals and standard errors.
- However, might not always provide exact p-values or the traditional statistics familiar from OLS regression.



Coding and classification

One of the most widely used tools in the social sciences is manually classifying observations into a set of categories that are determined before the analysis begins.

The models allow measuring virtually any latent aspect of data!



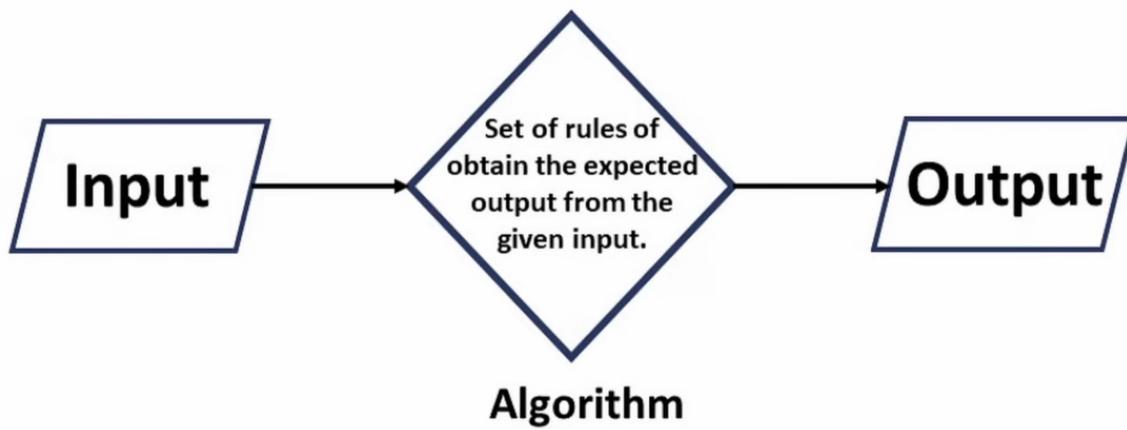
Causal inference from observational data

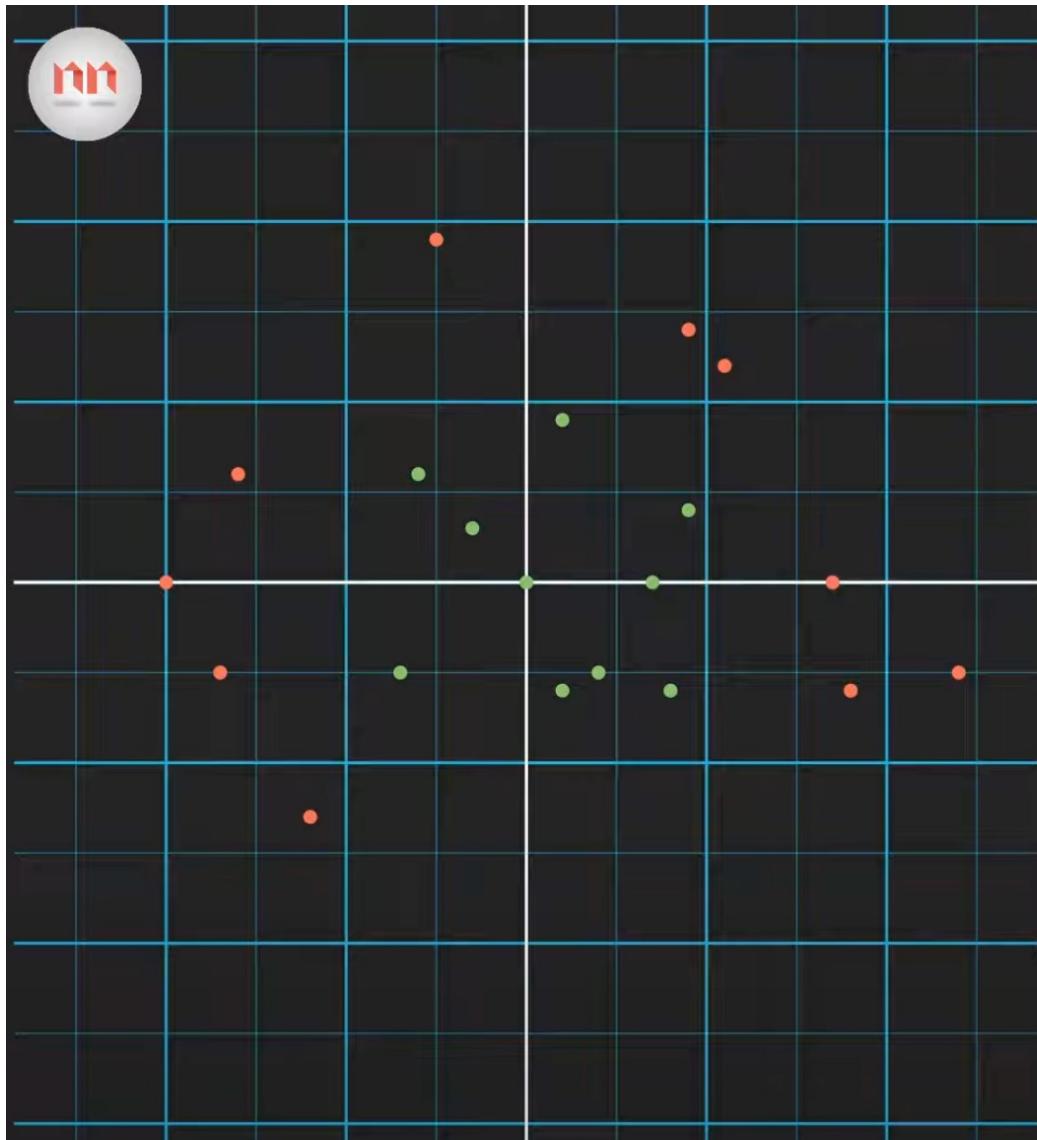
- Imagine we're interested in studying the effect of attending a private (treatment) vs a public high school (control) on college admission rates (outcome).
- Obviously, the problem is confounding factors. (E.g., having rich parents.)
- We calculate *the propensity score*: the probability of a student attending a private high school given their observed characteristics (covariates). We estimate this probability using statistical or machine learning models
- We match students based on their propensity scores.
- The effect of treatment is then = the average difference between the matched individuals



How do the models actually work?

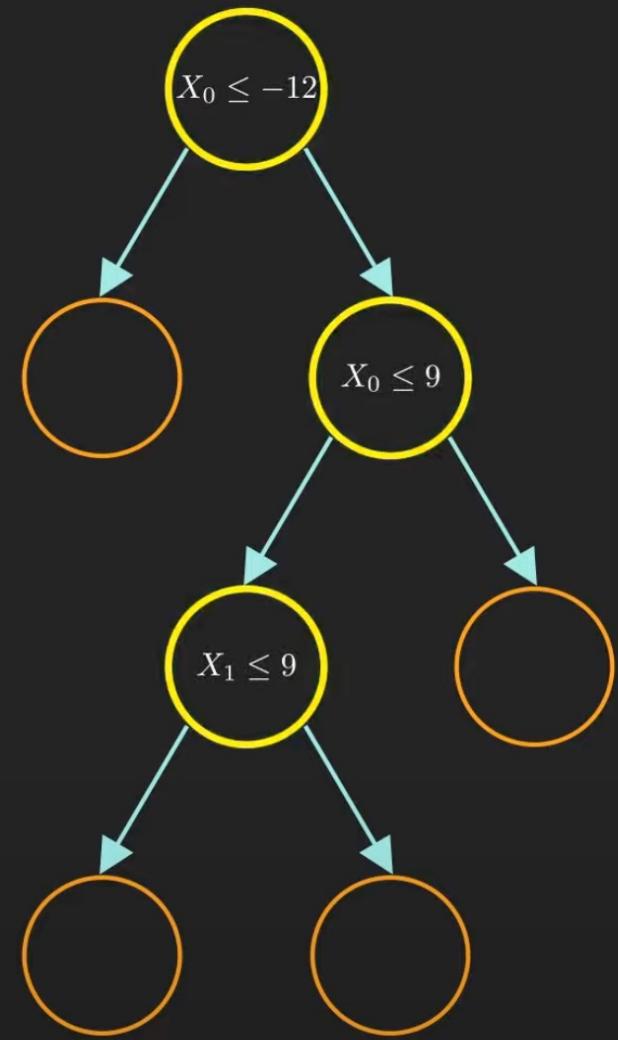
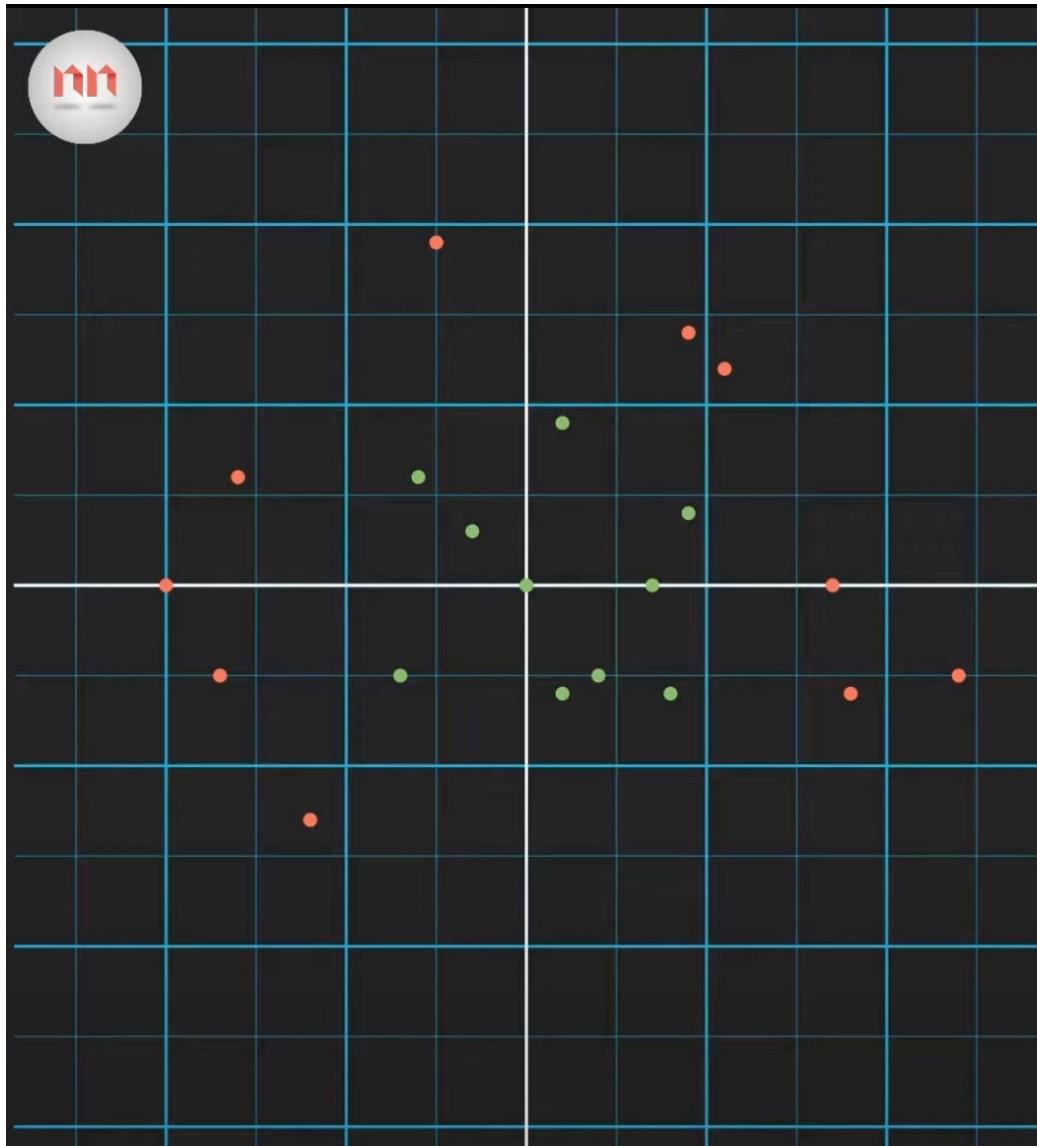
Let's get nerdy!

