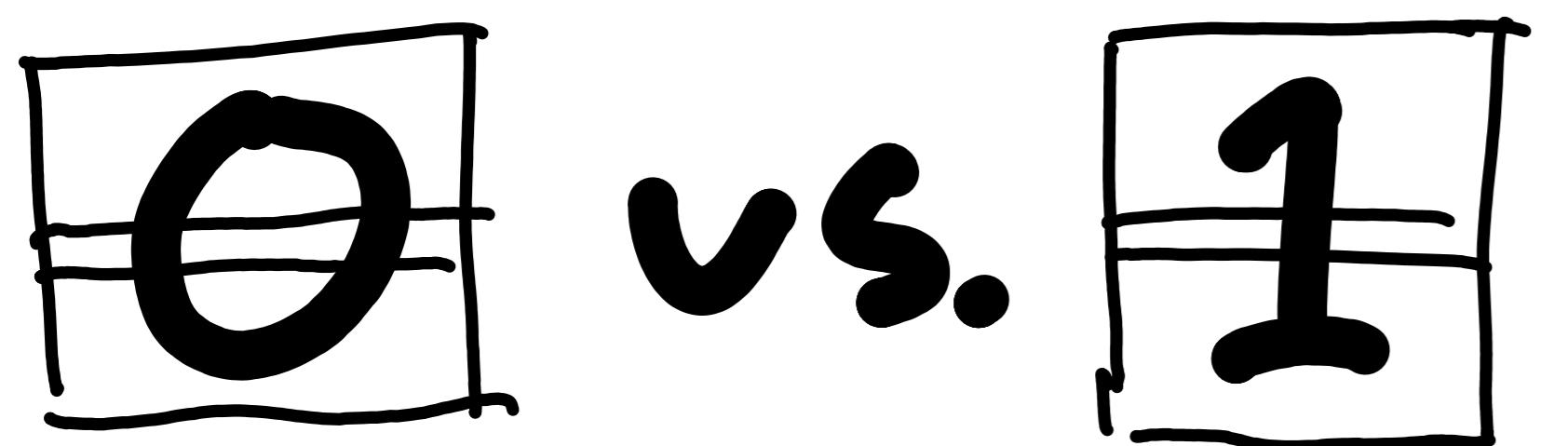


CSC2516

Neural Networks + Deep learning

Fall 2023



$n_dark = \emptyset$

for pixel in row:

 if pixel is black:

$n_dark += 1$

 if $n_dark > 1$:

 return \emptyset

else:

 return 1

0 1 2 3 4 5 6 7 8 9

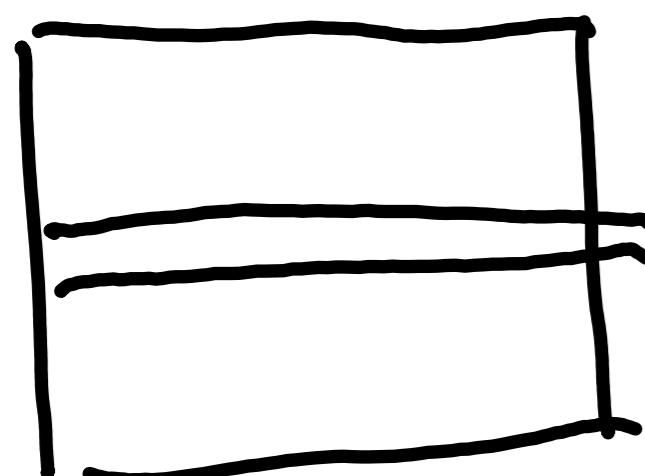
Machine Learning (ML):
program whose behavior is
"learned" from examples.

Old-fashioned ML:

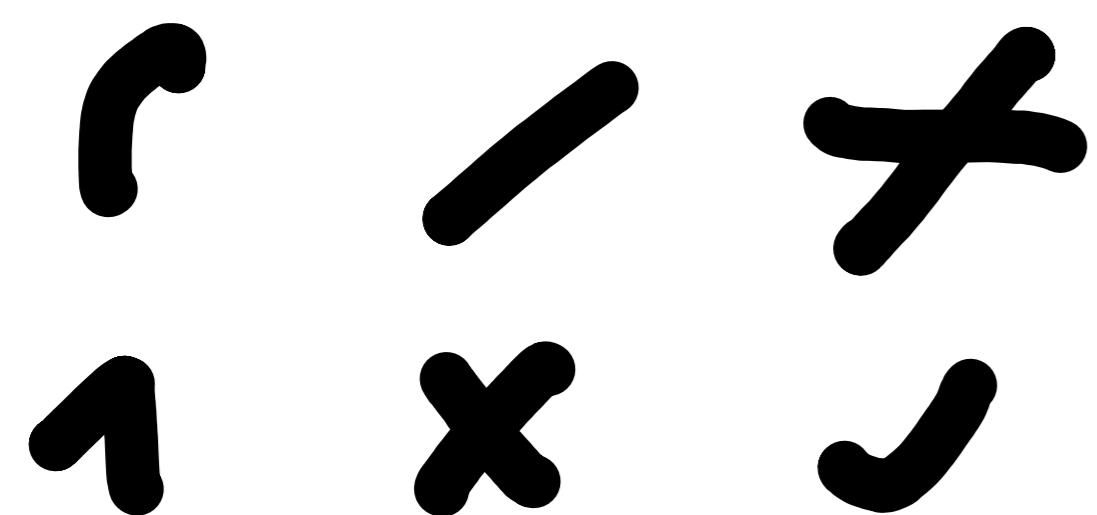
1. Design features
2. Train a "simple" model

0 0 1 | 8 8 7 7

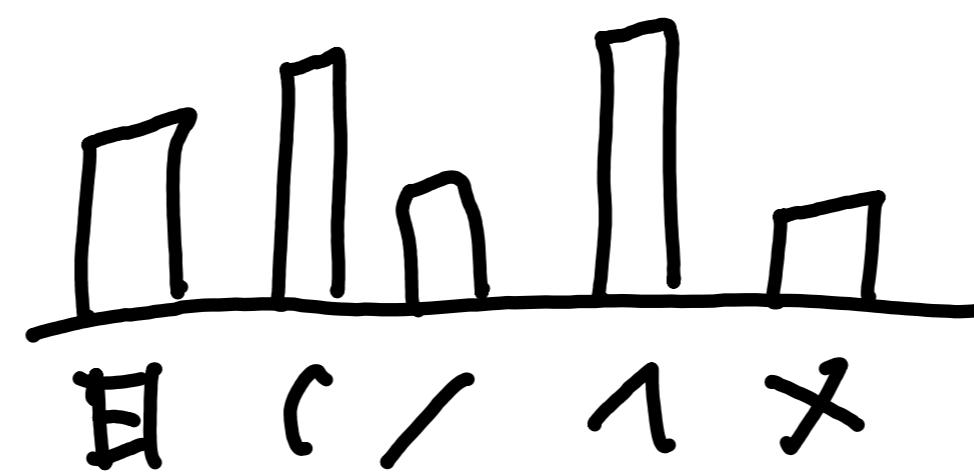
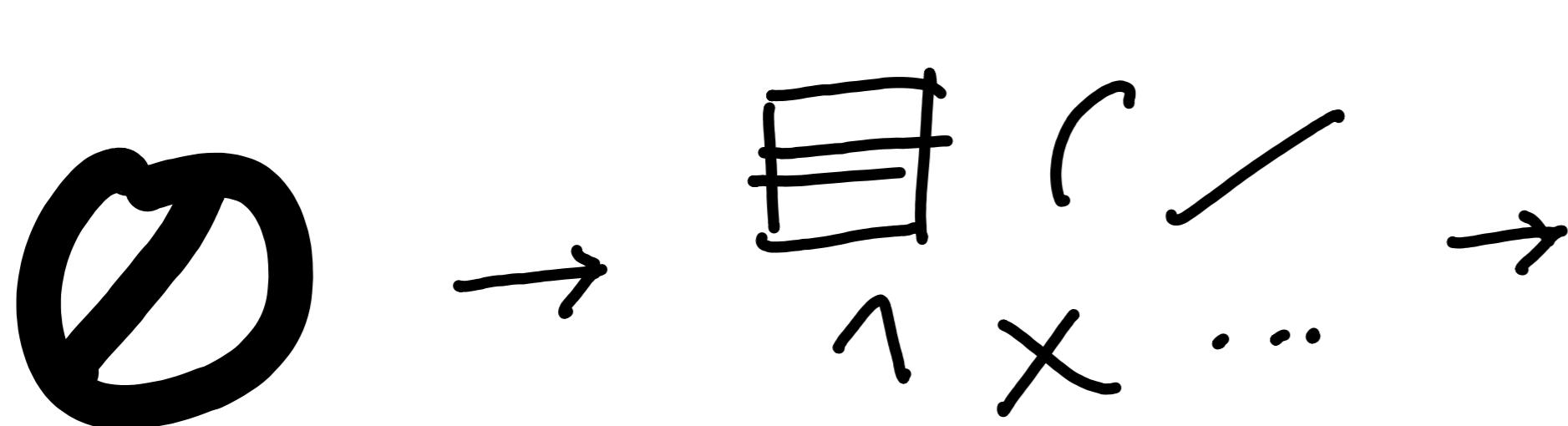
Feature: Does the image have some property (or how much)?



