

UNIVERSITÄT
HEIDELBERG



Machine Learning for Biochemistry

L2, Structural Bioinformatics

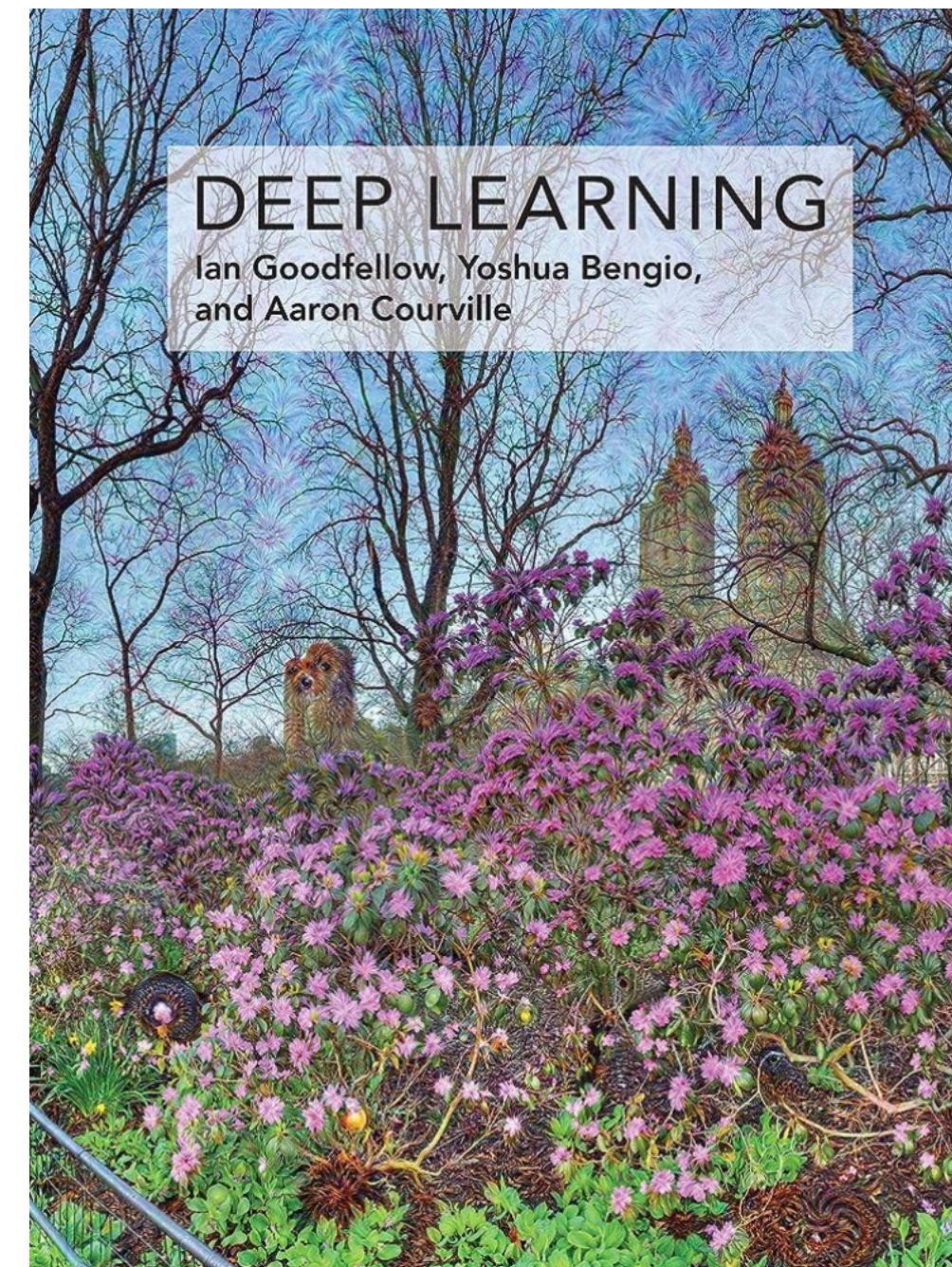
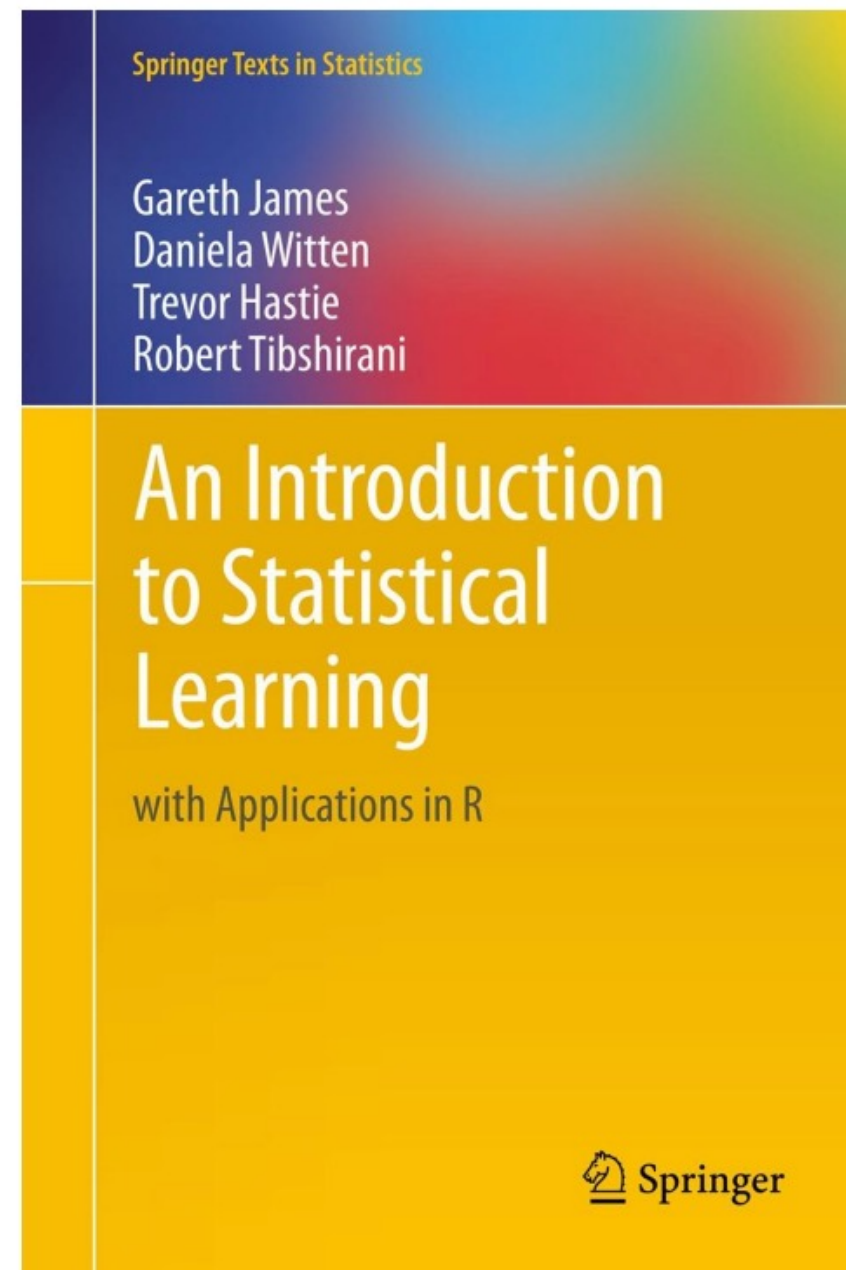
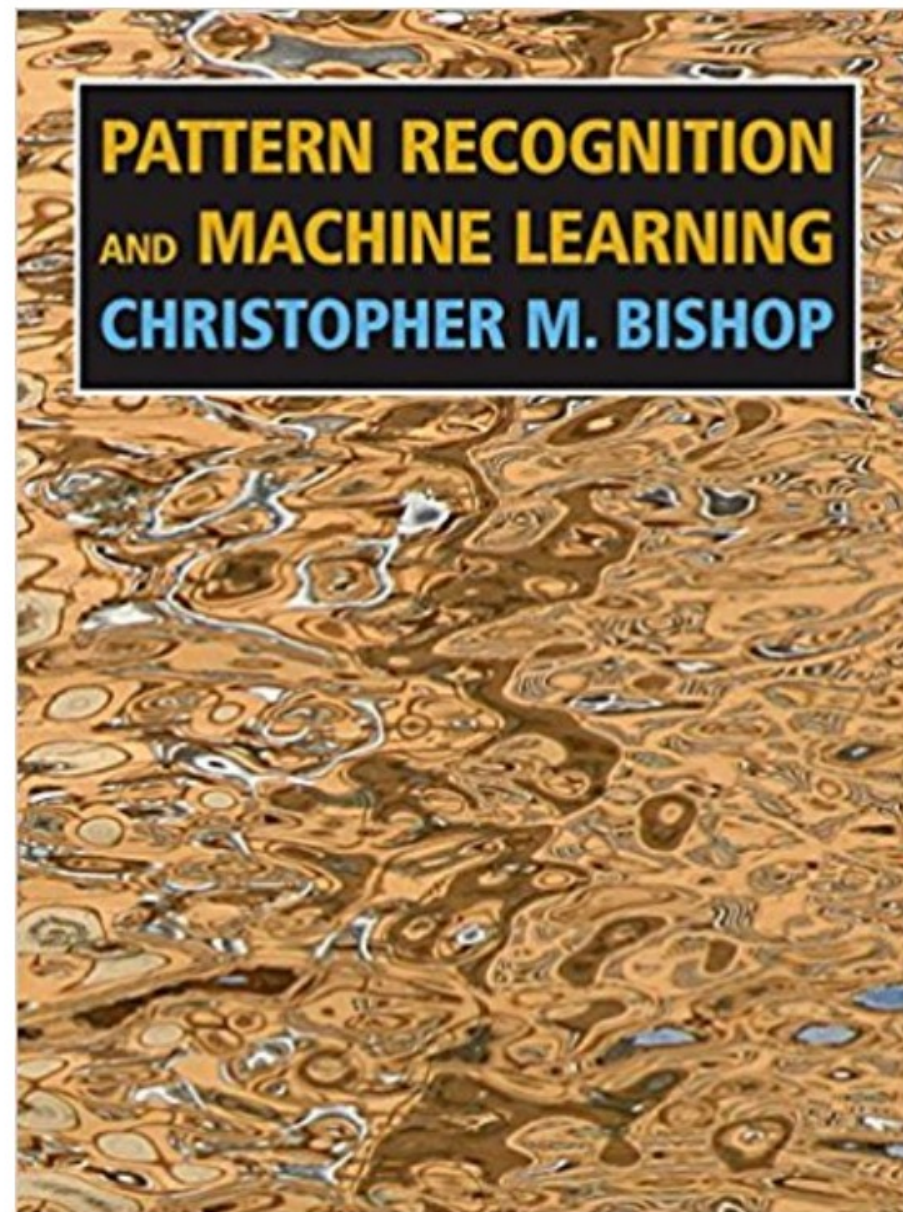
WiSe 2023/24, Heidelberg University

Overview

1. **Types of Machine Learning**
2. **Linear Models**
3. **Gradient Descent**
4. **Deep Learning**
5. **Outlook for what's to come**