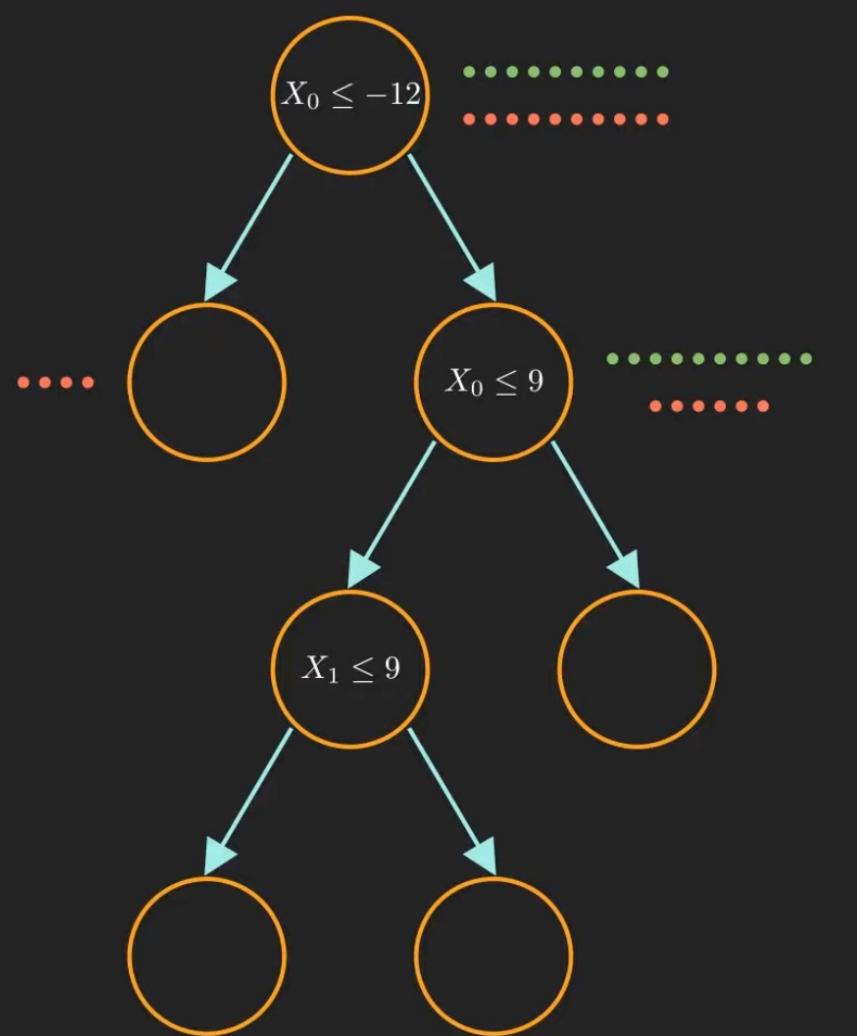
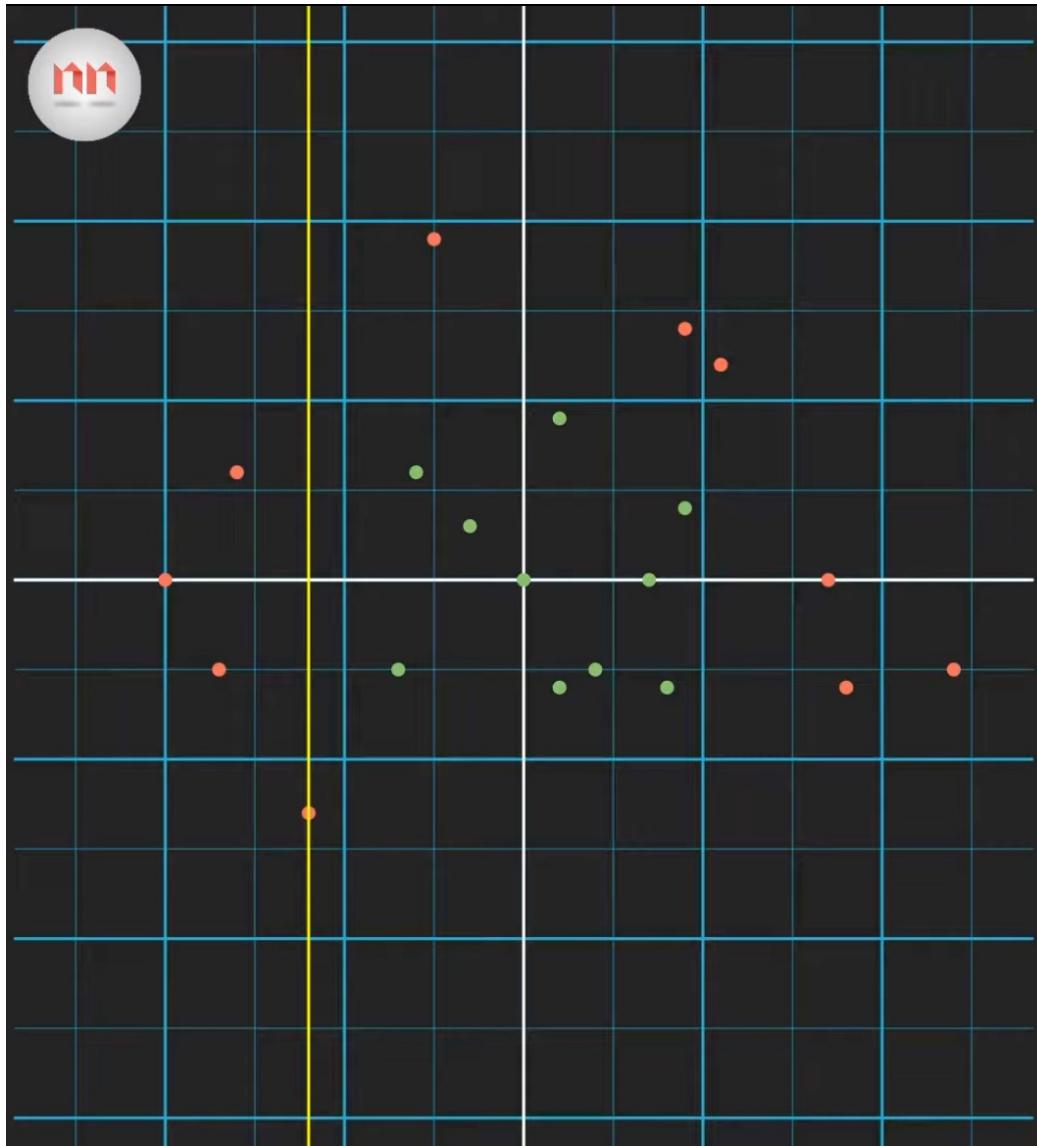