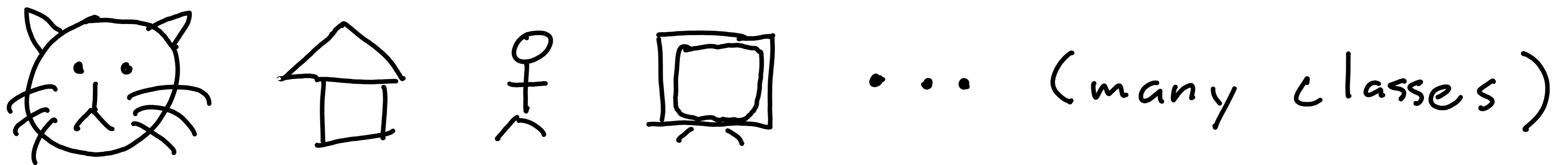
dark pixels



does the image have this shape?

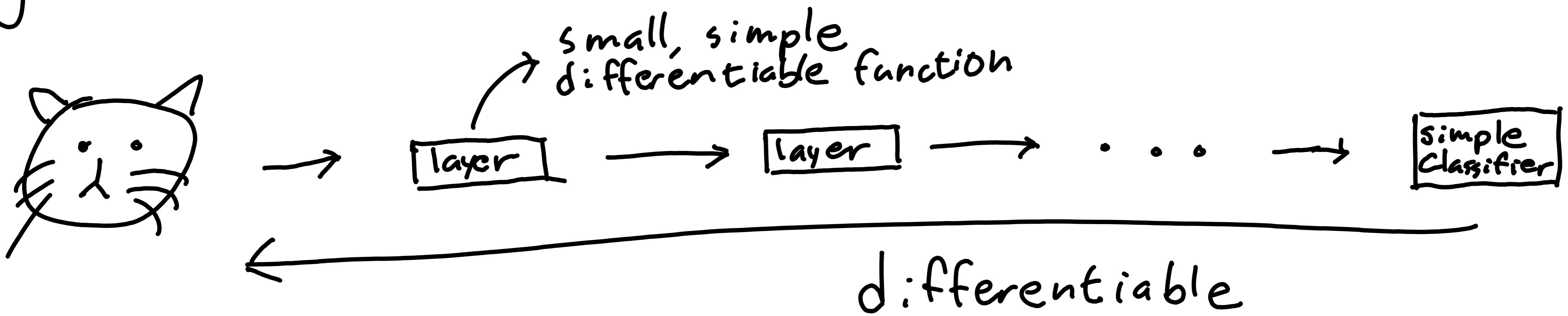


→ simple classifier

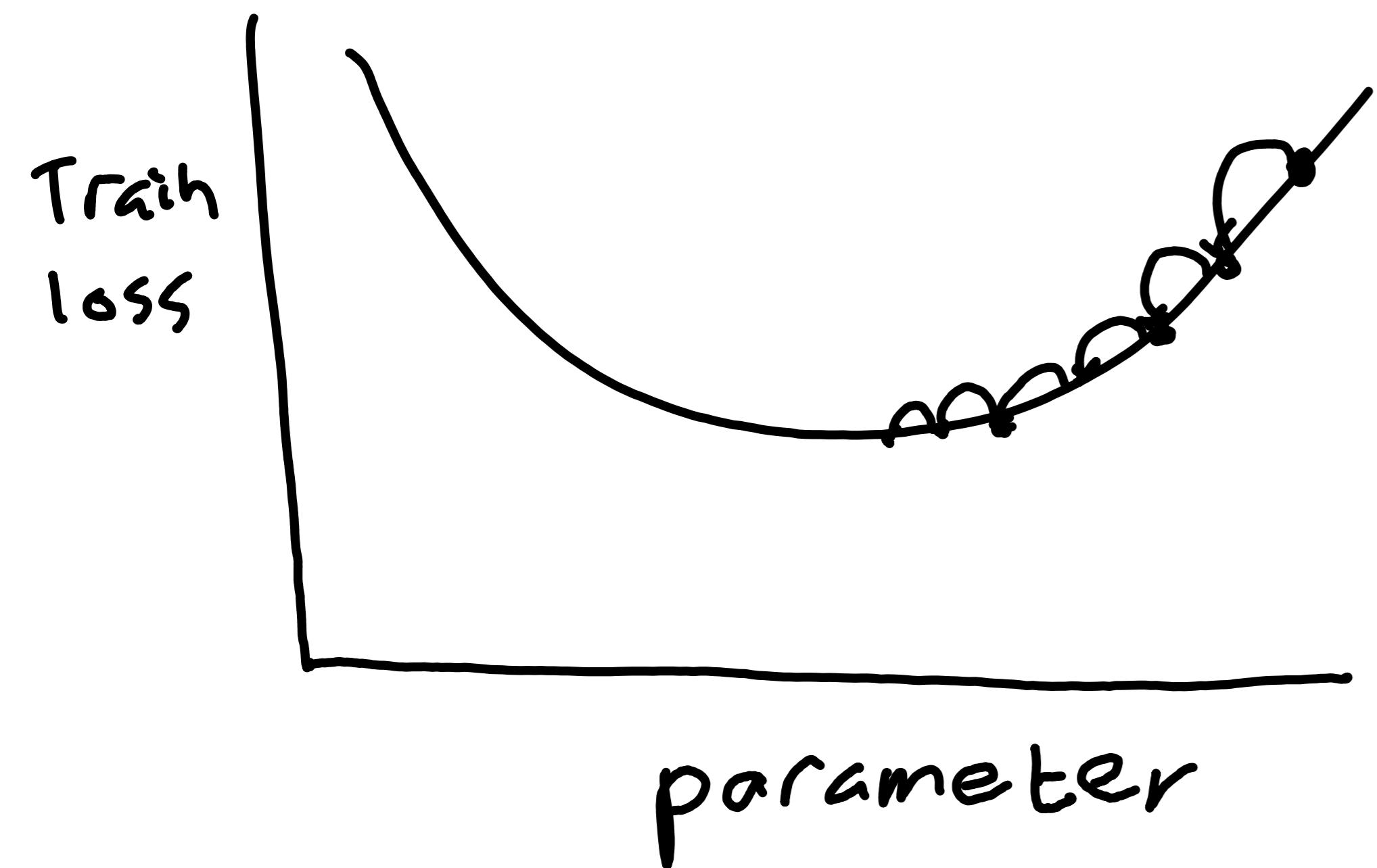


... (many classes)