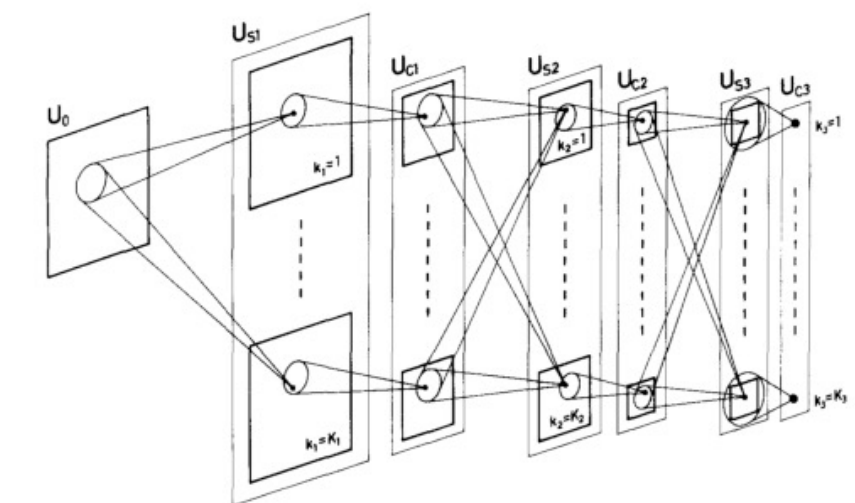
Book Recommendations

The Classics



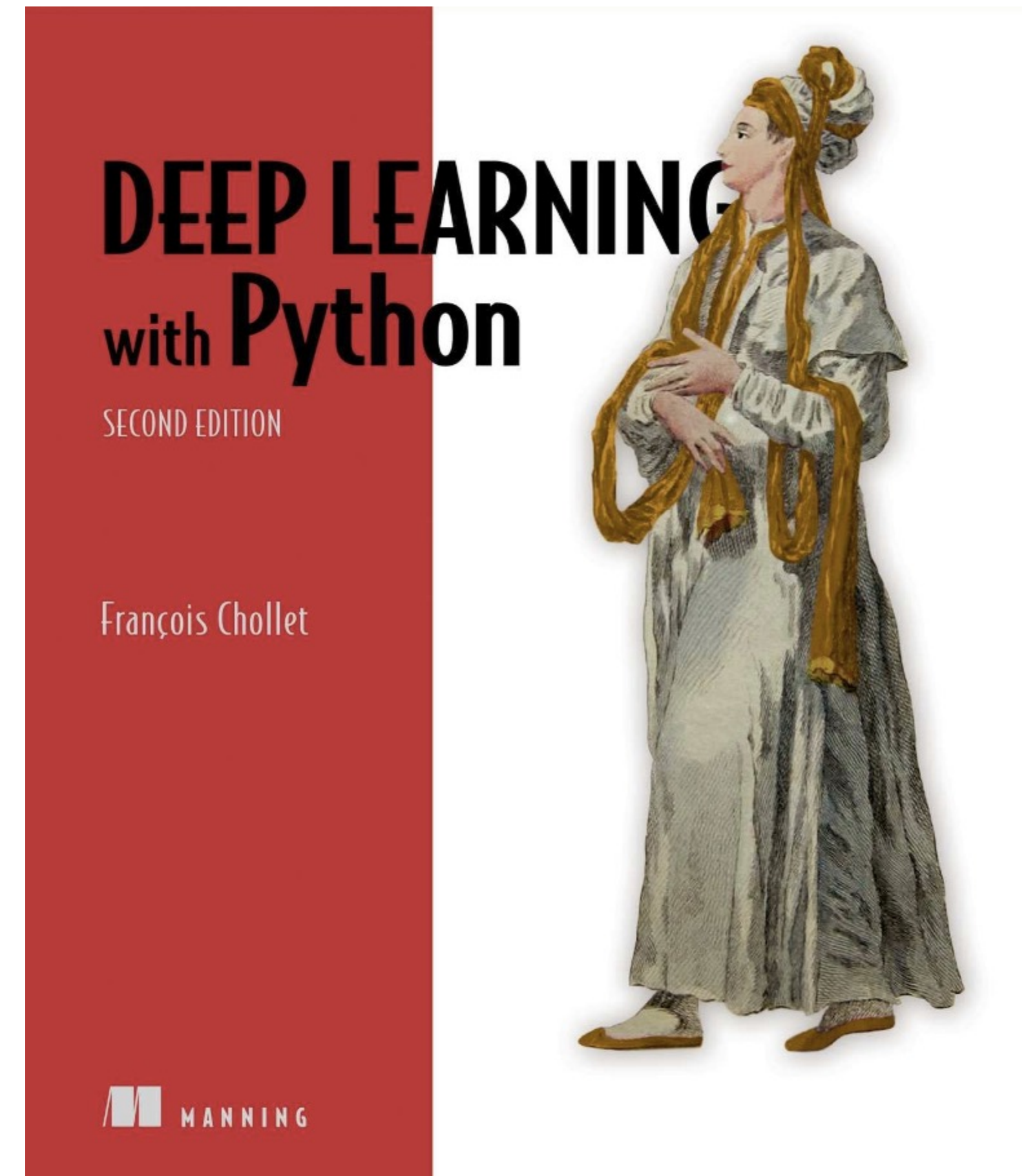
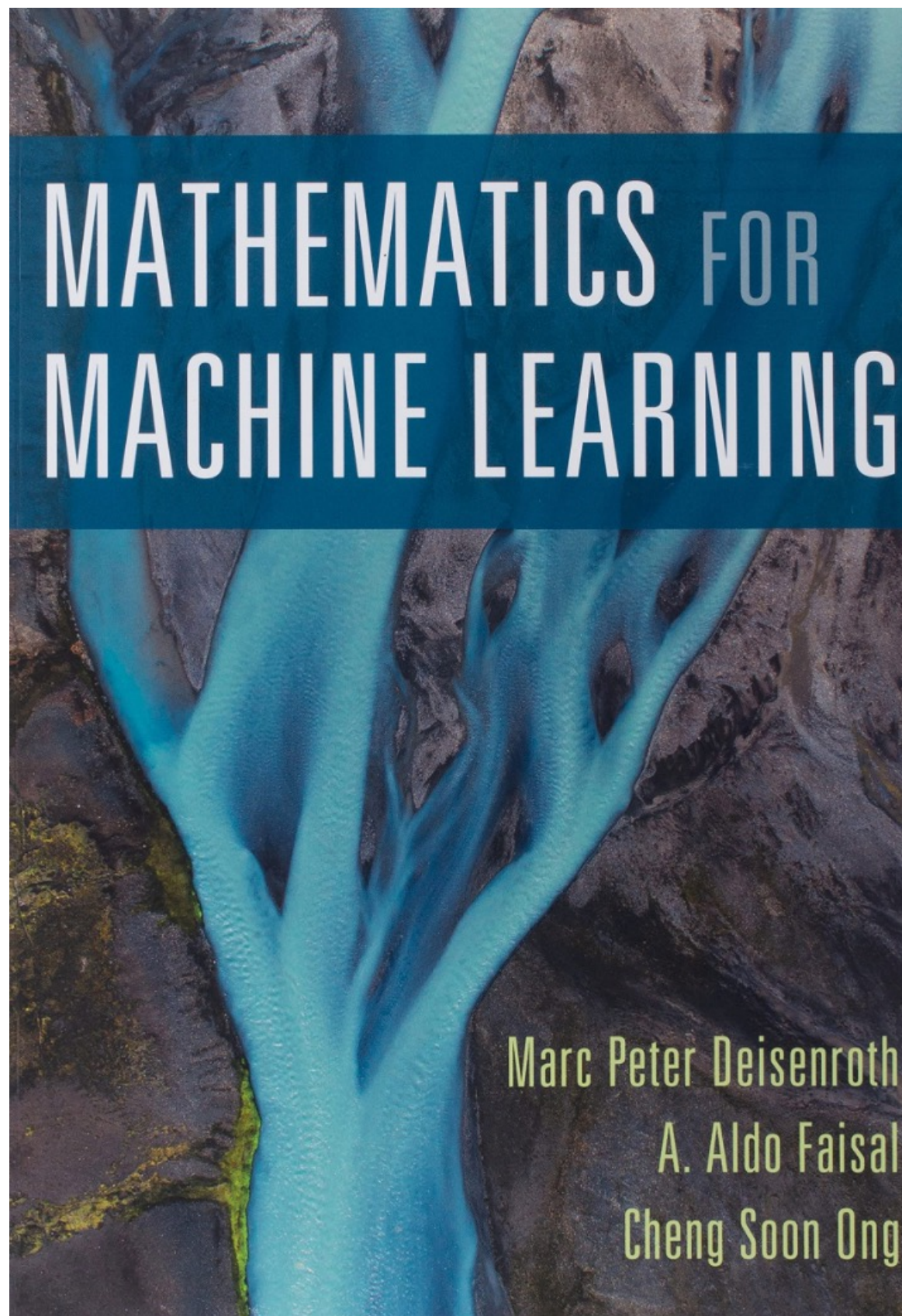
The Little Book of Deep Learning