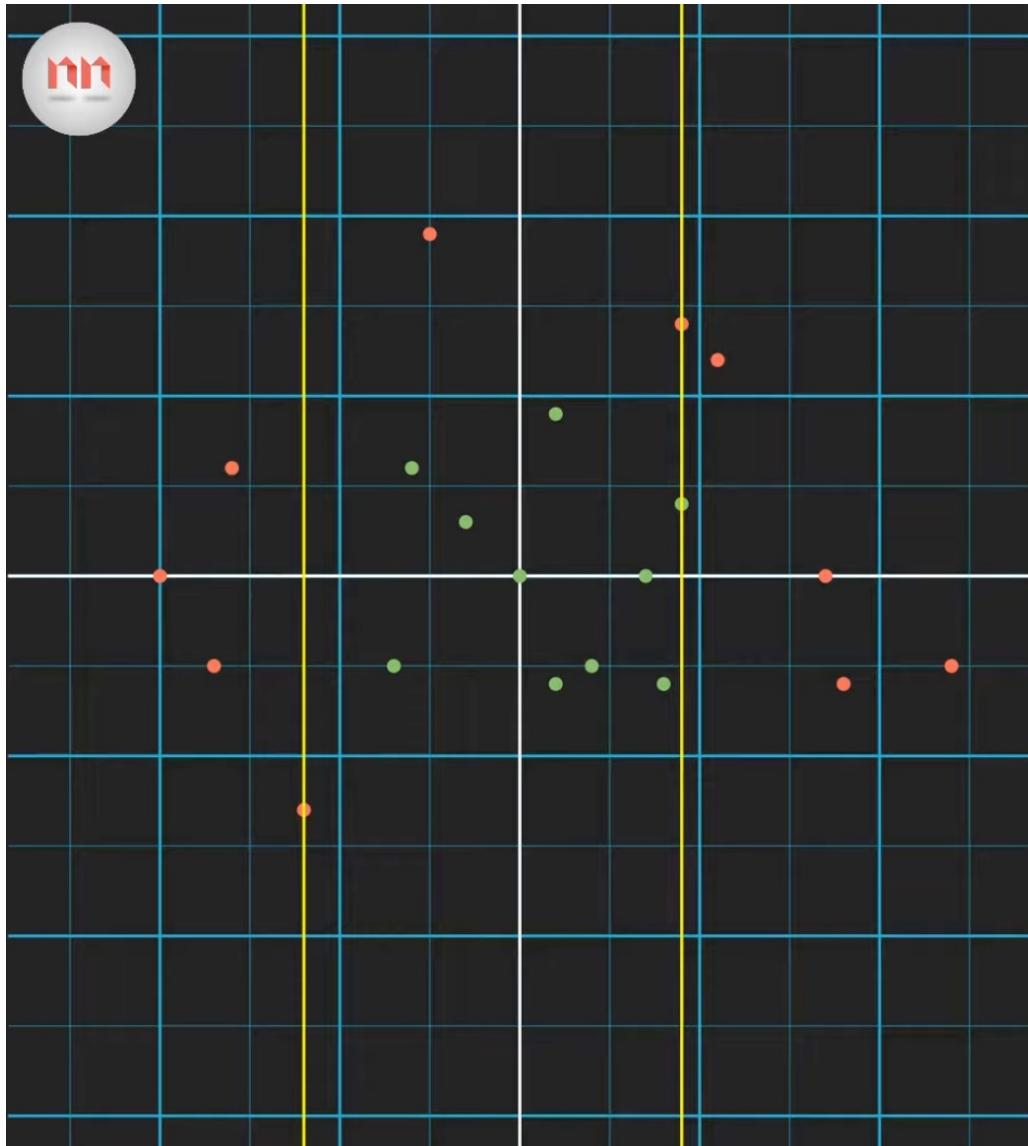


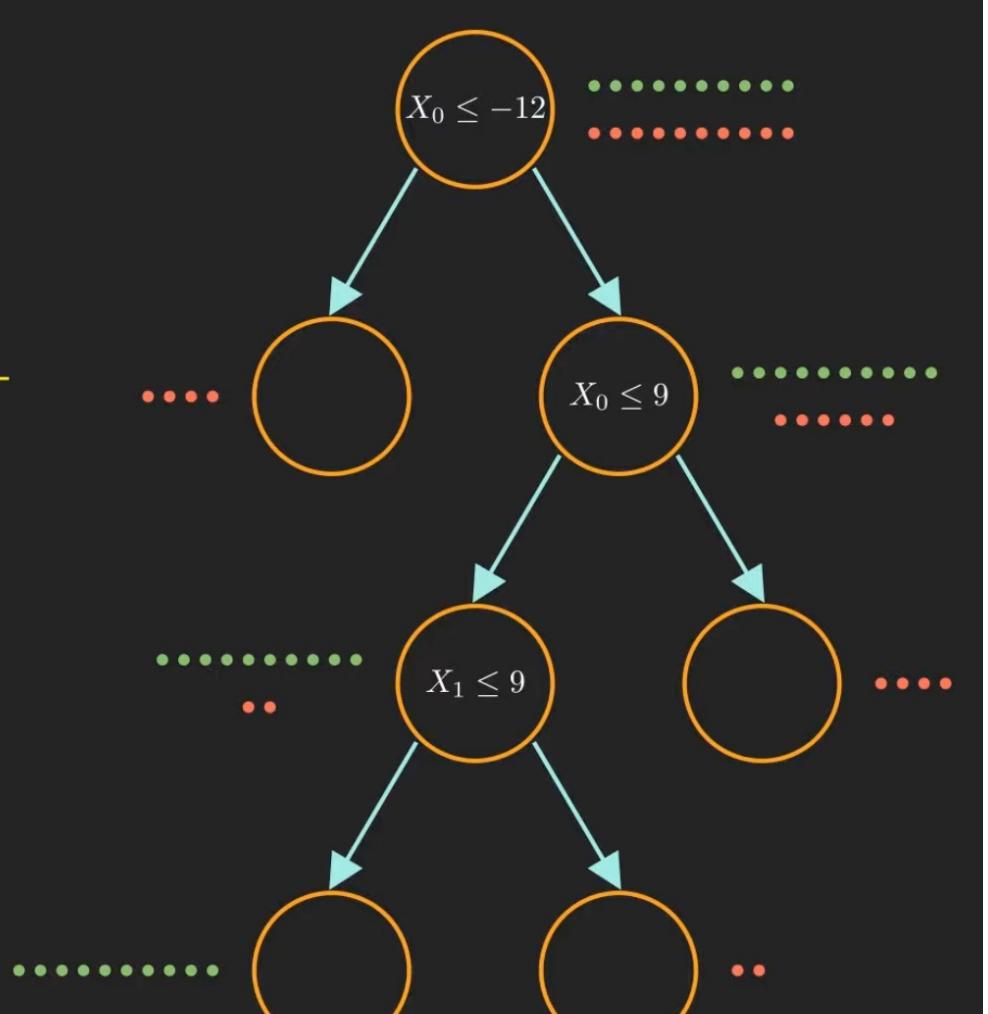
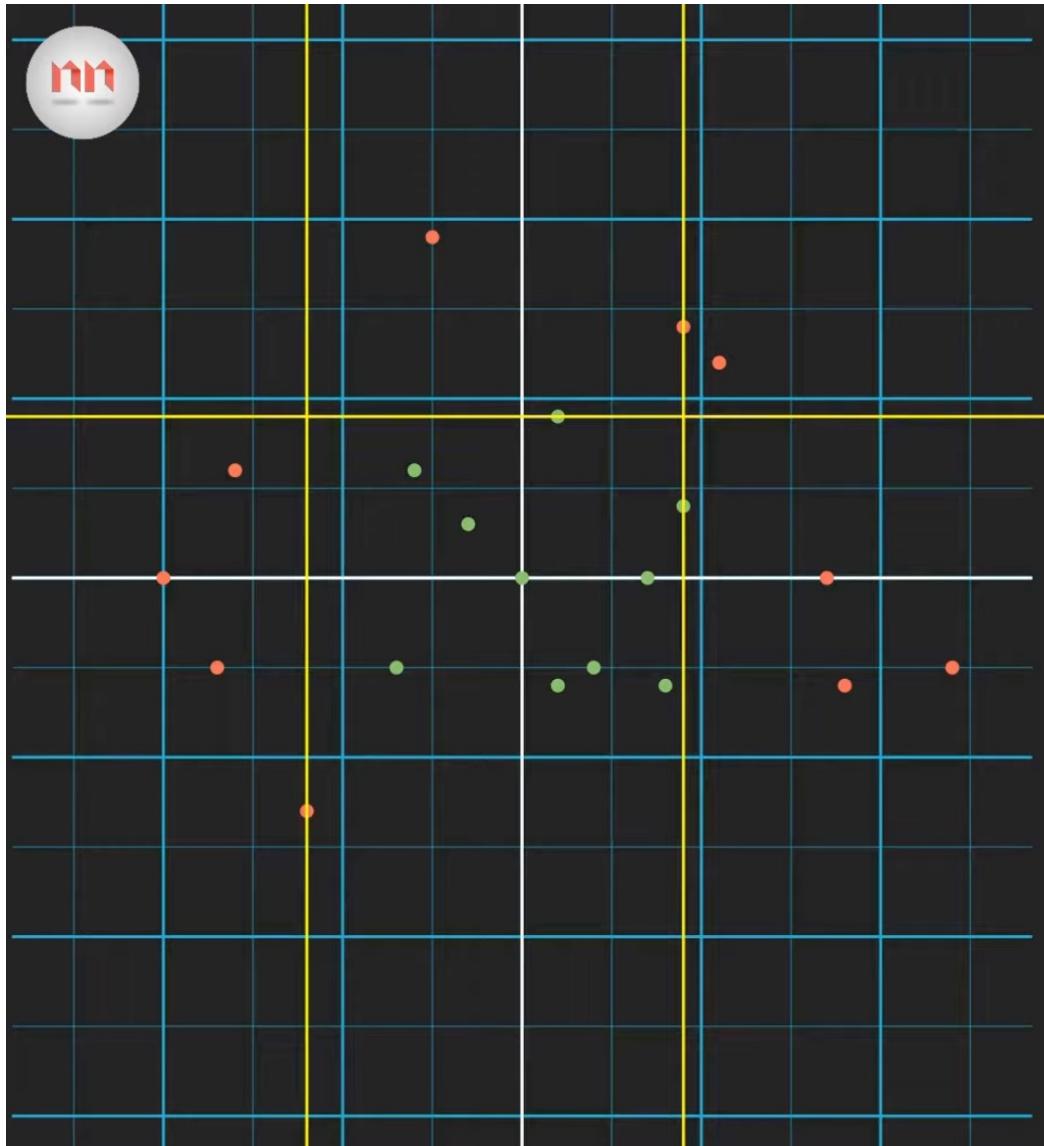
Decision trees