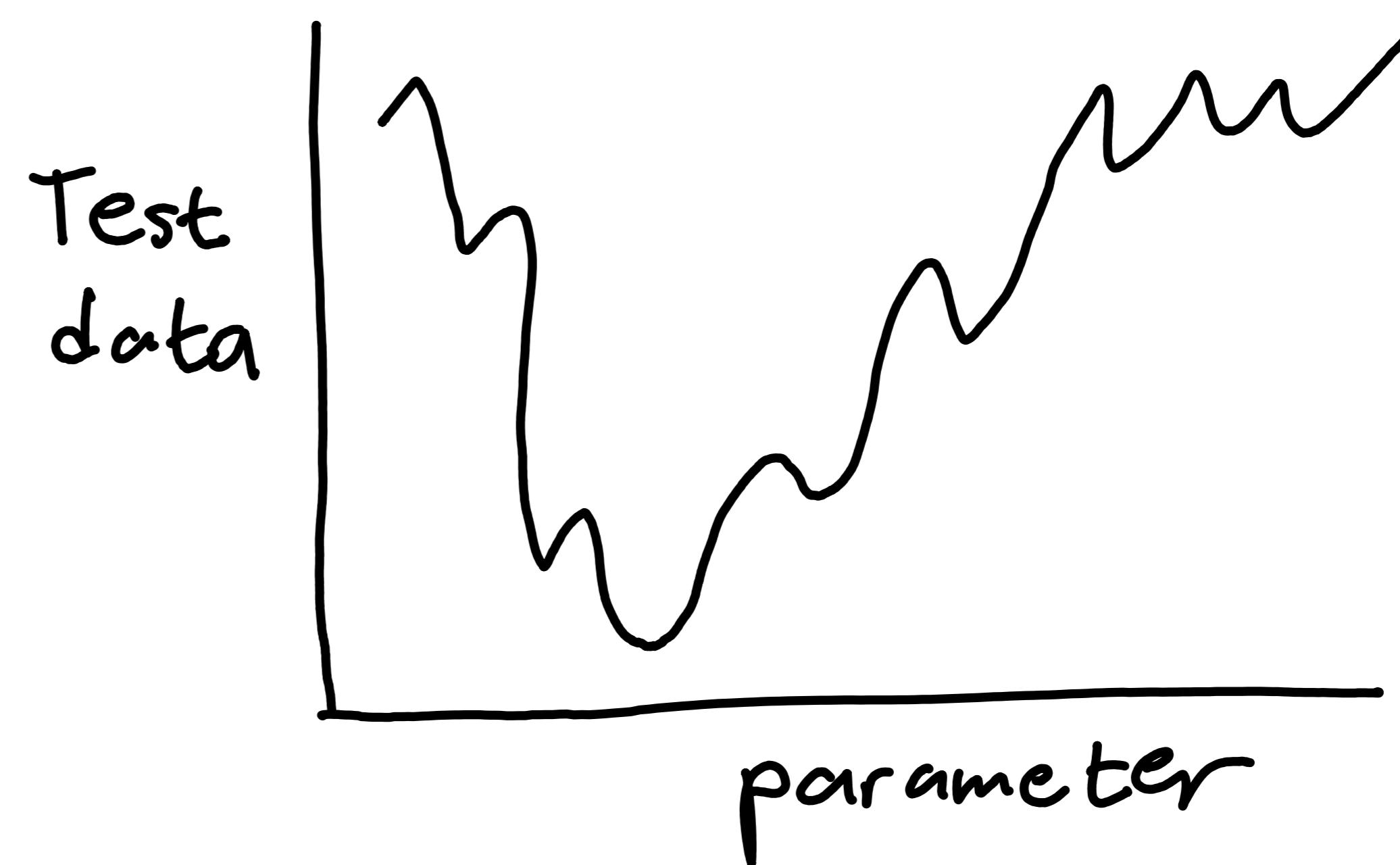
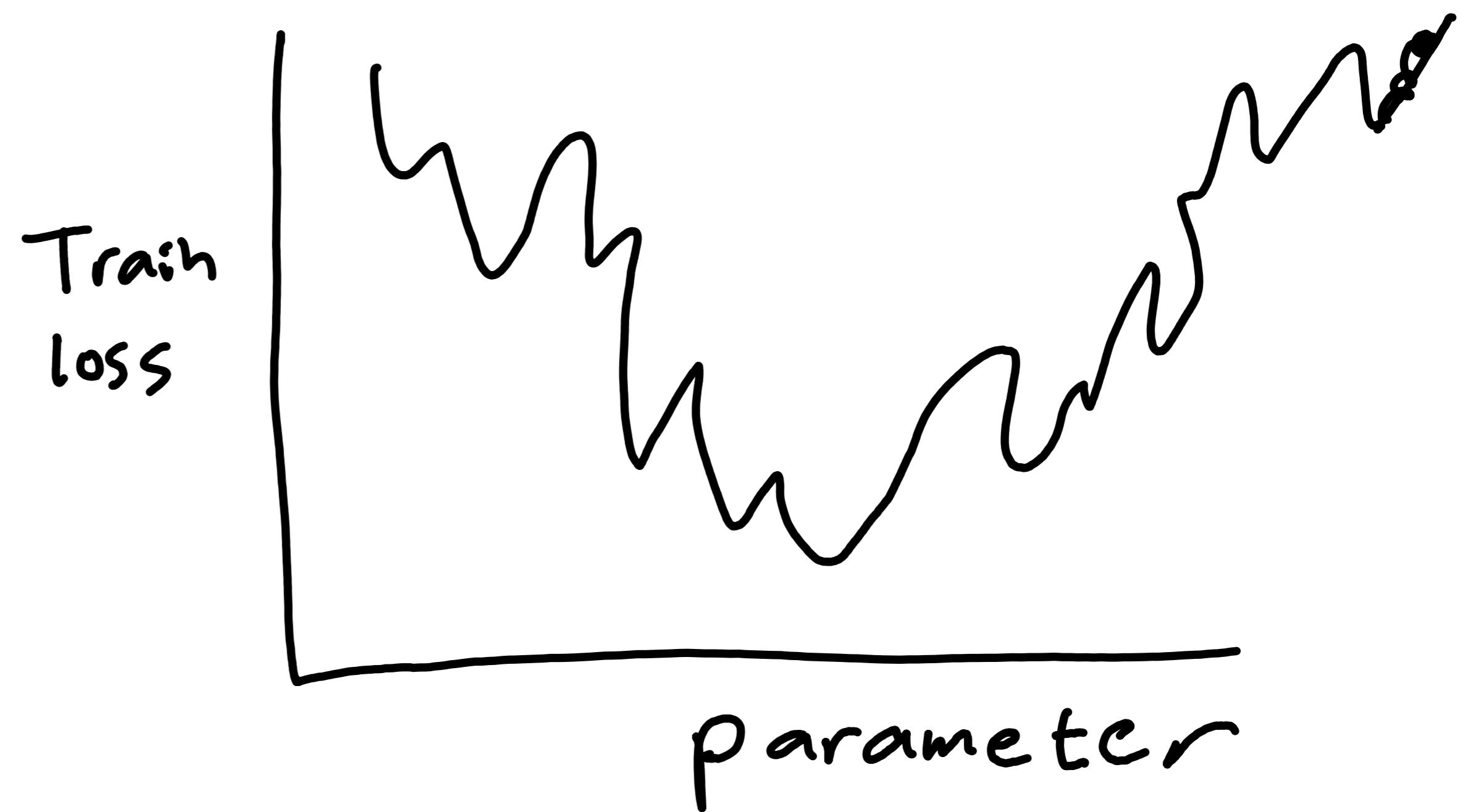
Deep learning: Learn features from data using big, complex, differentiable models that can be trained end-to-end with gradient descent.