François Fleuret

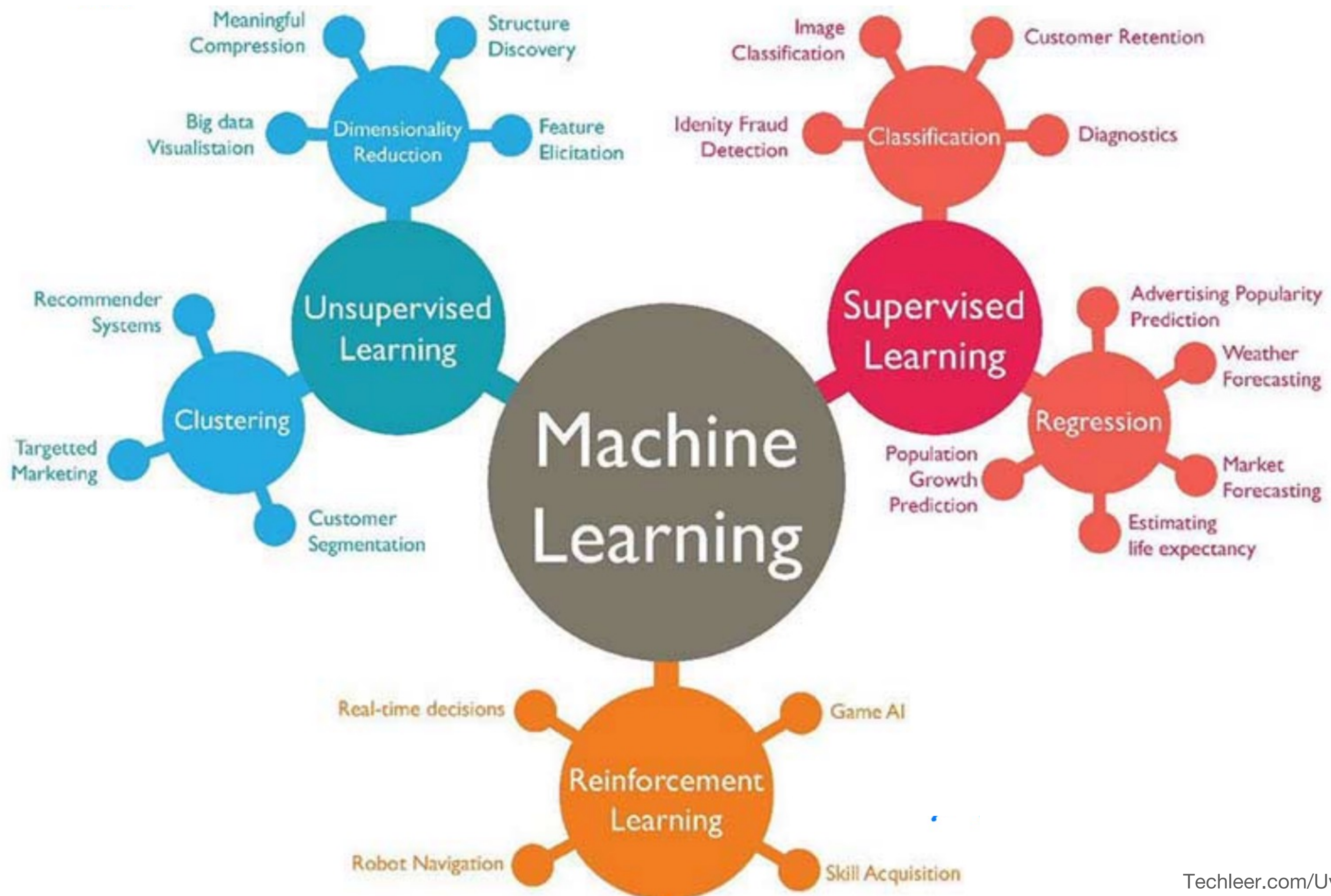


Book Recommendations

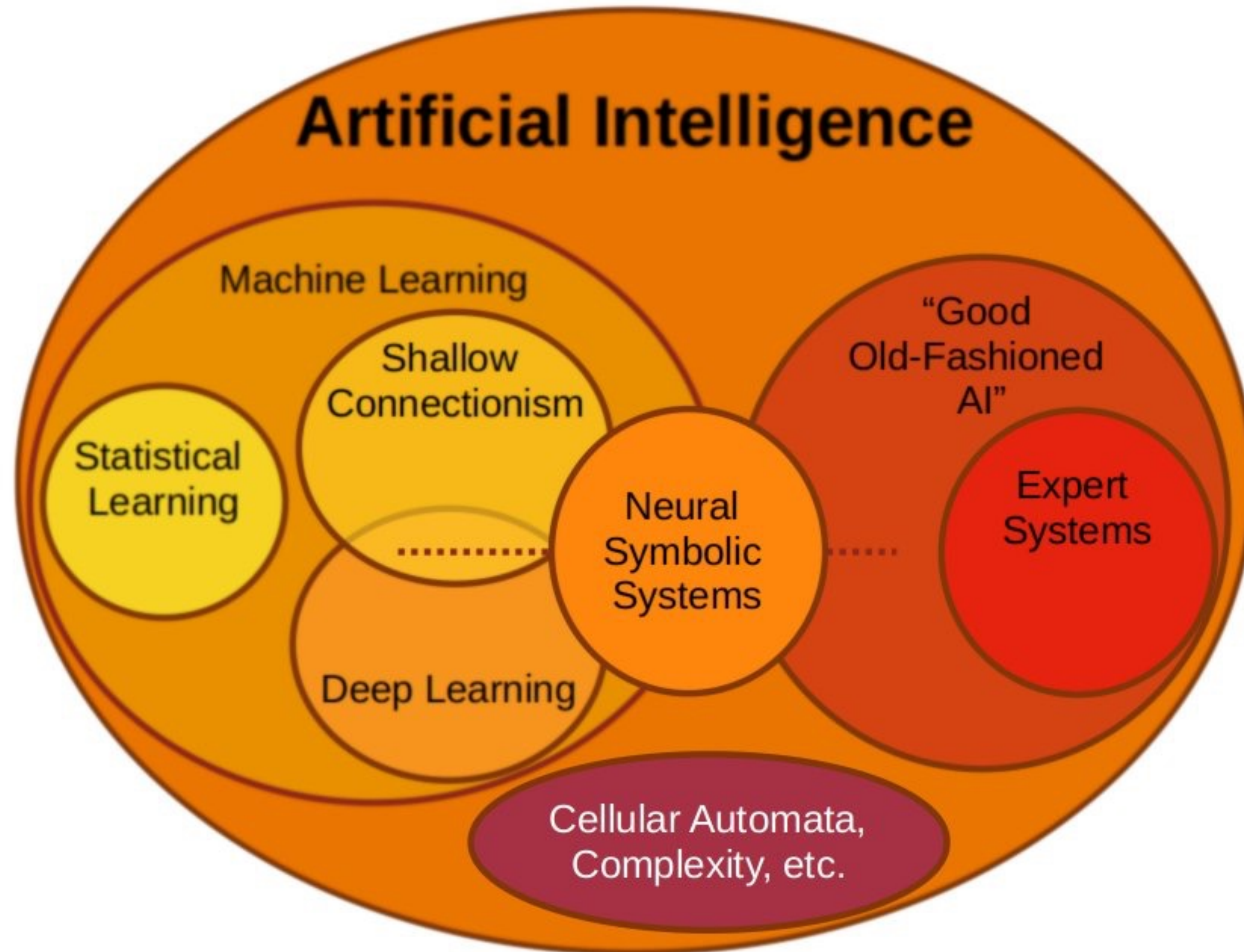
Background and more introductory books



1 Types of Machine Learning



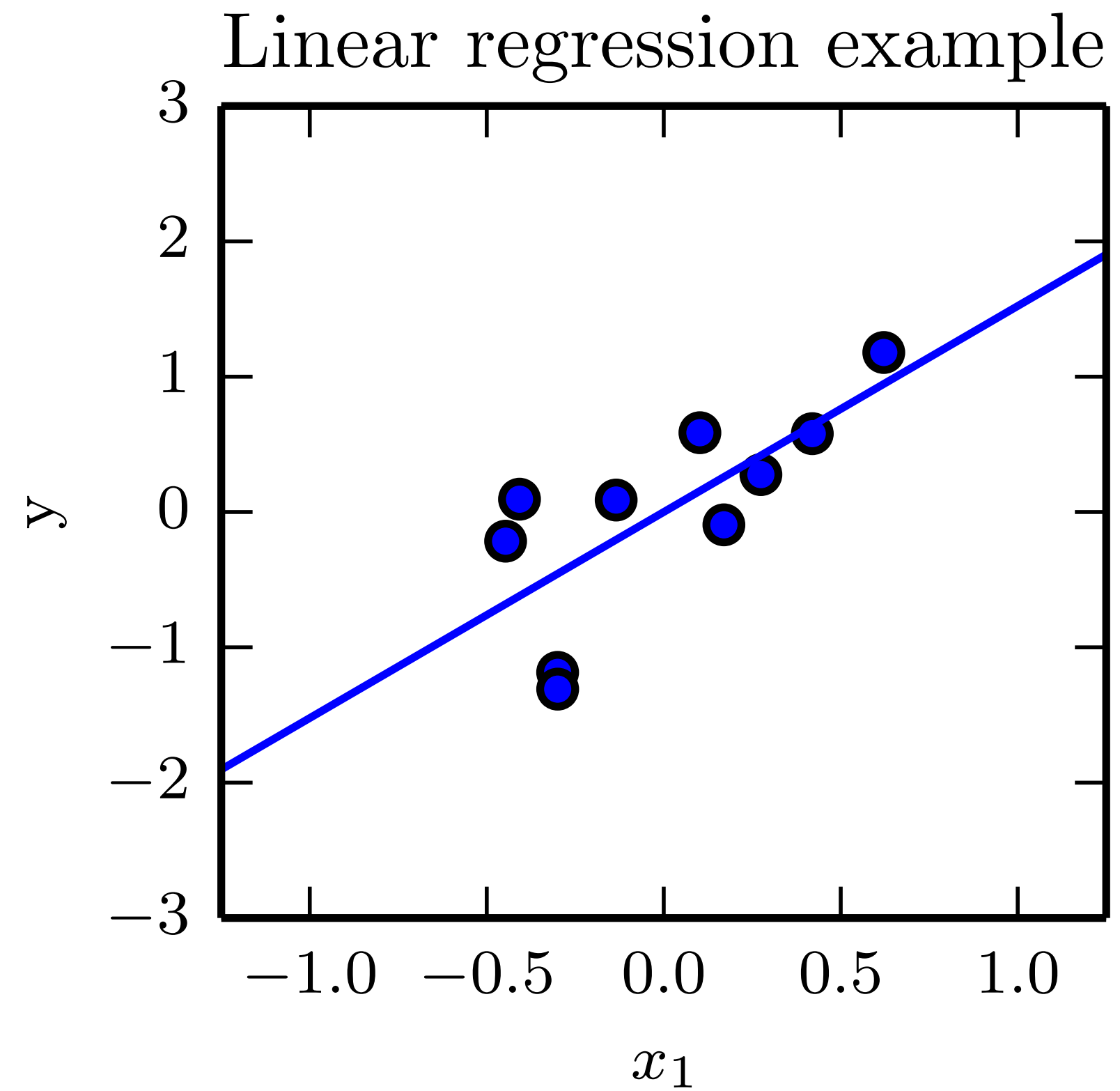
AI is more than just Deep Learning



2 Linear Models

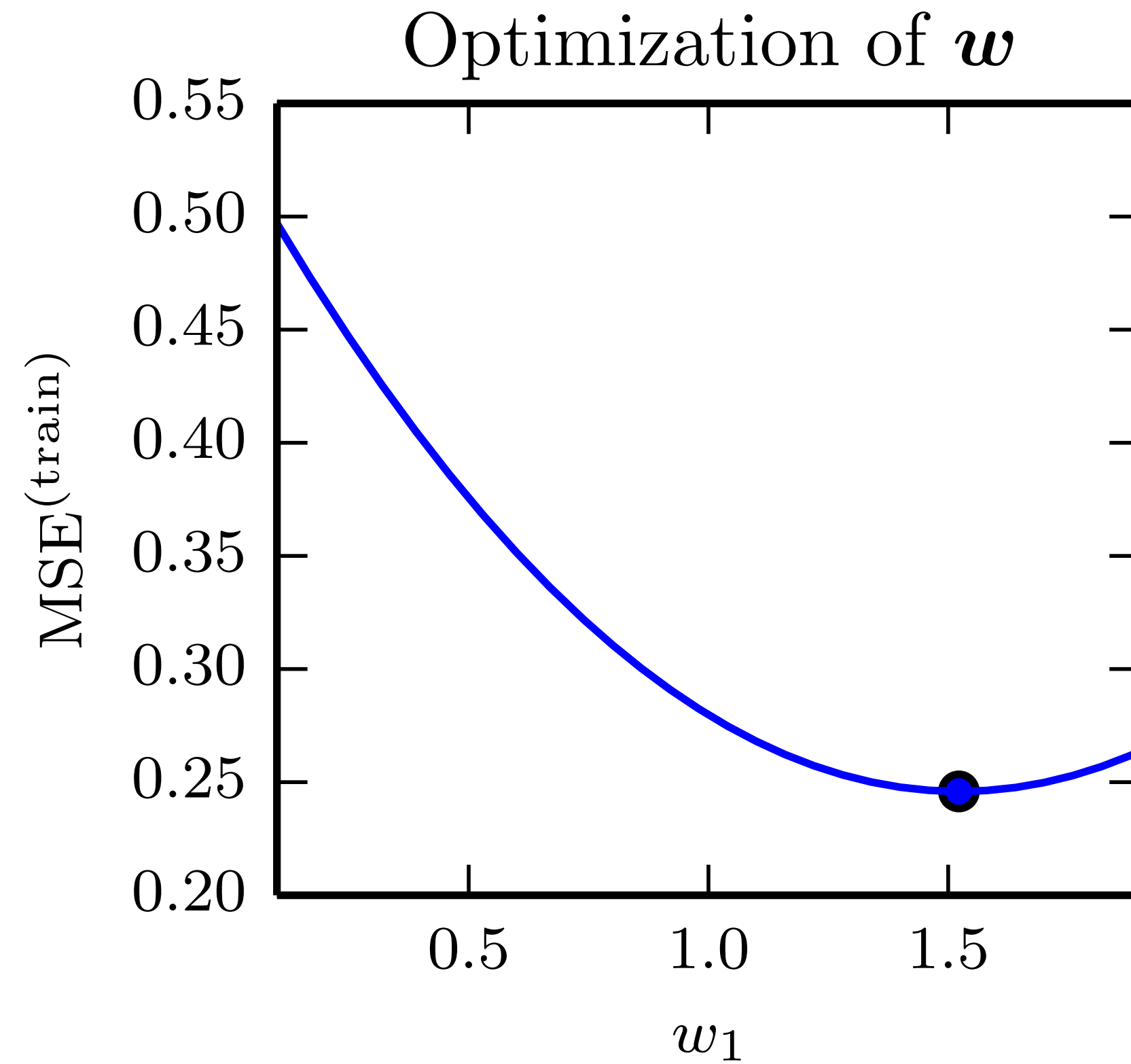
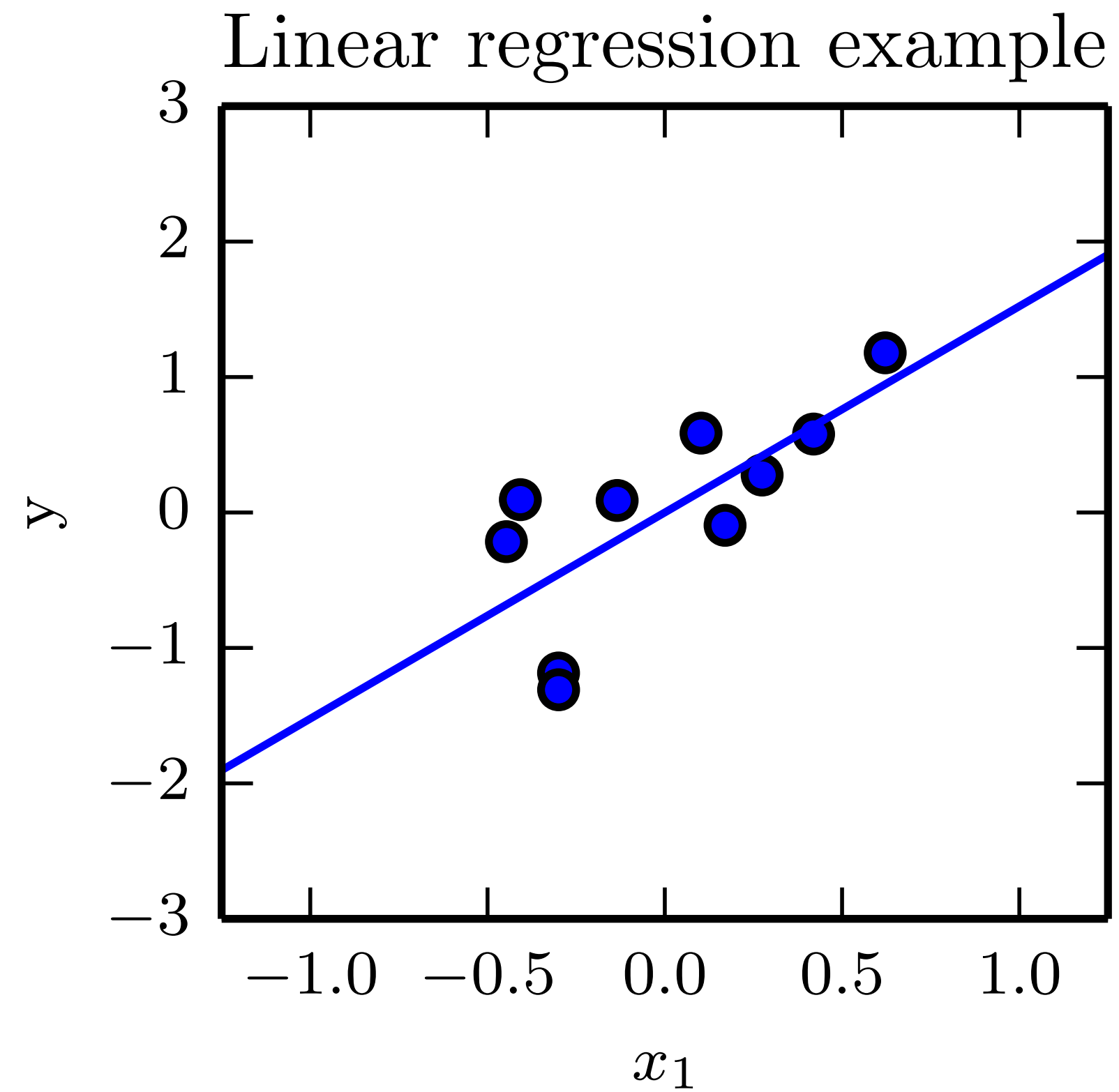
What is a linear model?