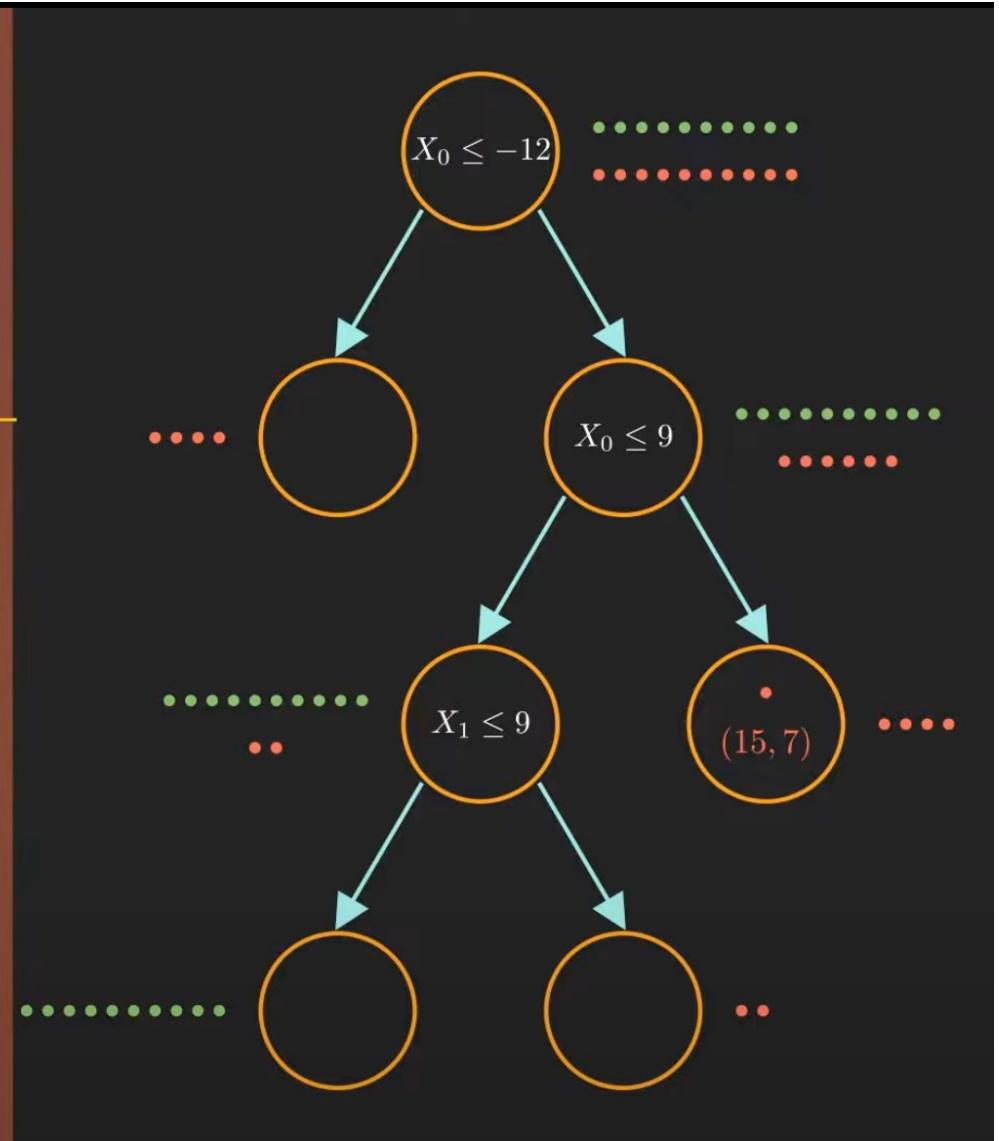
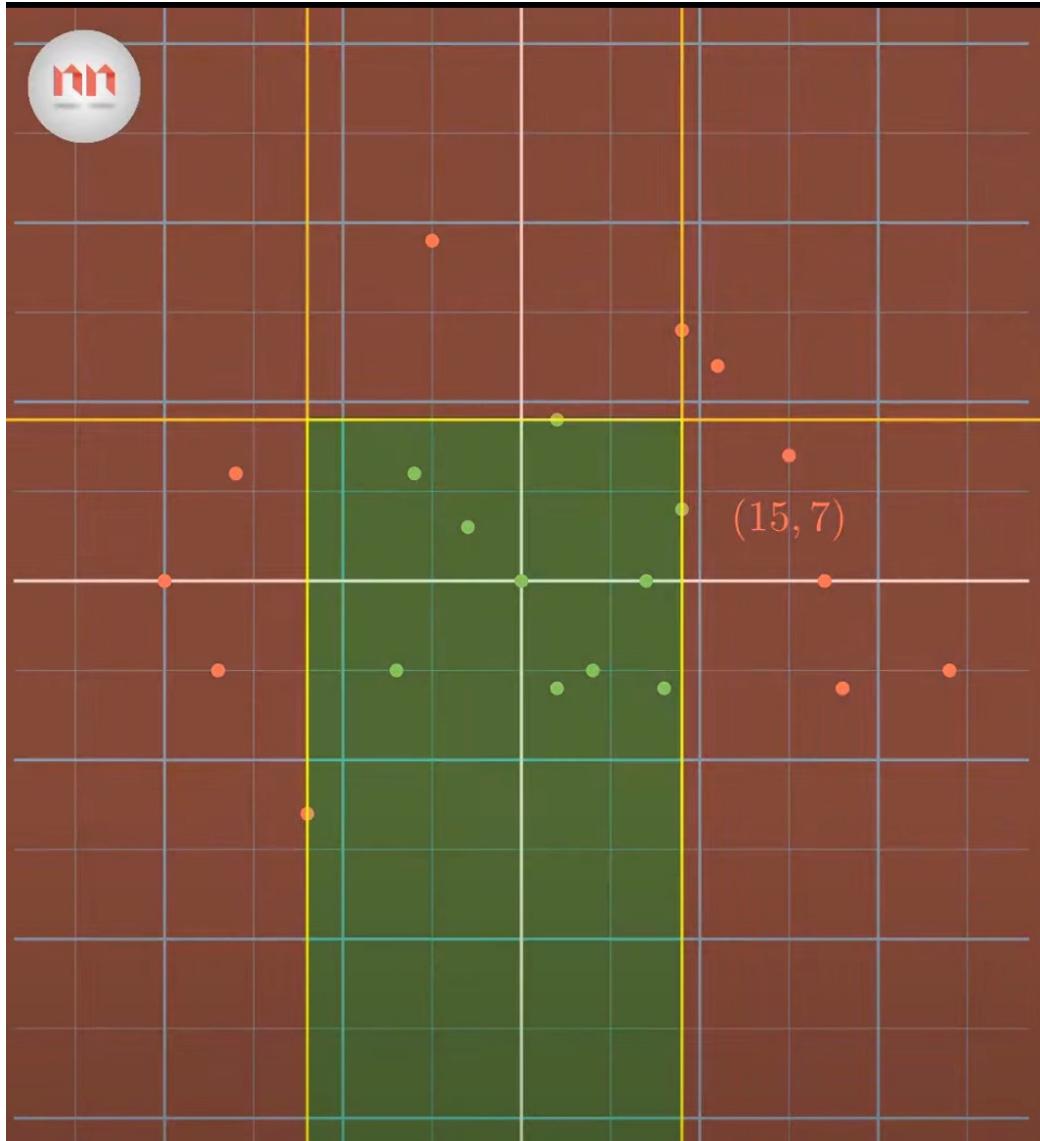
- Among the simplest models
- We want to separate the two classes of dots
- Not possible with a straight line











Training decision trees

- In each step, the model uses the split that maximizes the information gain – minimizes the *entropy*
 - How can I draw the line to *most separate* the two classes?
 - Compares all possible splits.
-
- Downside with Decision Trees: they are highly sensitive to training data. They can fail to generalize.

id	x_0	x_1	x_2	x_3	x_4	y
0	4.3	4.9	4.1	4.7	5.5	0
1	3.9	6.1	5.9	5.5	5.9	0
2	2.7	4.8	4.1	5.0	5.6	0
3	6.6	4.4	4.5	3.9	5.9	1
4	6.5	2.9	4.7	4.6	6.1	1
5	2.7	6.7	4.2	5.3	4.8	1

id
2
0
2
4
5
5

id
2
1
3
1
4
4

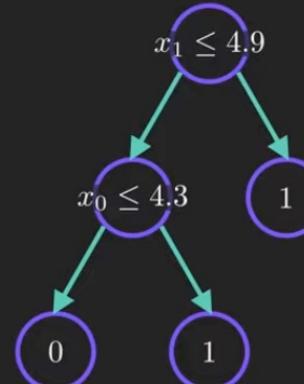
id
4
1
3
0
0
2

id
3
3
2
5
1
2

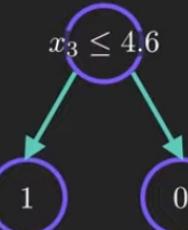
Random Forest

- A forest is several trees.
- Each tree is trained on a random subset of the training data and makes predictions independently.
- The final prediction is obtained by *averaging* or *voting* the predictions.