Old - fashioned ML



Deep Learning



regularization, inductive
bias, pre-training, ...

Linear Regression

Linear model that predicts a scalar from some inputs.

↳ a single real number

~~ex~~ Predict house value from

$$\text{price} = w_1 \text{age} + w_2 \text{sqft} + w_3 \text{nbr} + b$$

↑ features
↑ parameter

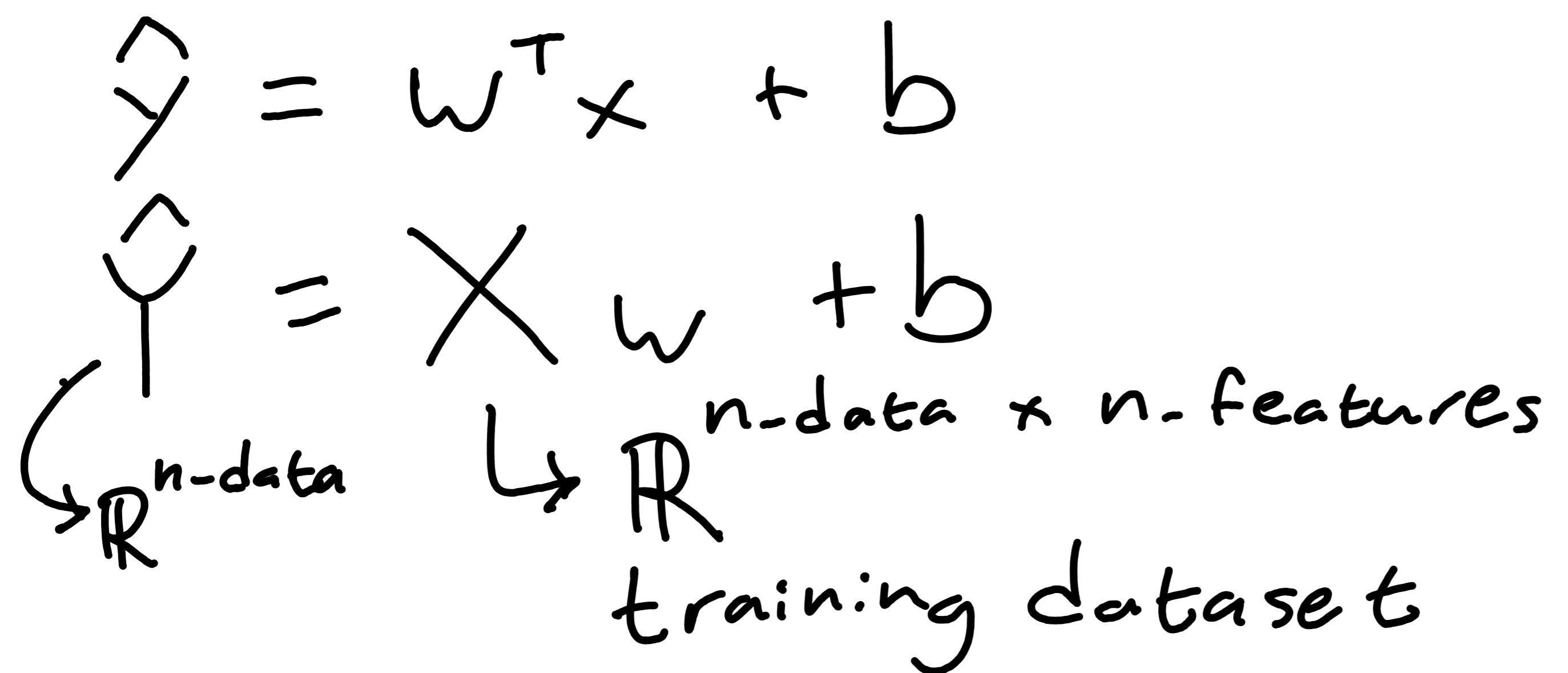
age, sq. ft., # bedrooms, lot size, ...