Adjust your free parameter based on some loss



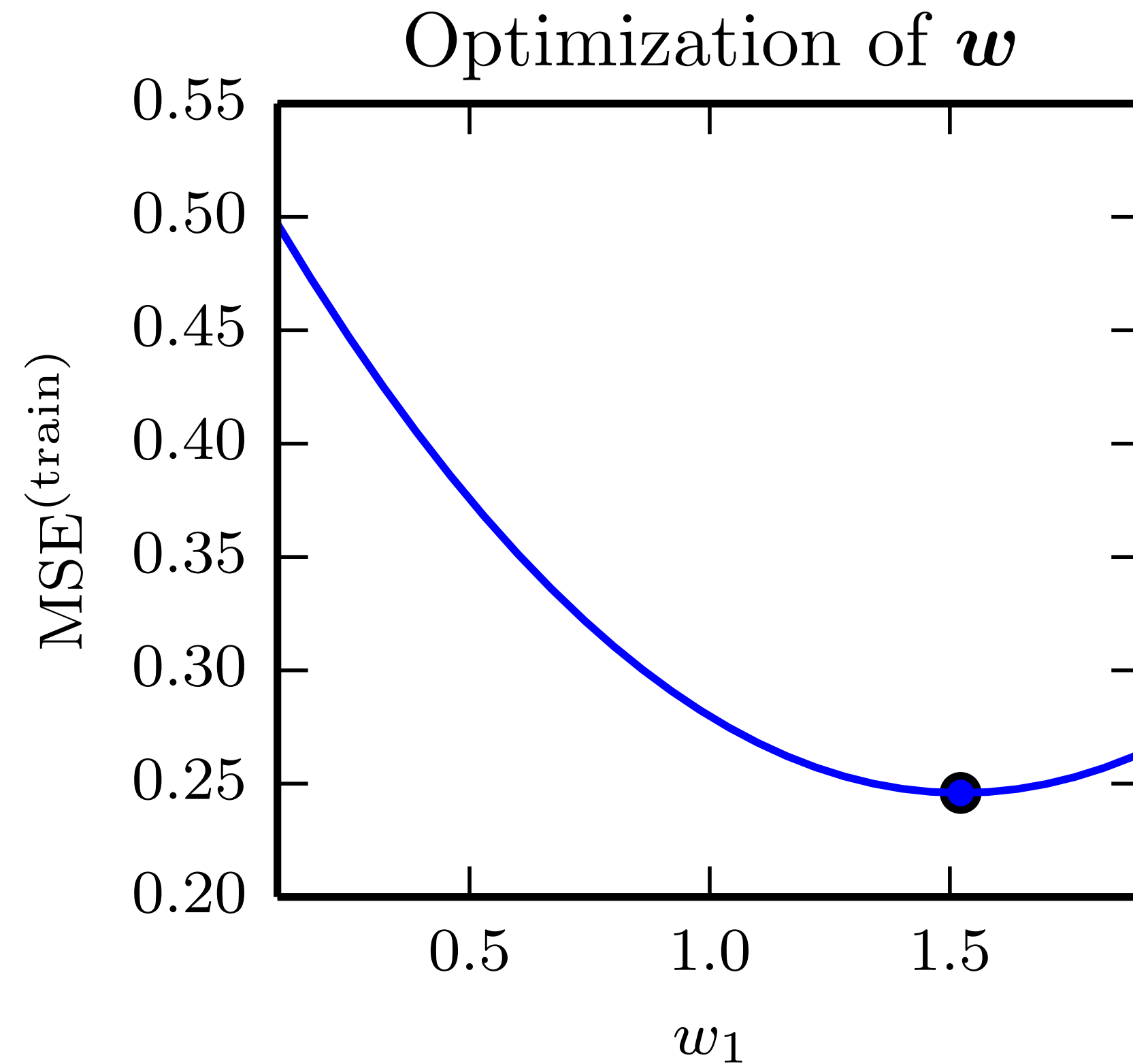
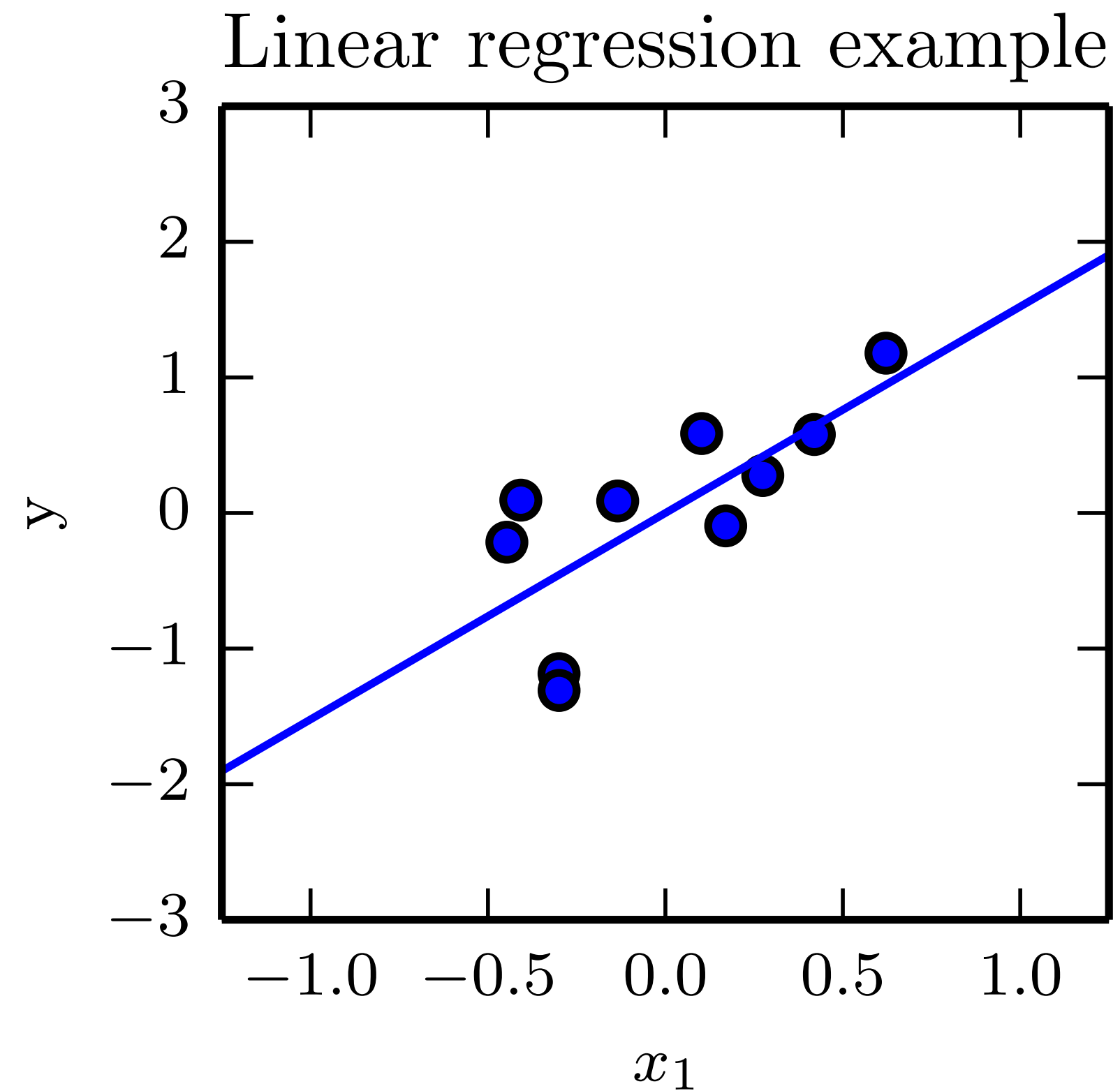
What is a linear model?

Adjust your free parameter based on some loss



What is a linear model?

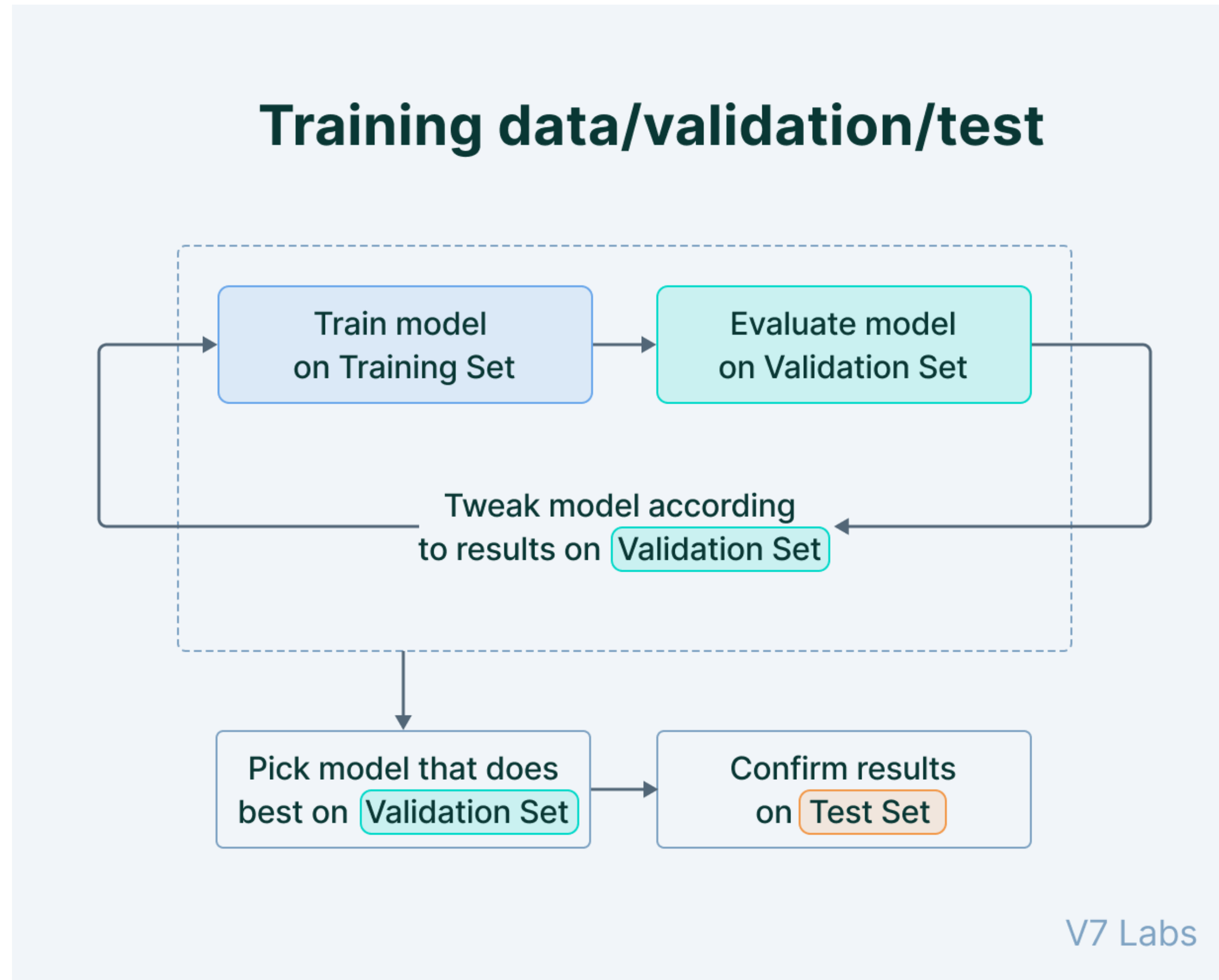
Adjust your free parameter based on some loss



MSE (Mean-squared error): $L[f] := \sum_{x \in X} (f(x; \theta) - y(x))^2$

What is a train versus a test dataset?

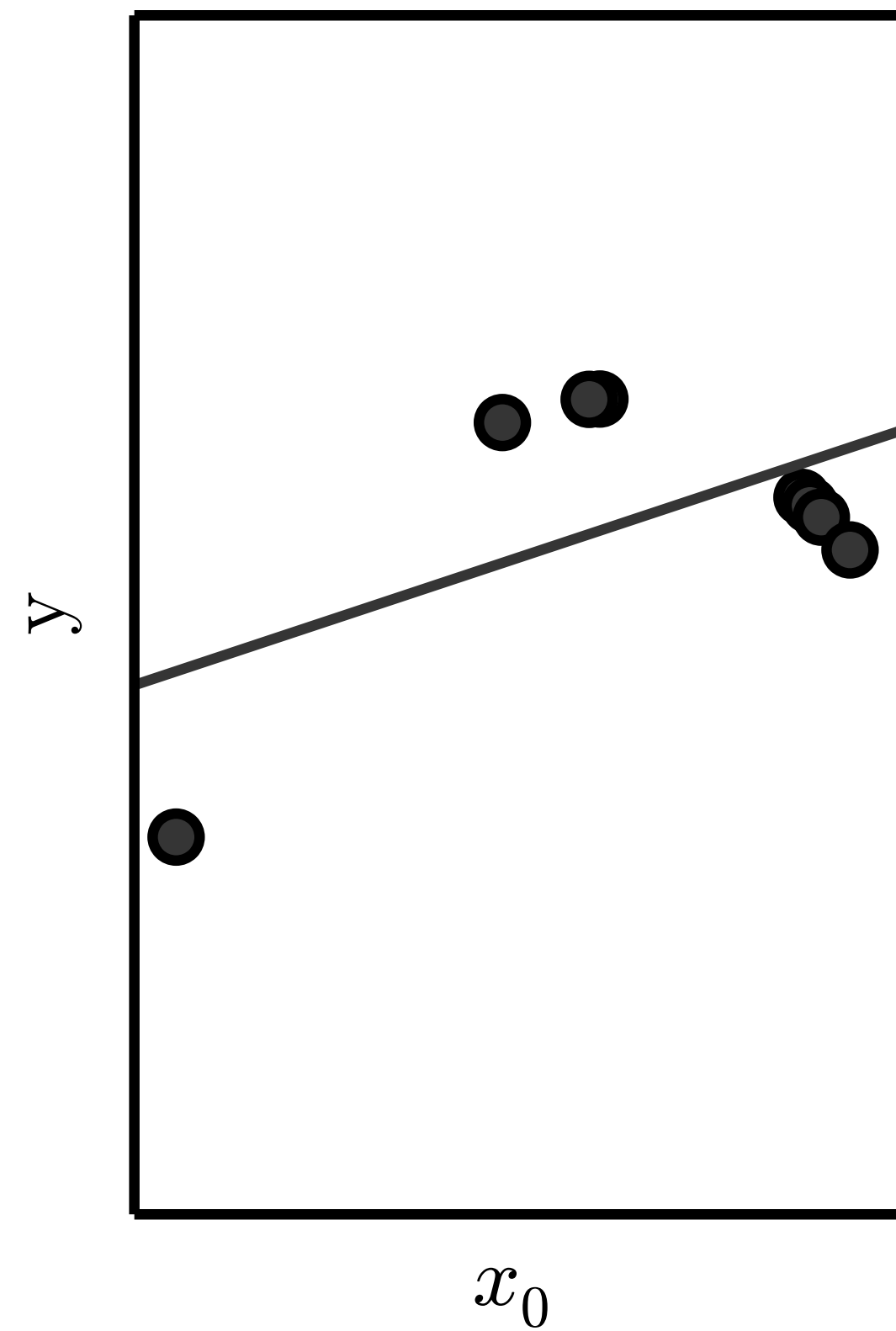
Evaluate how well your model generalises