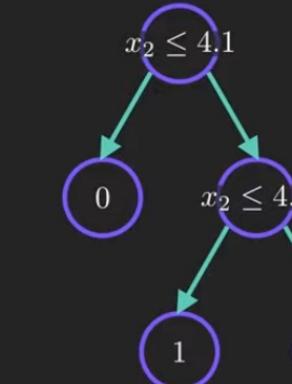
x_0, x_1



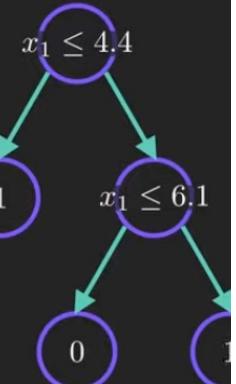
x_2, x_3



x_2, x_4



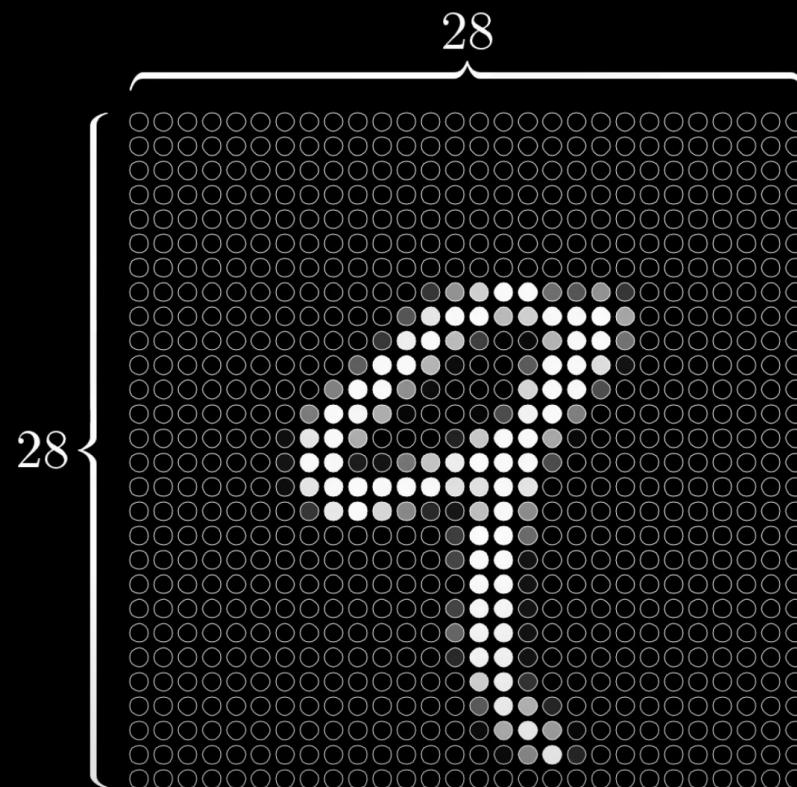
x_1, x_3



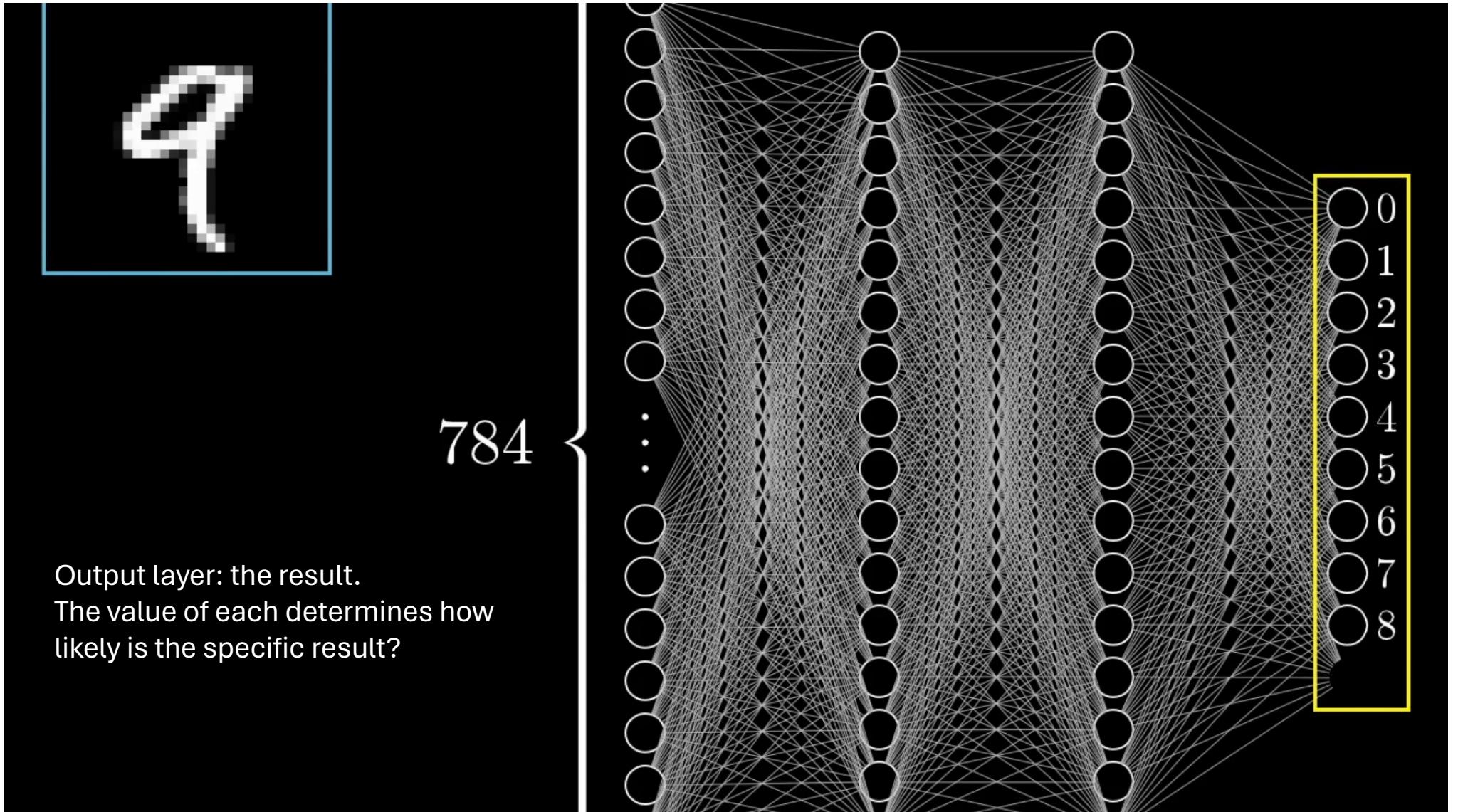
Artificial Neural Networks

Task: identify the handwritten number

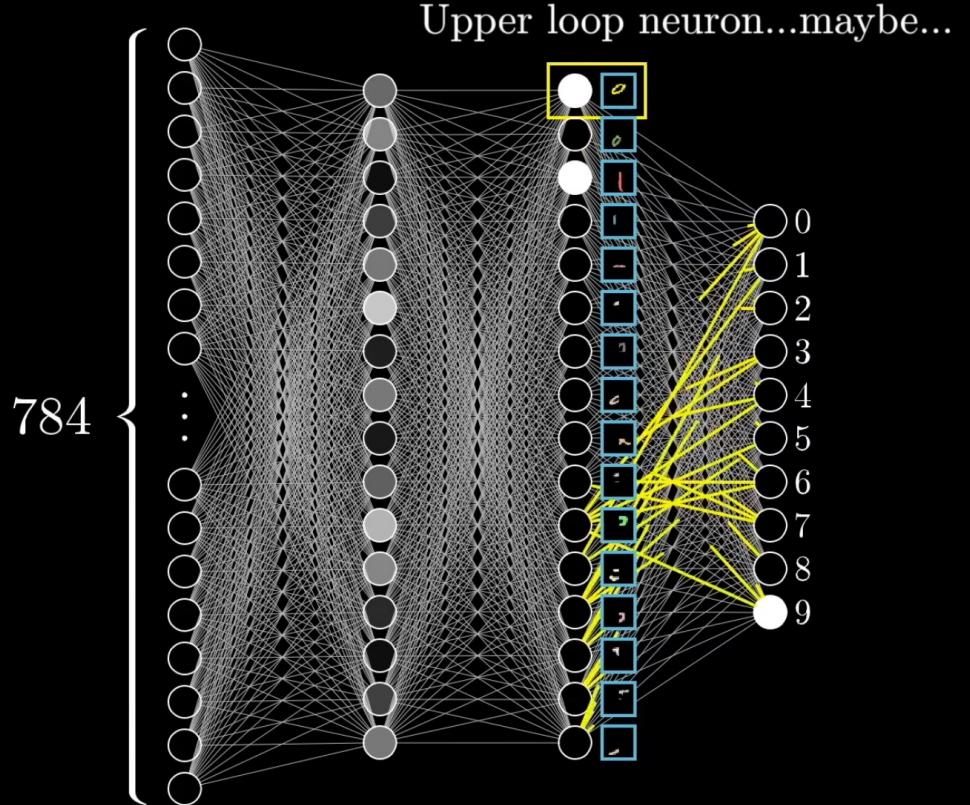
Greyness is activation:
from 0 to 1.