$$x = [\text{age}, \text{sqft}, \text{nbr}]$$

$$w = [w_1, w_2, w_3]$$

$$\text{price} = w^T x + b$$

$$\hat{y} = w^T x + b$$



 ↴ $R^{n\text{-data}}$ ↴ $R^{n\text{-data} \times n\text{-features}}$
 training dataset

Goal: fit the parameters (w, b)
using our training data.

How well are we doing?

How closely are predicted \hat{y}
to the true values in the training set?

Loss function:

$$\frac{1}{2} (\hat{y} - y)^2$$

prediction target

squared
error

Goal: Change parameters to minimize
the loss over the training dataset.

~~Solve?~~

$$\hat{w}^* = \cancel{(X^T X)}^{-1} X^T \cancel{y}$$

Use gradient descent!

Gradient descent:

1. Initialize parameters Θ (w, b)
2. Repeat:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} L$$

new param values old param values learning rate gradient of loss (L) w.r.t. θ

until some criterion is met.

ex loss is zero. (rare)

ex loss stops decreasing.

ex performance on held-out data stops improving.

ex run out of compute or patience.

$$\hat{y} = w^T x + b$$

$$L = \frac{1}{2} (\hat{y} - y)^2$$

(x, y) is a single training datapoint.

We need $\nabla_{\theta} L$ ($\nabla_w L$ and $\nabla_b L$).

Use the chain rule!

Functions: $L(\hat{y})$, $\hat{y}(w)$ (or $\hat{y}(b)$)

$$\frac{dL}{dw} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dw}$$

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{db}$$

$$\hat{y} = \omega^T x + b$$

$$L = \frac{1}{2} (\hat{y} - y)^2$$

$$\frac{dL}{d\hat{y}} = \hat{y} - y$$

$$\frac{dL}{dx} = x$$

$$\frac{dL}{db} = 1$$

}

$$\nabla_w L$$

$$\frac{dL}{dw} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dw} = (\hat{y} - y)x$$

$$\nabla_b L$$

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{db} = \hat{y} - y$$

$$w \leftarrow w - \gamma \nabla_w L$$

use 1 example \rightarrow stochastic gradient descent

"batch" X , a matrix of examples

$$X \in \mathbb{R}^{n\text{-features}}, X \in \mathbb{R}^{n\text{-data} \times n\text{-features}}$$

$$w \leftarrow w - \frac{\gamma}{n_{\text{data}}} \sum_{\text{ndata}} \nabla_w L$$

$\xrightarrow{\text{random subset of training set}}$
"minibatch stochastic
gradient descent"

$$w \leftarrow w - \frac{\eta}{n_{\text{data}}} \sum_{\text{n}_{\text{data}}} \nabla_w L$$

Hyper Parameters:

Values set "by hand" instead of G.D.

1. η learning rate (step size)

2. "batch size" (n_{data})

3. Stopping criterion

4. Initialization scheme