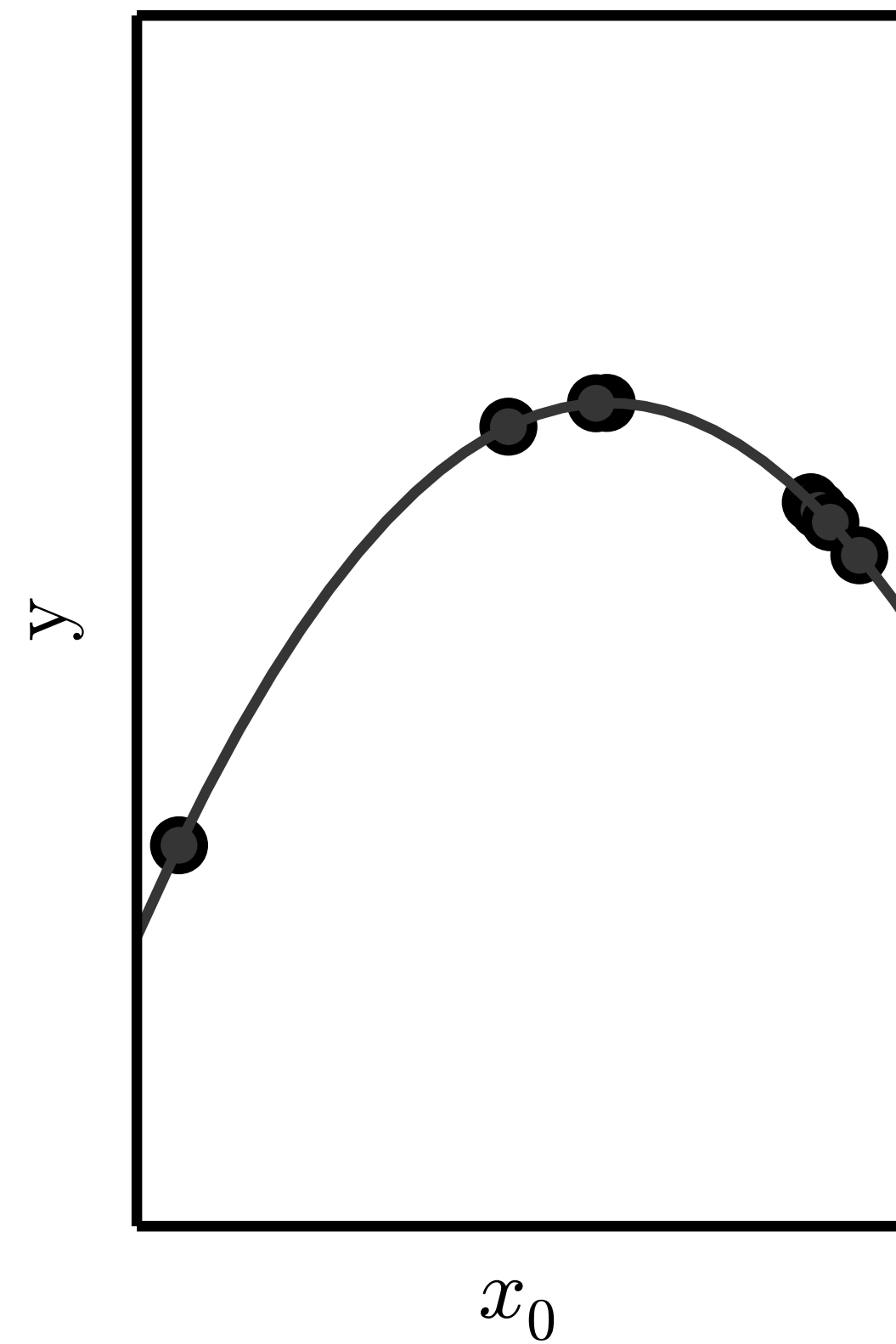
Why limit yourself to linear models?