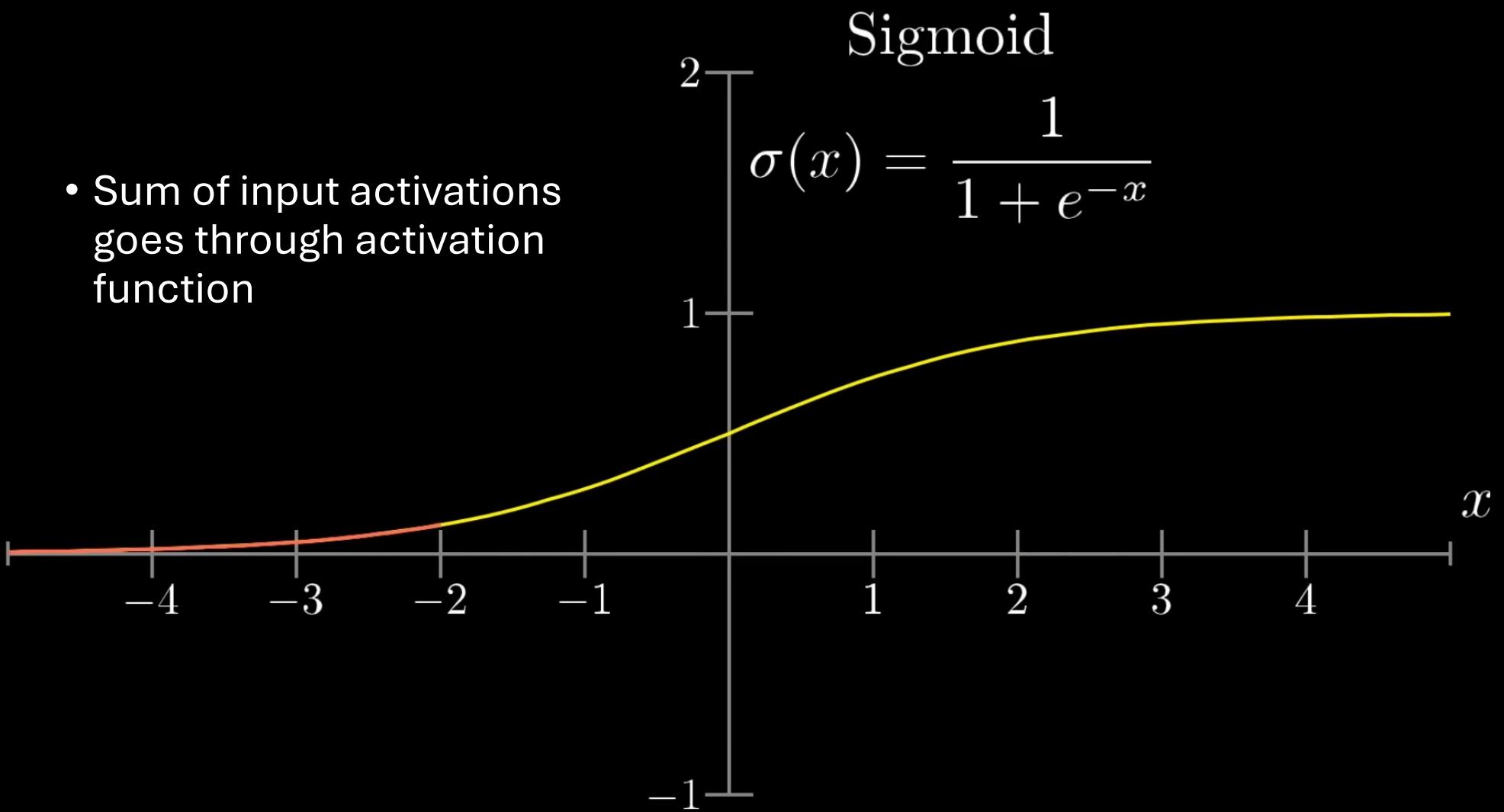
$$28 \times 28 = 784$$



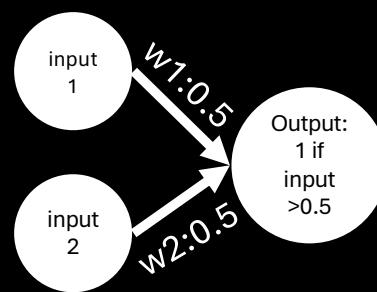
- Activation of each layer depends on the previous layer
- The neurons are connected with *weights*
- Each layer learns ‘features’ - subcomponents

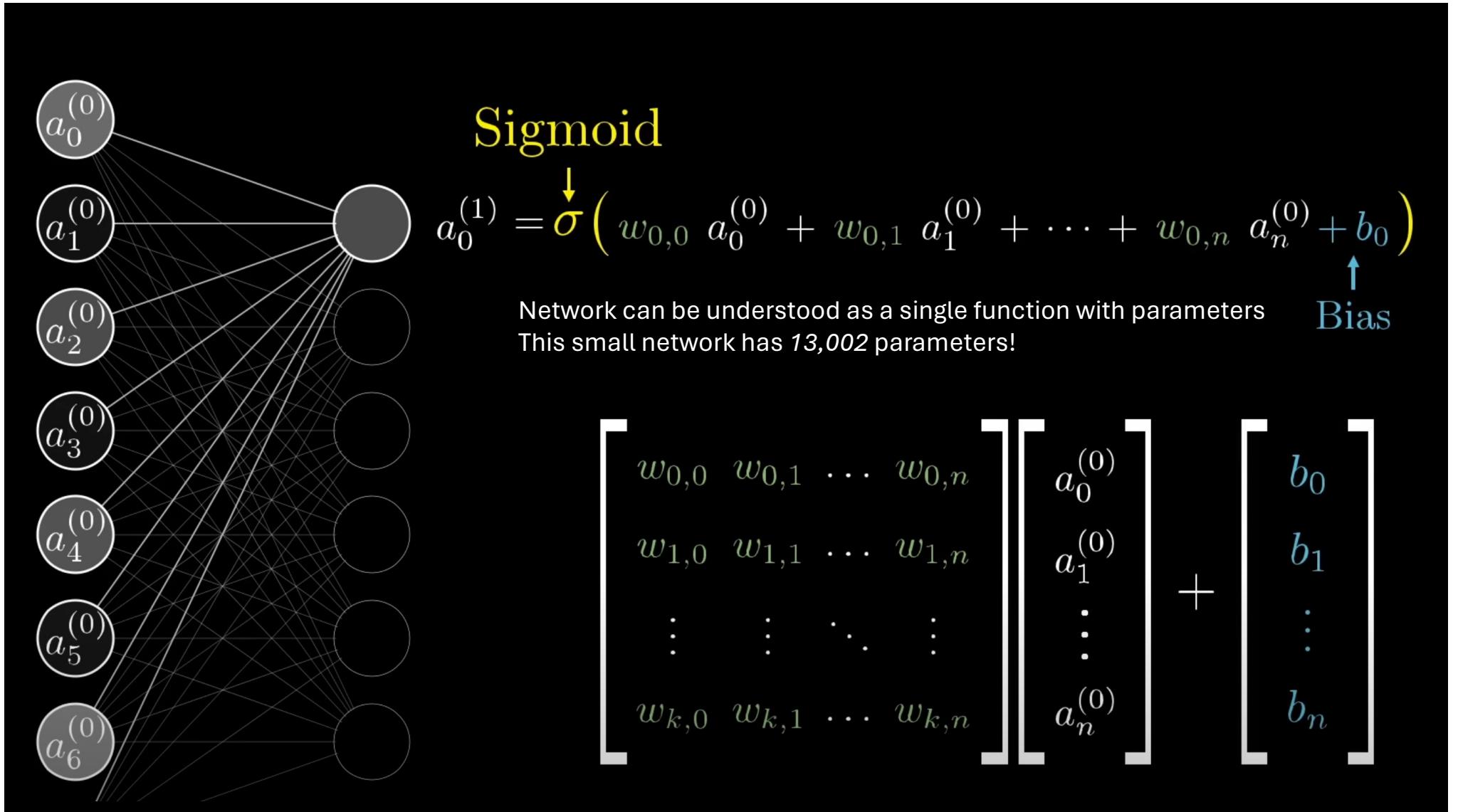


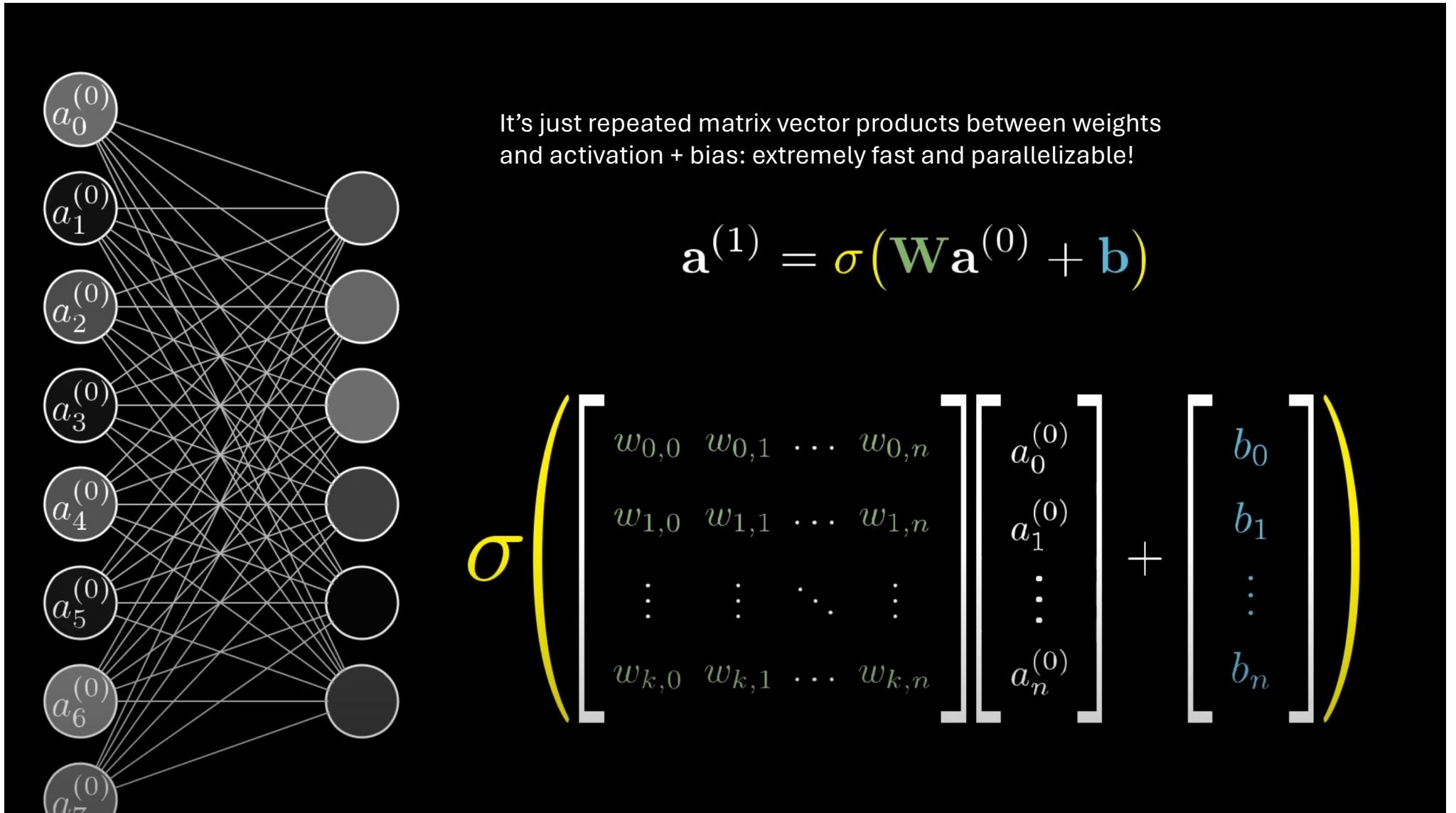
- Sum of input activations goes through activation function



What logical function is implemented by this ANN?

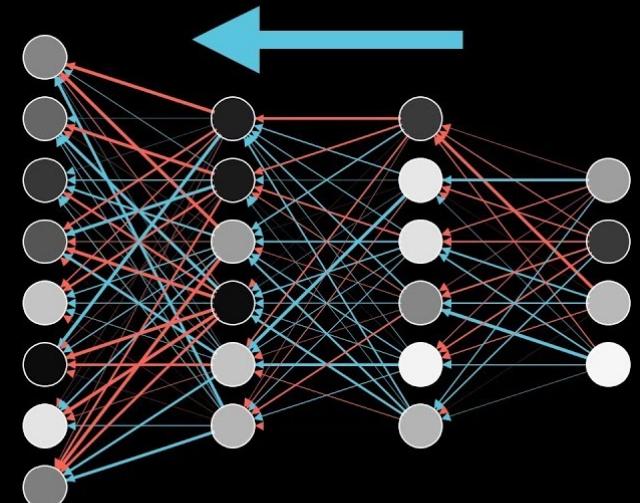






Training: Backpropagation

- The error is fed backwards in the network from the output.
- Error = difference between your answer and the correct answer.
- The error moves backwards in the network, adjusting weights by *gradient descent*
- How should I change this weight to most reduce the error?





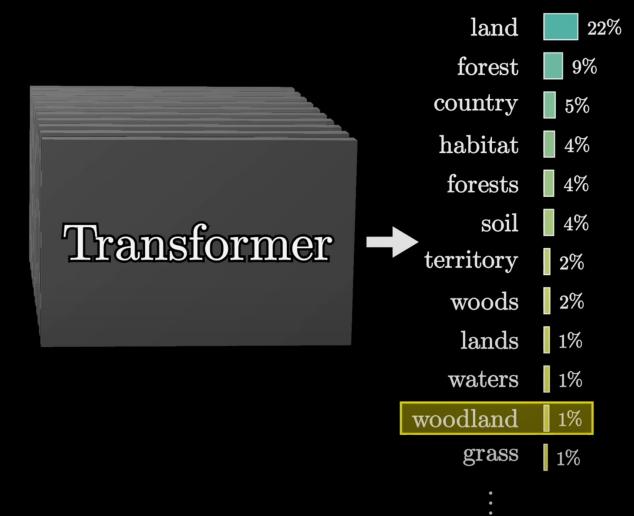
Transformers: fancy neural networks for text

Transformers

(as in GPT – Generative Pretrained *Transformers*)

- Transformers are a particular neural network structure.
- The task of predicting the next word given a word sequence.
- The easiest way of predicting the next word is to build an understanding of the world and how it works.

Behold, a wild pi creature,
foraging in its native
woodland _____





To date, the cleverest thinker of all time was

???

Tokenization

To date, the cleverest thinker of all time was

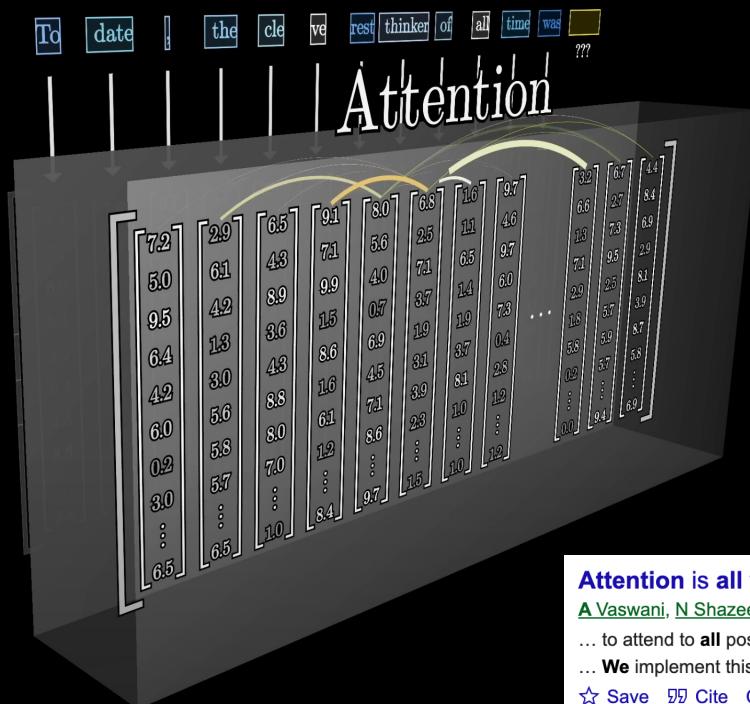
???

Word embeddings

To	date	,	the	cle	ve	rest	thinker	of	all	time	was	???
5.4	7.8	9.7	2.6	3.6	5.6	1.6	9.7	3.2	6.7	4.4		
7.1	5.2	7.9	7.7	4.3	4.3	1.1	4.6	6.6	2.7	8.4		
6.0	5.6	4.6	4.5	6.9	9.8	6.5	9.7	1.3	7.3	6.9		
5.4	9.2	7.7	5.6	0.6	1.0	1.4	6.0	7.1	9.5	2.9		
4.2	0.7	1.2	0.2	6.6	2.1	1.9	7.3	2.9	2.5	8.1		
6.4	0.9	6.3	6.1	6.6	1.6	3.7	0.4	1.8	5.7	3.9		
4.3	0.2	1.4	6.1	2.1	6.5	8.1	2.8	5.8	5.9	8.7		
8.8	8.2	9.4	6.1	1.3	2.5	1.0	1.2	0.2	5.7	5.8		
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮		
3.8	8.6	4.1	6.8	3.6	2.4	1.0	1.2	0.0	9.4	6.9		

Attention-block

The words talk to each others and give meaning to one another.
Which words change the meaning of each other, and how?

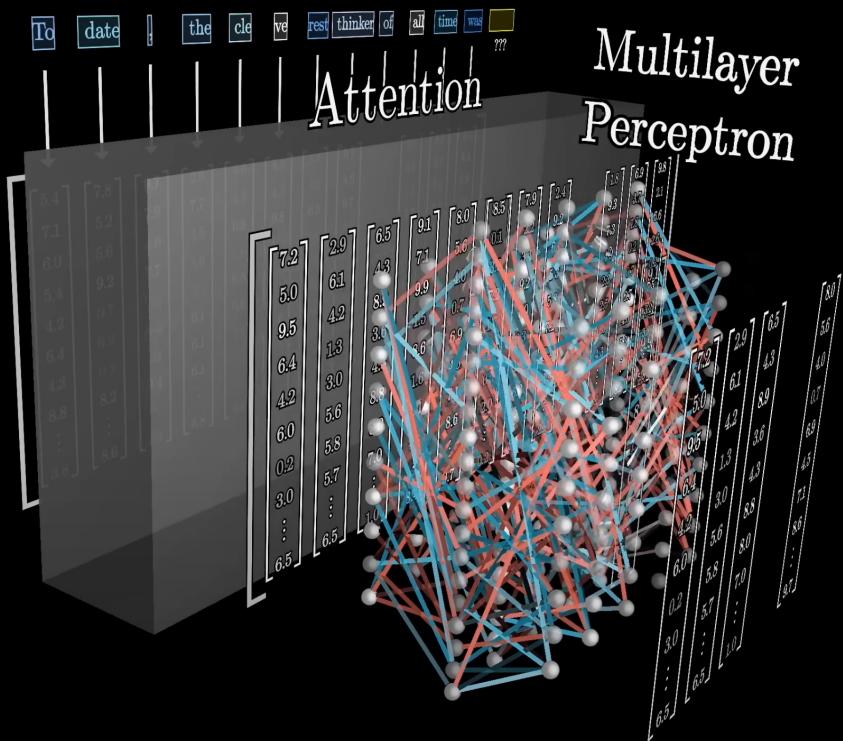


Attention is all you need

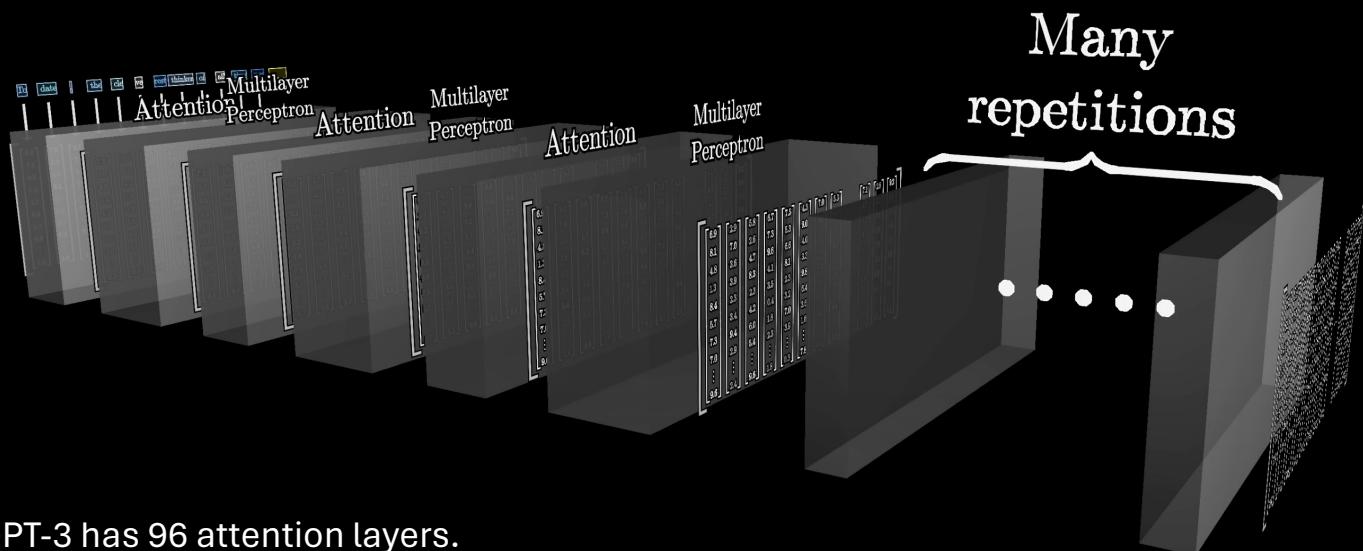
[A Vaswani, N Shazeer, N Parmar... - Advances in neural ...](#), 2017 - proceedings.neurips.cc
... to attend to **all** positions in the decoder up to and including that position. **We need** to prevent
... **We** implement this inside of scaled dot-product **attention** by masking out (setting to $-\infty$) ...
☆ Save ⚡ Cite Cited by 119393 Related articles All 87 versions ☰