Varying the degree of basis function results in different fits

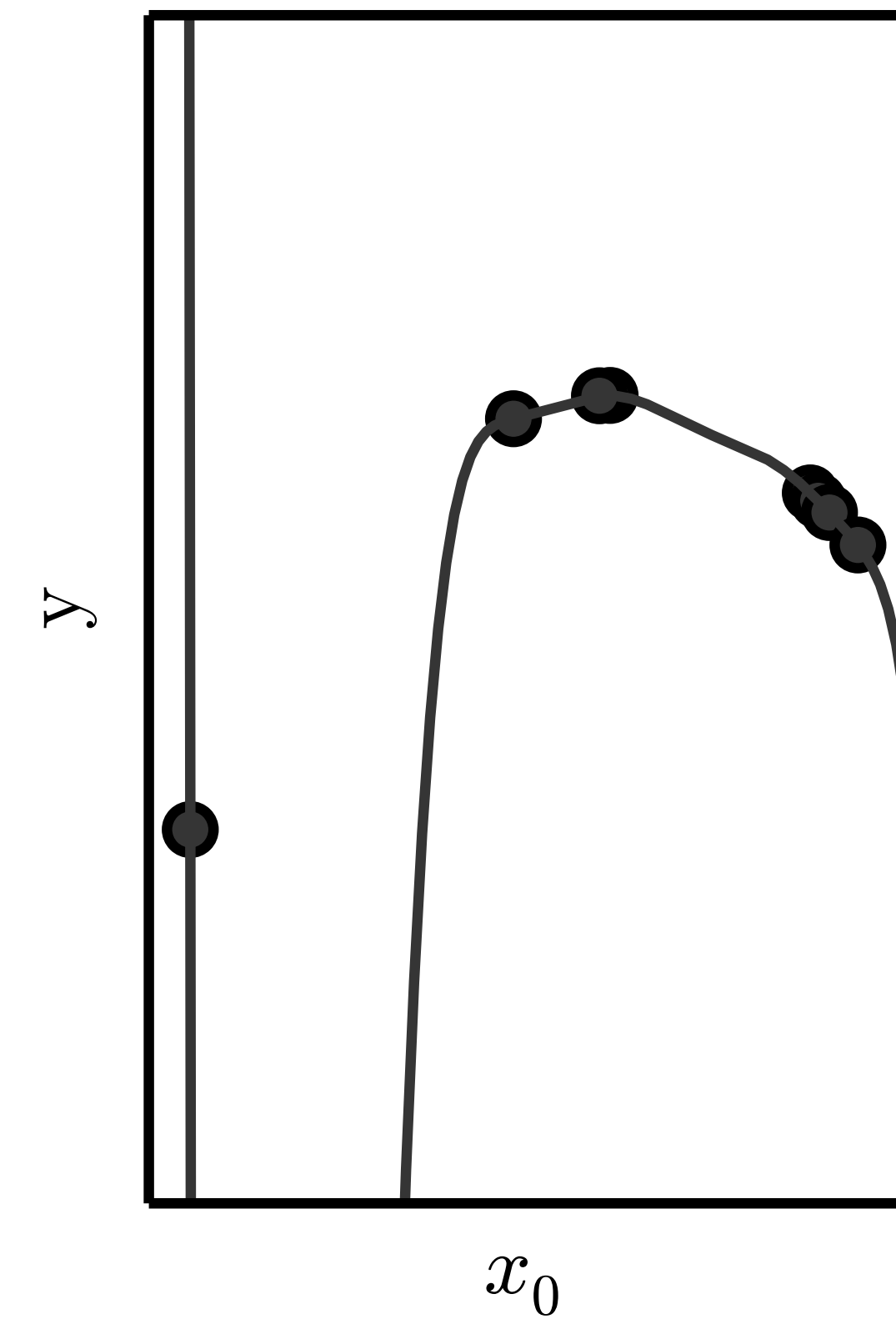
Underfitting



Appropriate capacity

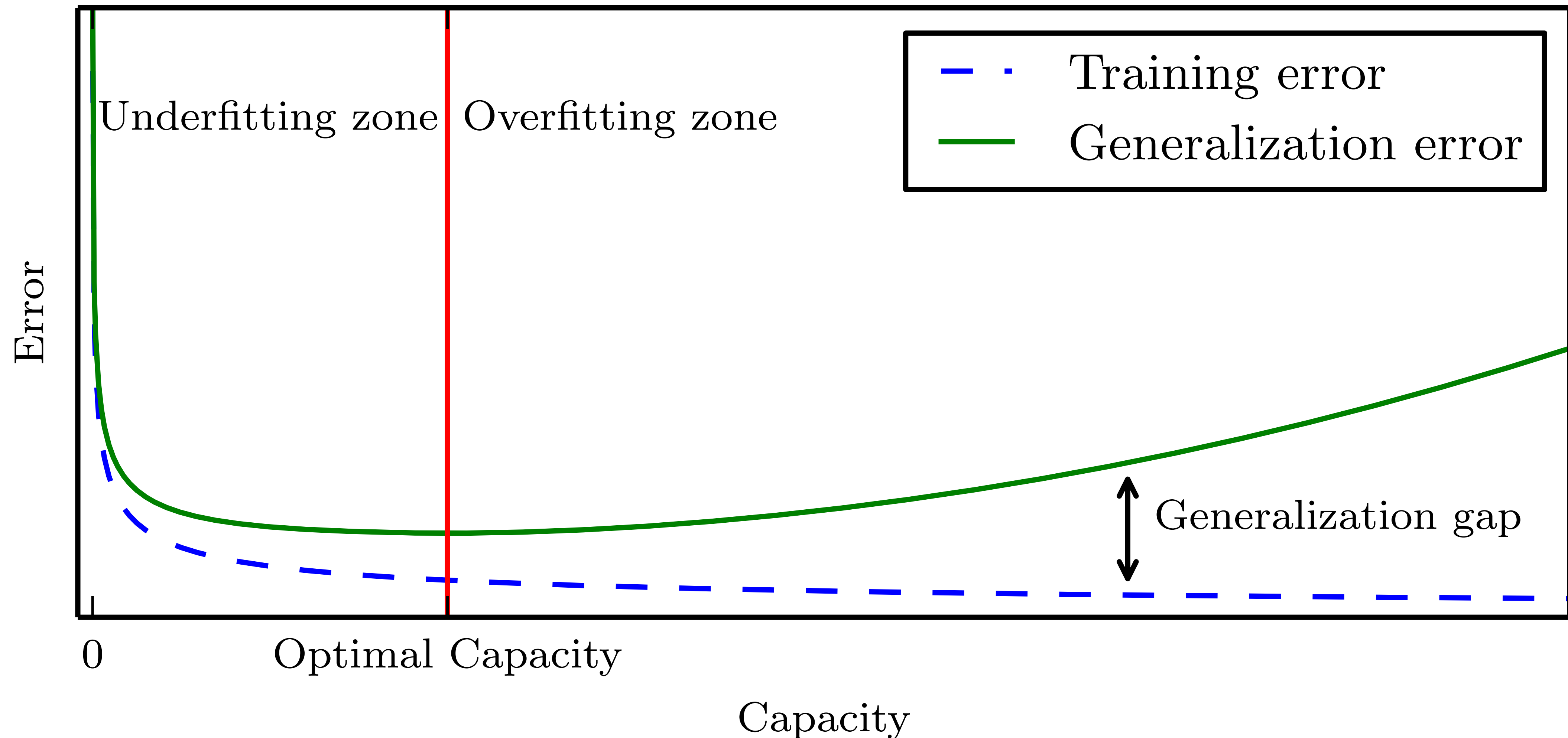


Overfitting



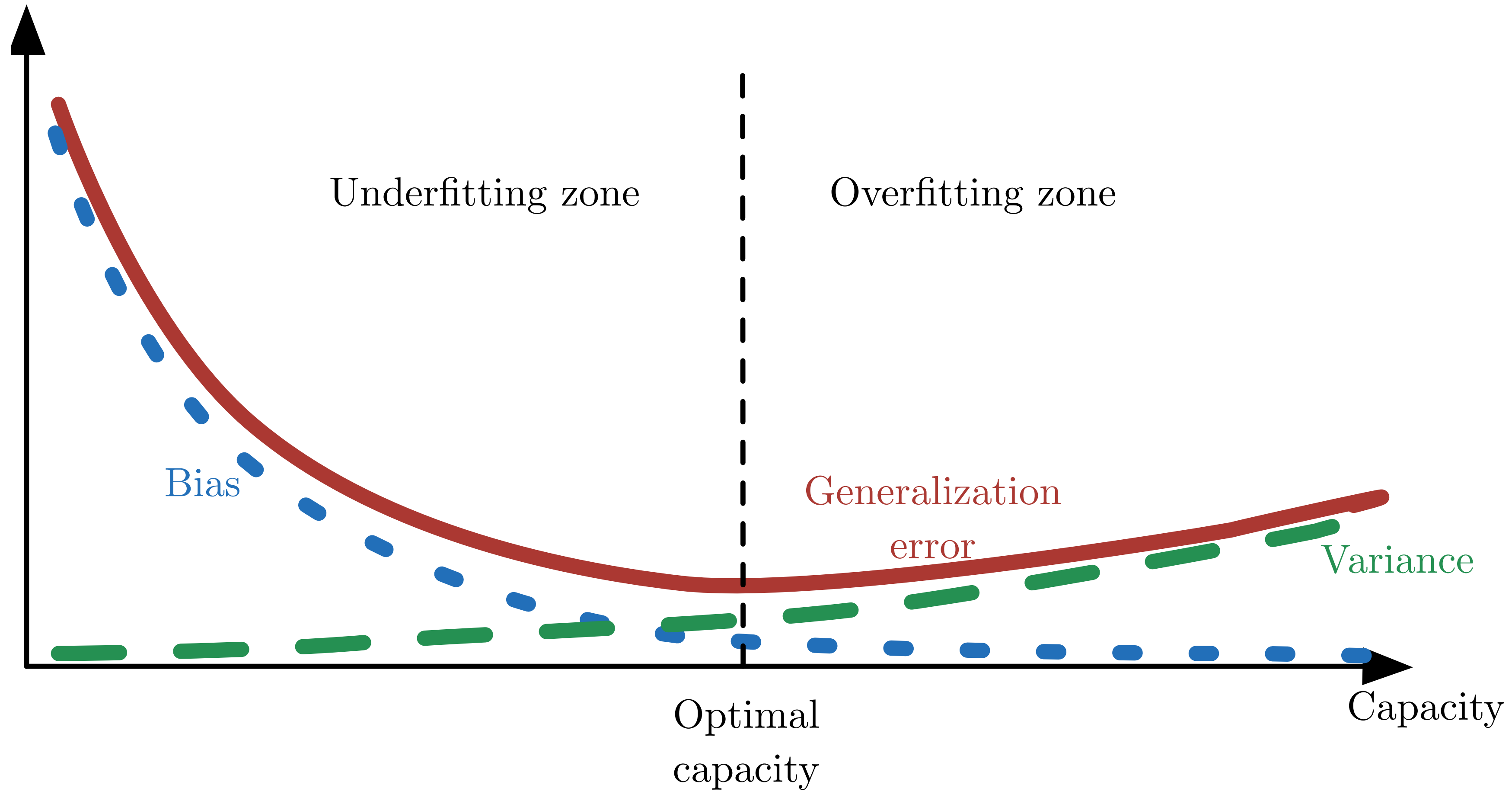
Always look at your test set!

Machine learning models like to overfit your data



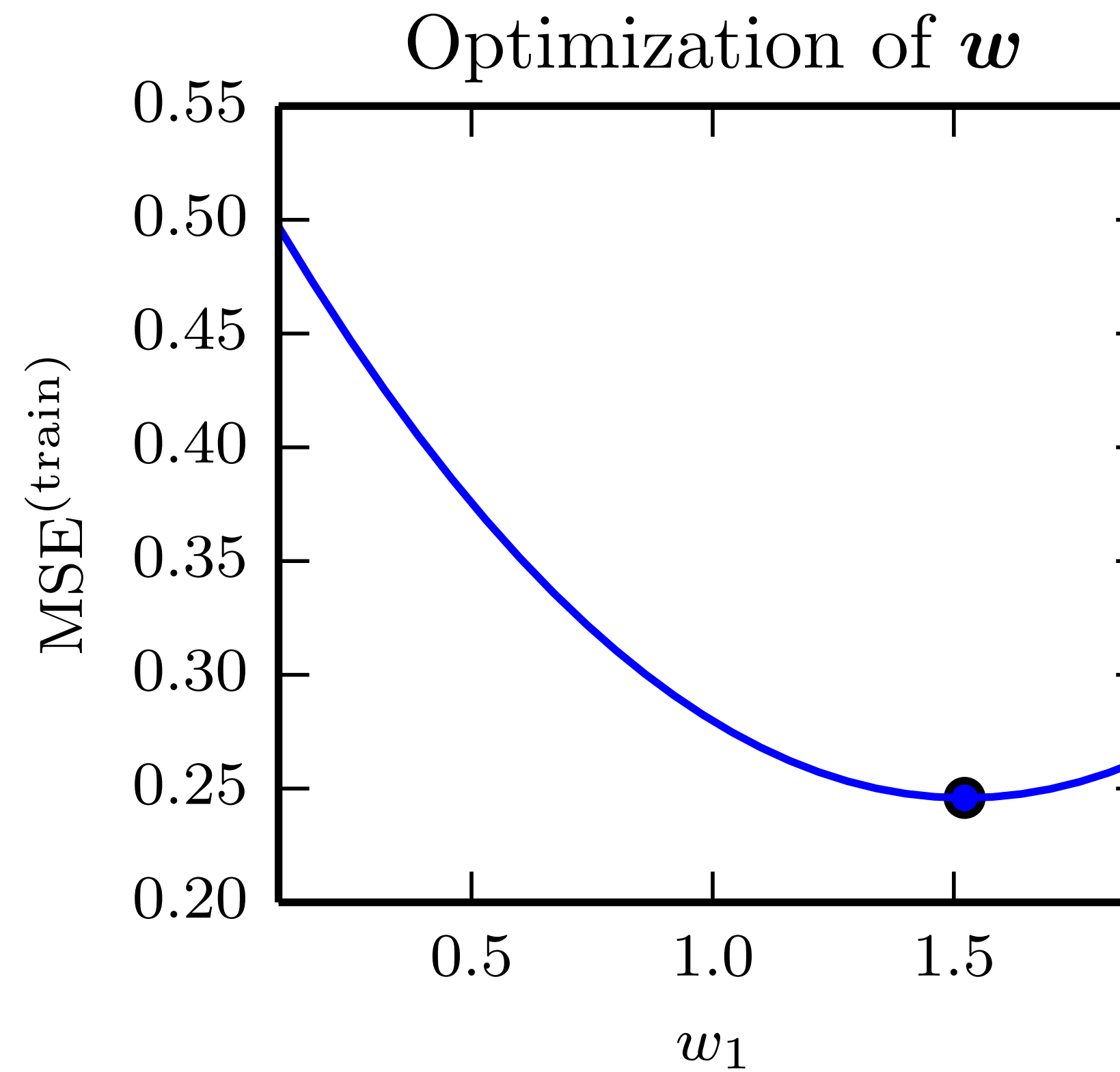
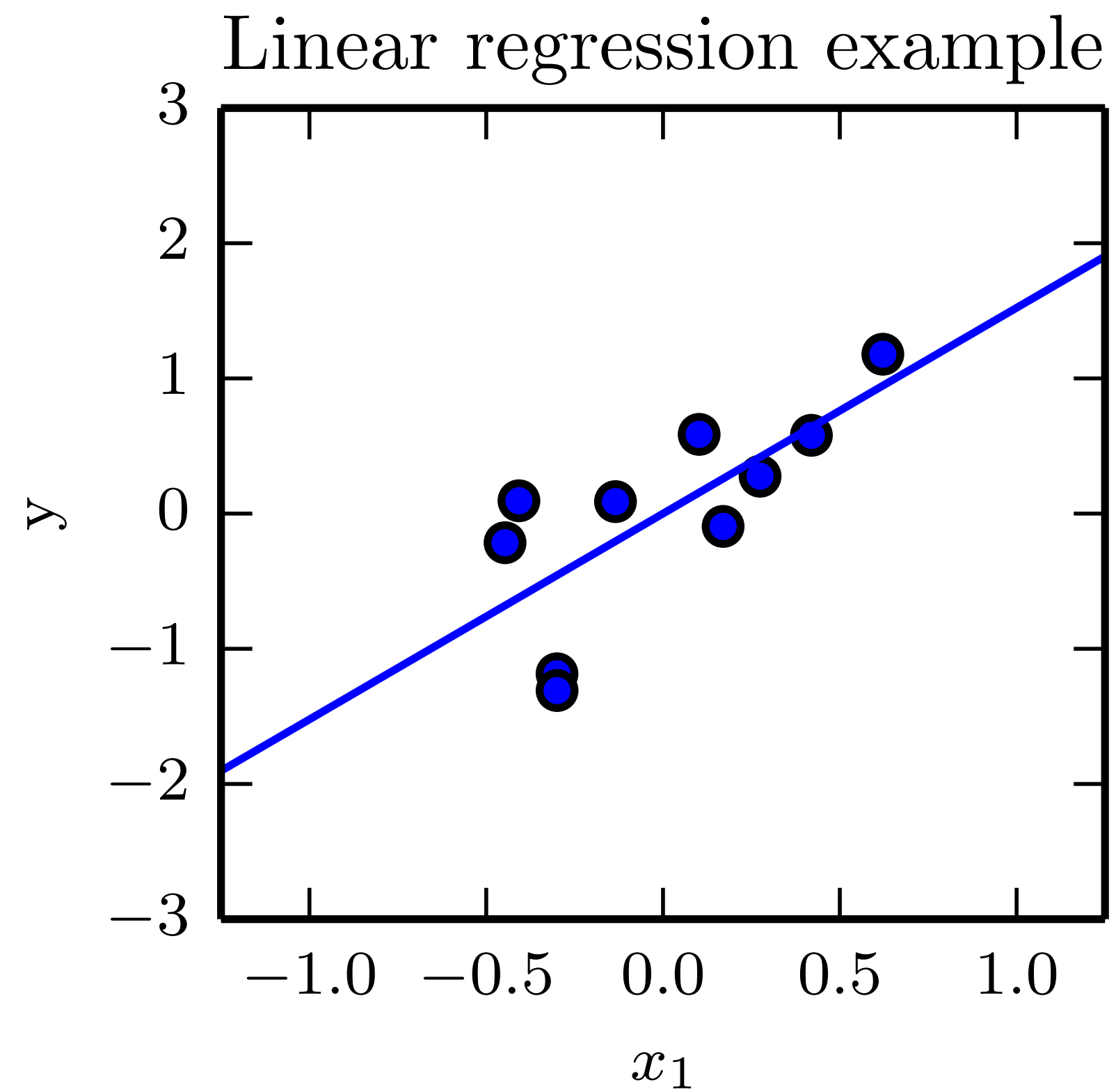
The Bias-Variance Trade-Off

B



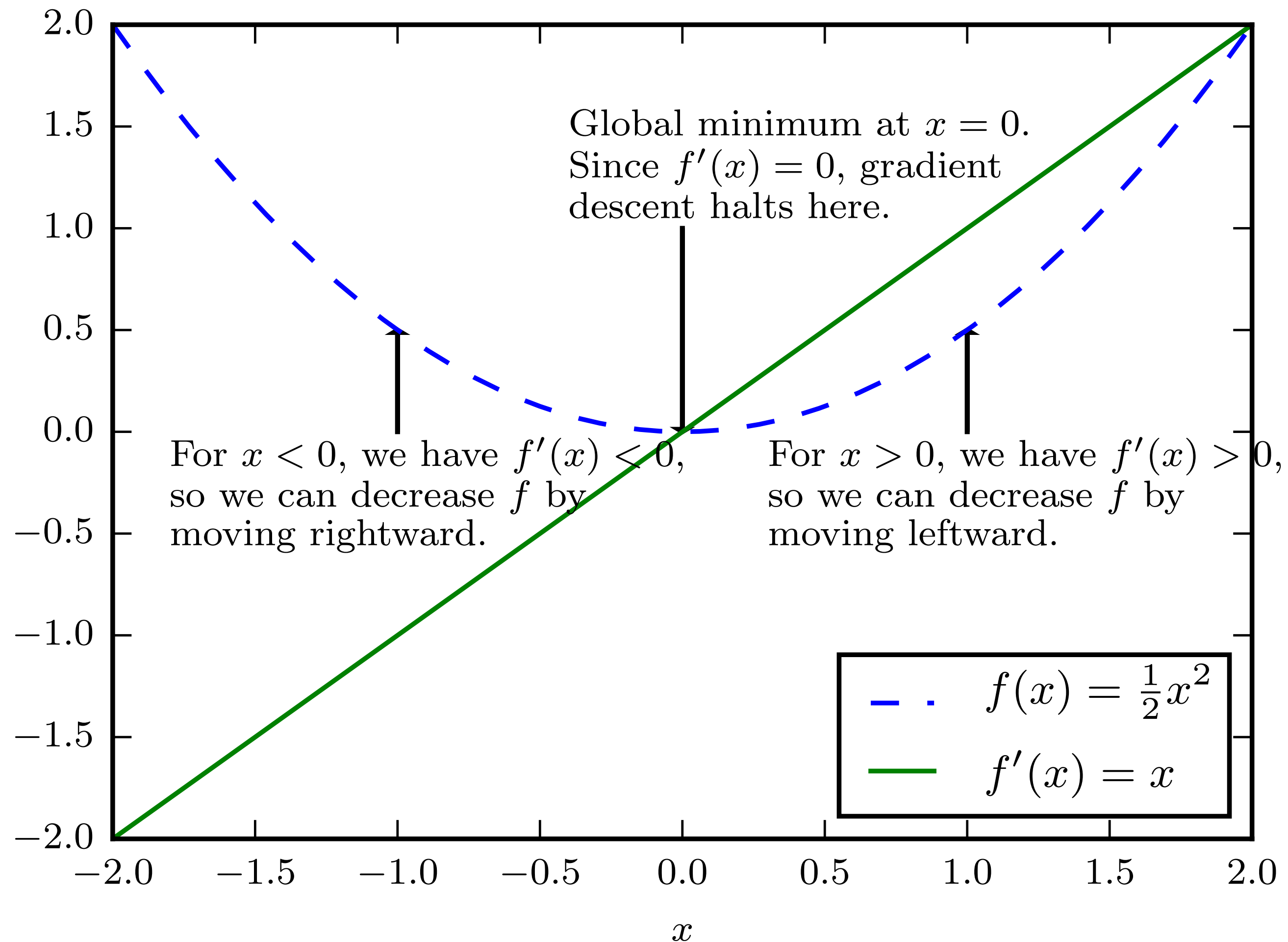
3 Gradient Descent

How do we optimize our model?



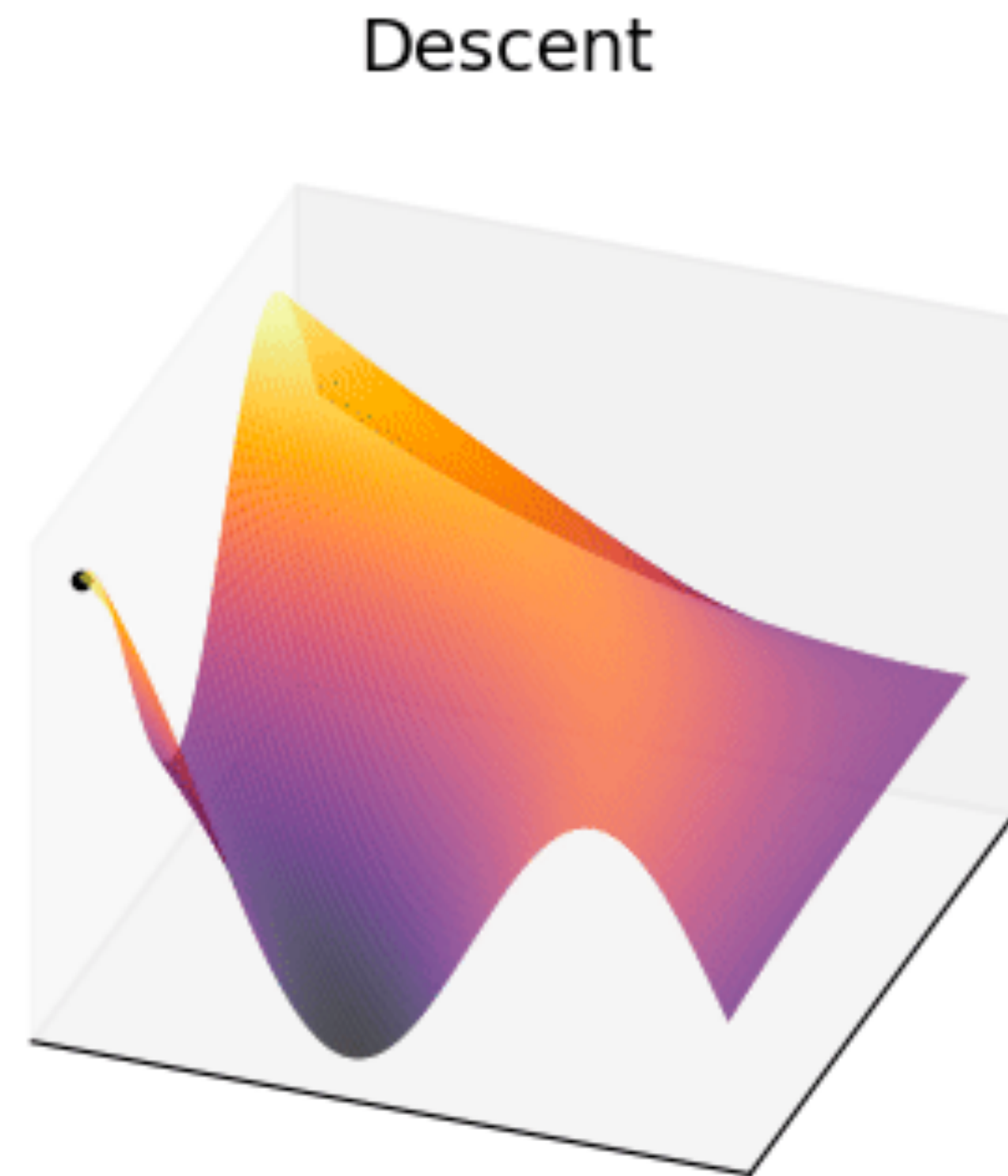
How do we optimize our model?

Follow the gradient



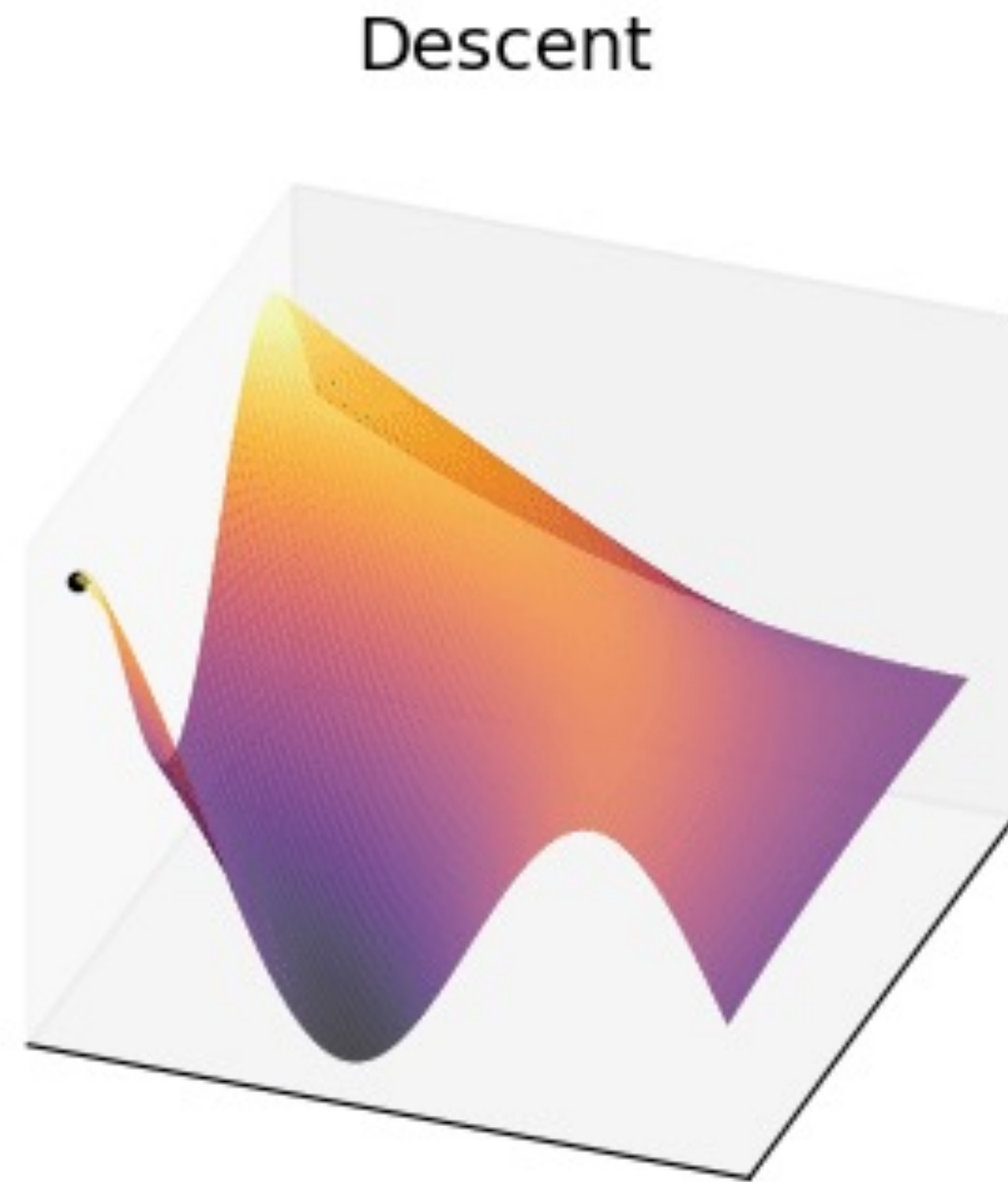
How can I imagine that?

Think about a ball rolling down the loss landscape



How can I imagine that?

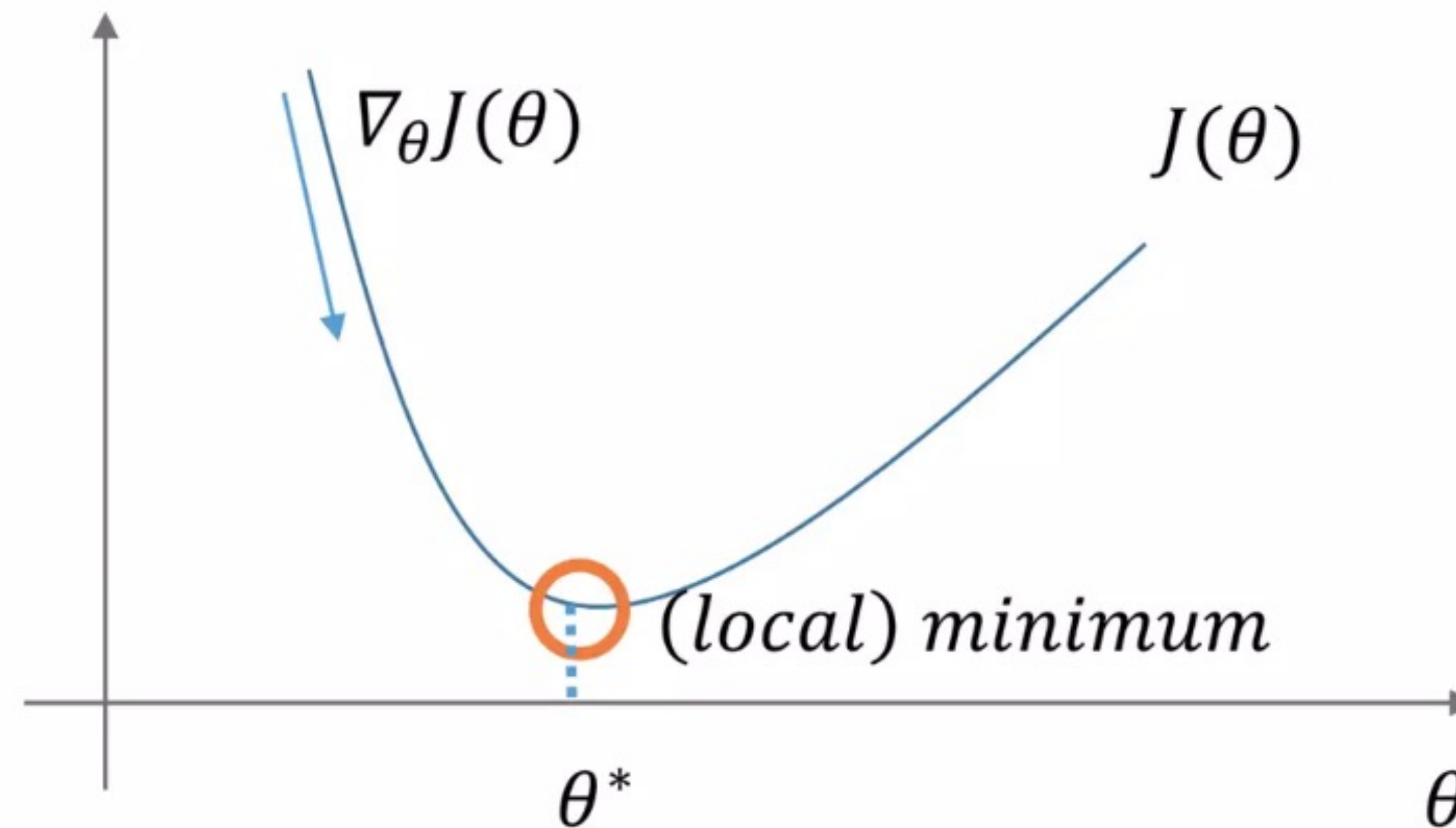
Think about a ball rolling down the loss landscape