Multi-layer perceptron / Feed-forward layer

A normal neural network: does the actual processing

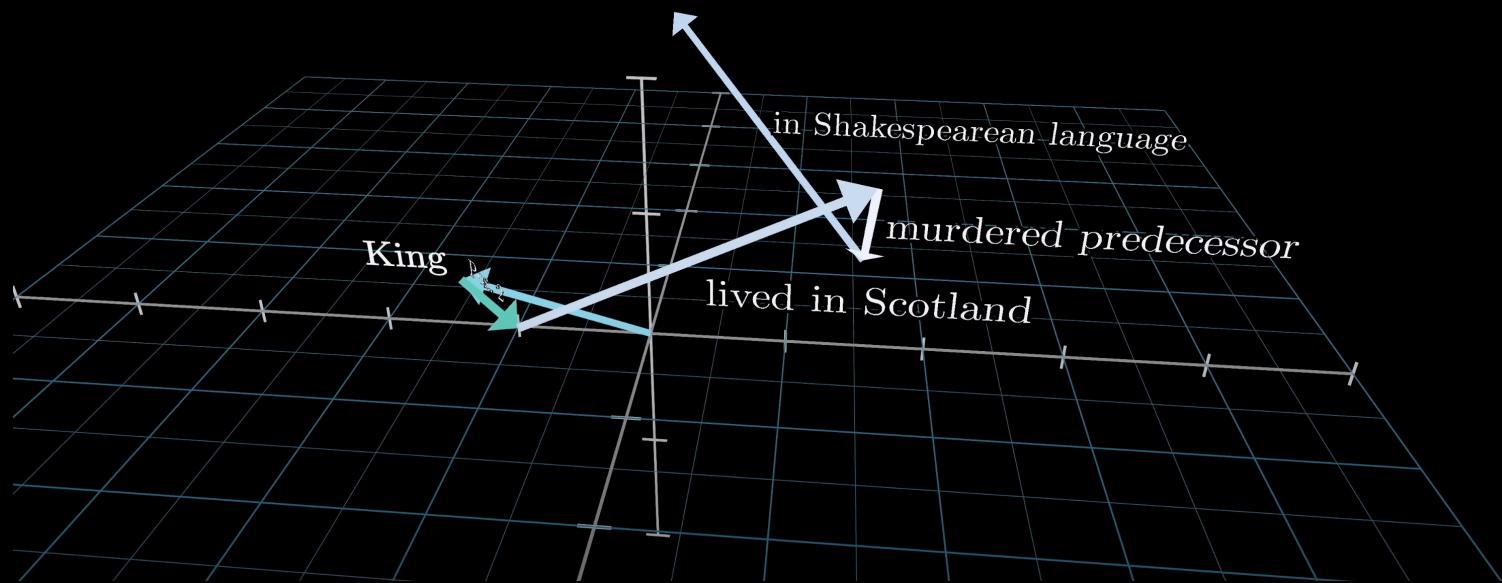


Rinse and repeat: attention x multilayer

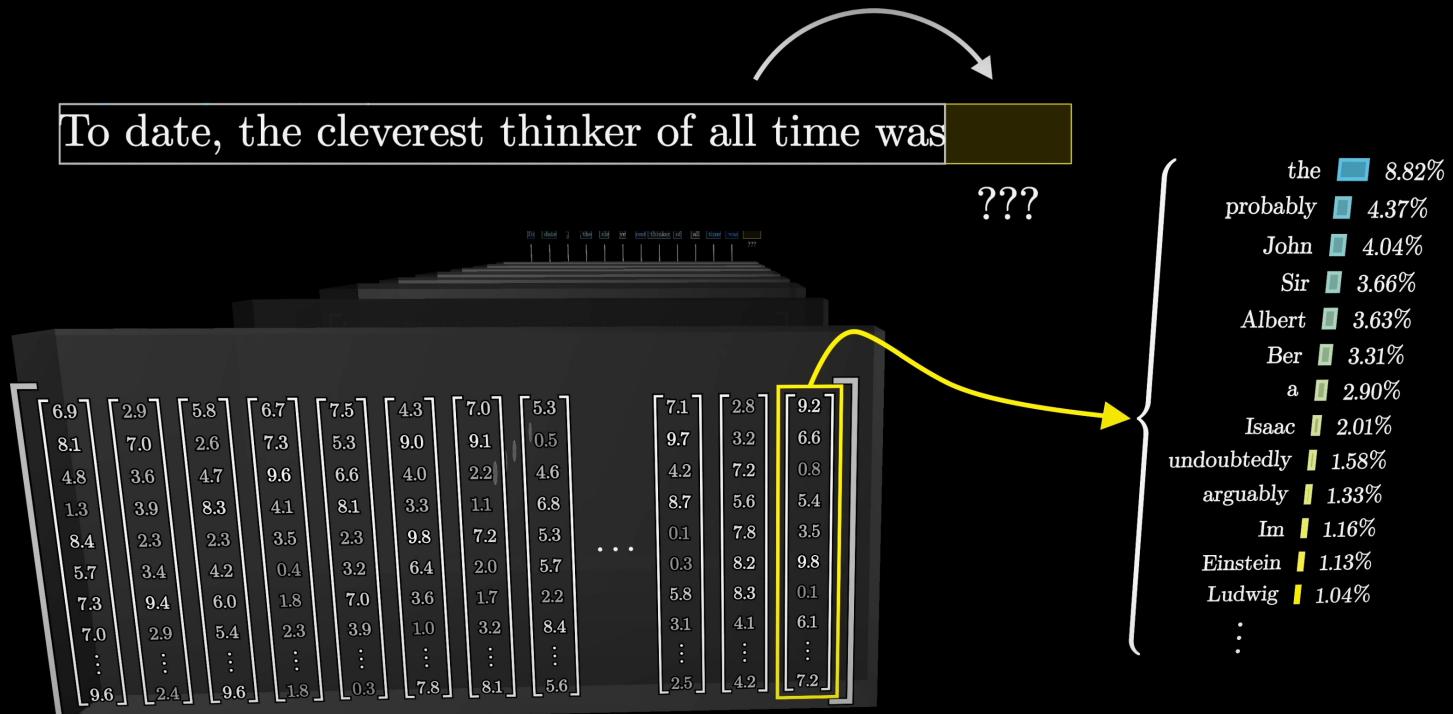


- ChatGPT-3 has 96 attention layers.
- Instead of just a word embedding, it builds a much richer representation of contextual meaning

The King doth wake tonight and takes his rouse ...

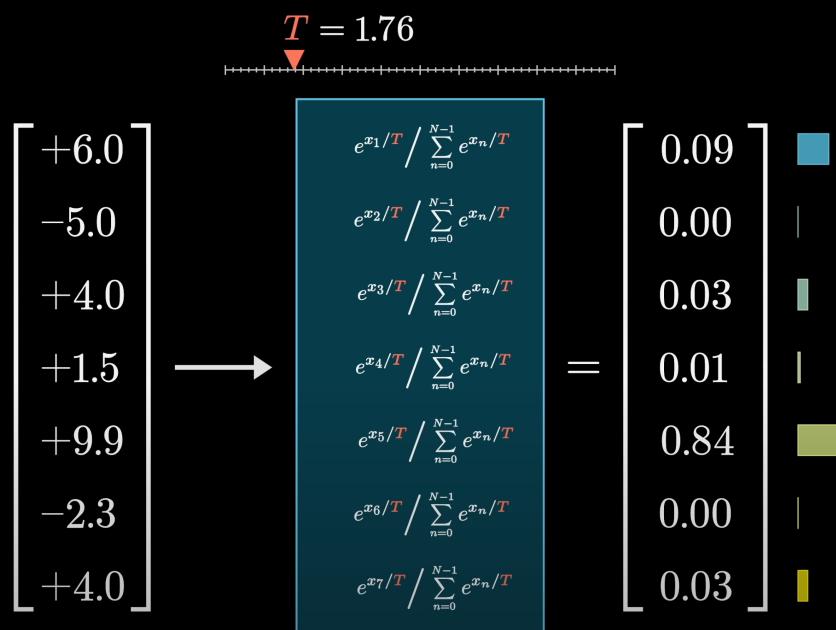


At the end, the last vector is a probability distribution over words
We sample over this distribution



The sampling is done with softmax

softmax with temperature



A standard way of turning numbers into a probability distribution

Temperature controls evenness of the distribution

$$e^{x_1/T} / \sum_{n=0}^{N-1} e^{x_n/T}$$



BERT

(Bidirectional Encoder
Representations from Transformers)

- Pre-trained to understand language and context: trained both to guess next word, and to guess a masked word.
- Gives it an understanding of language and context
- We can finetune it to problems we want to solve using additional layers that are trained from scratch

This is the crème of supervised text models!