Gradient Descent in formulas

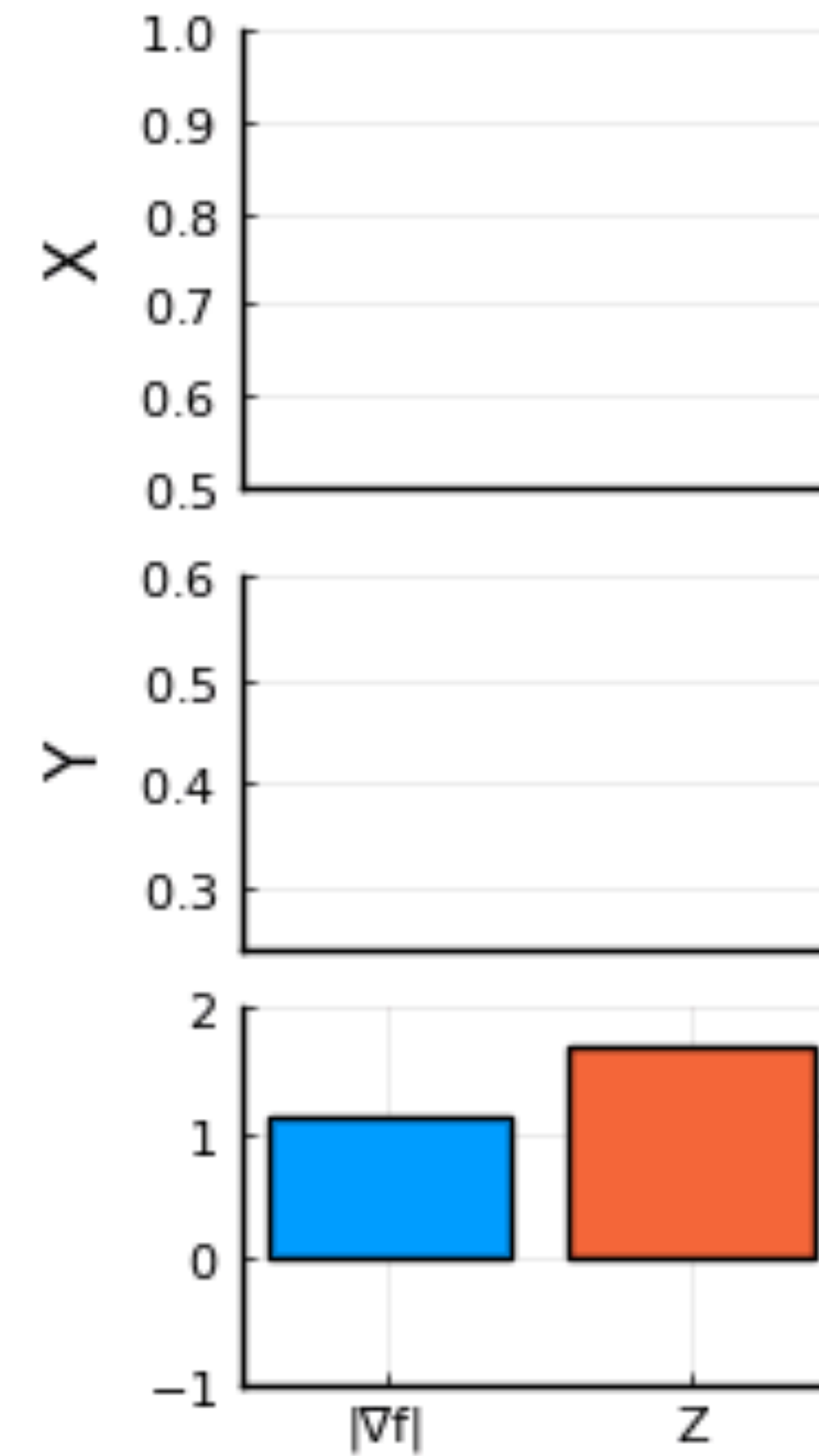
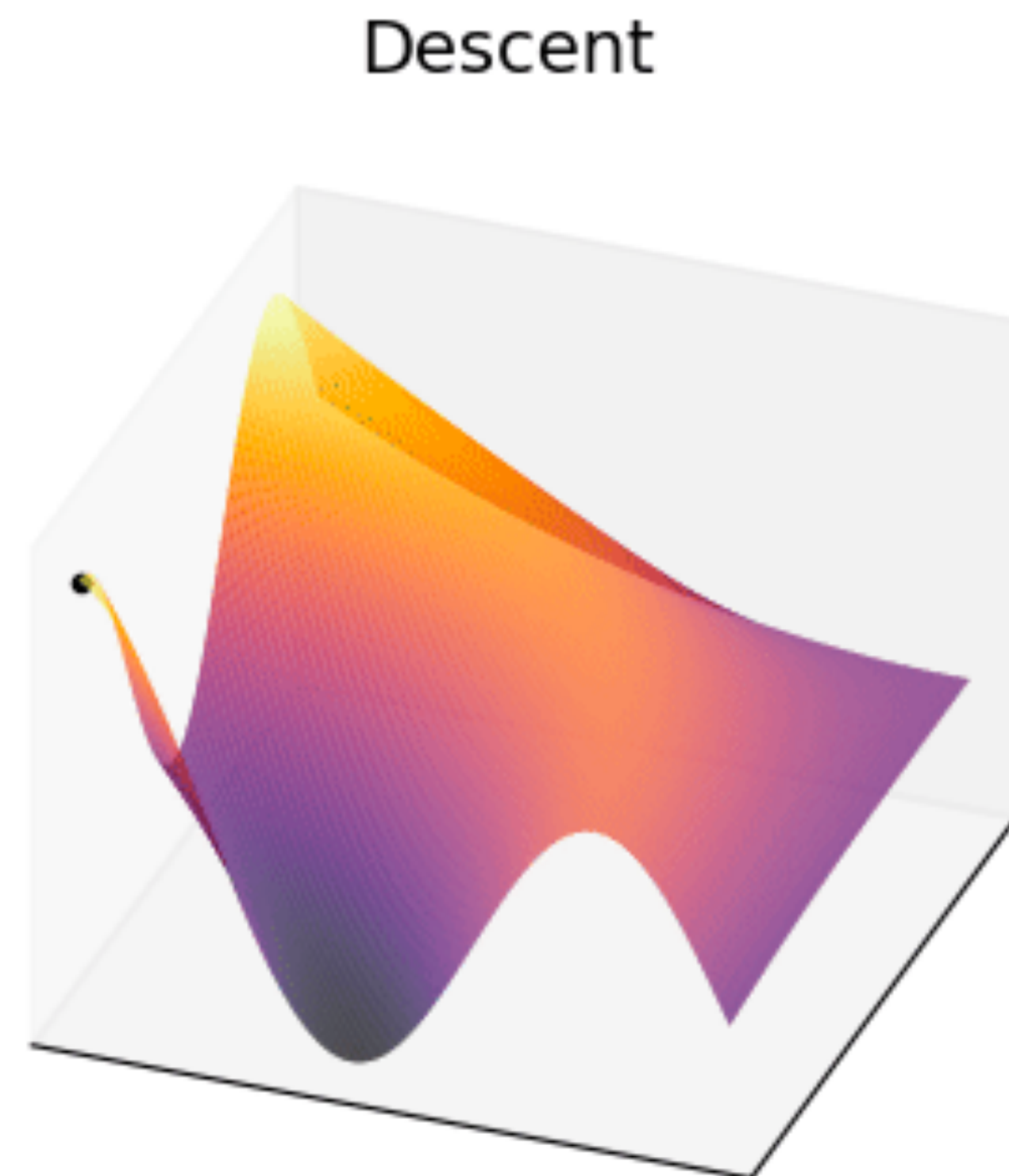
Going in the opposite direction

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta).$$



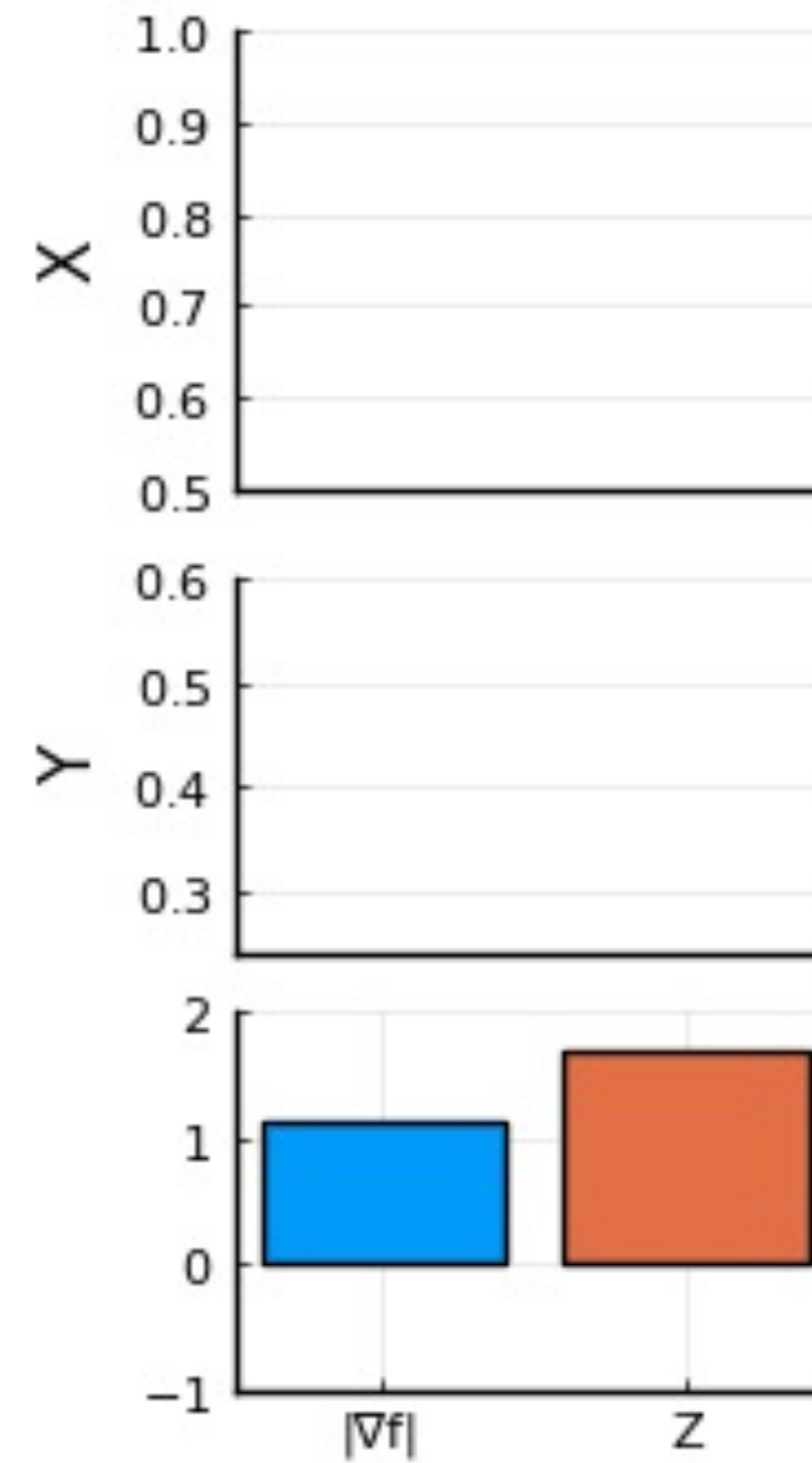
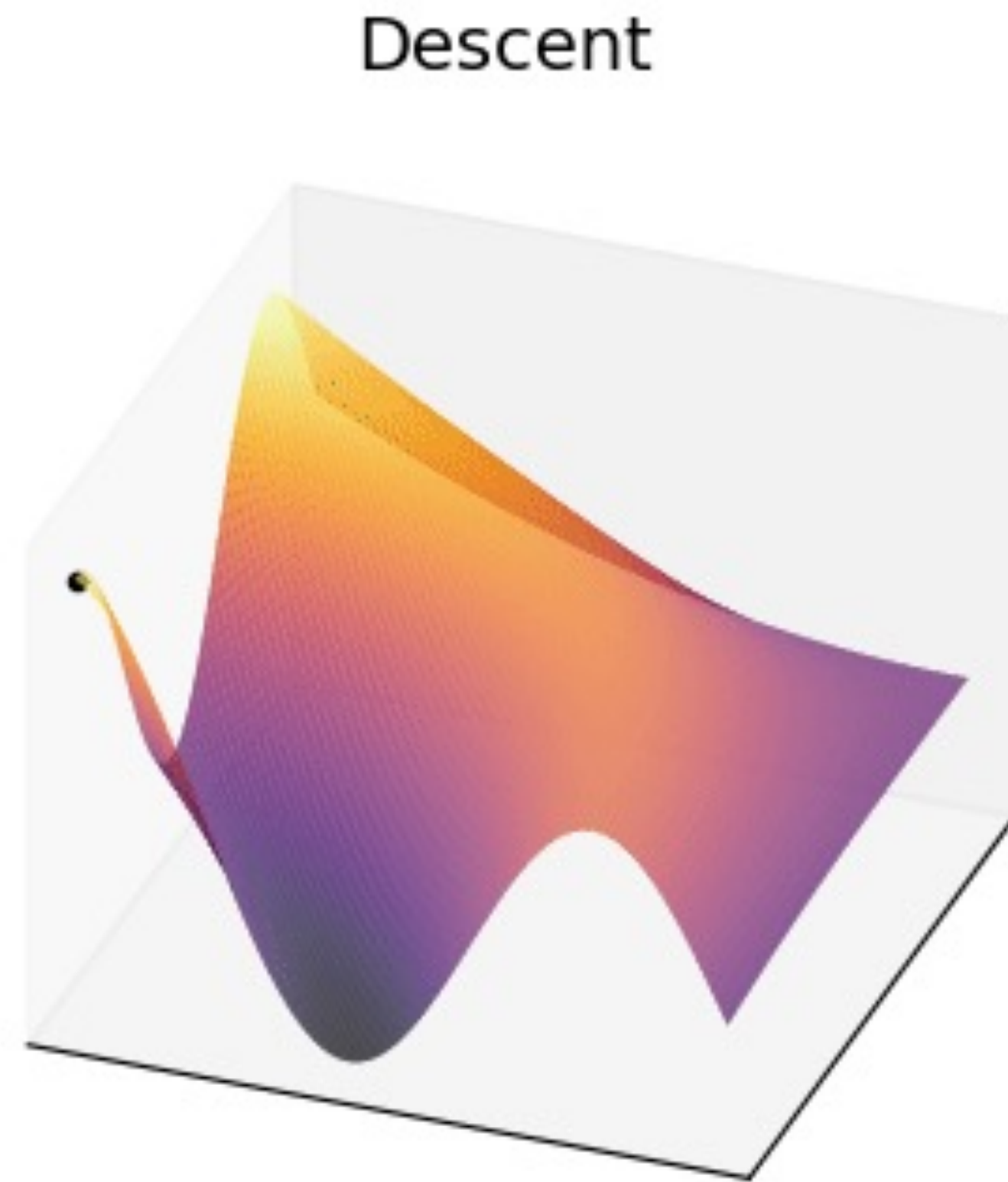
Going down the gradient

Local Minima can become a problem



Going down the gradient

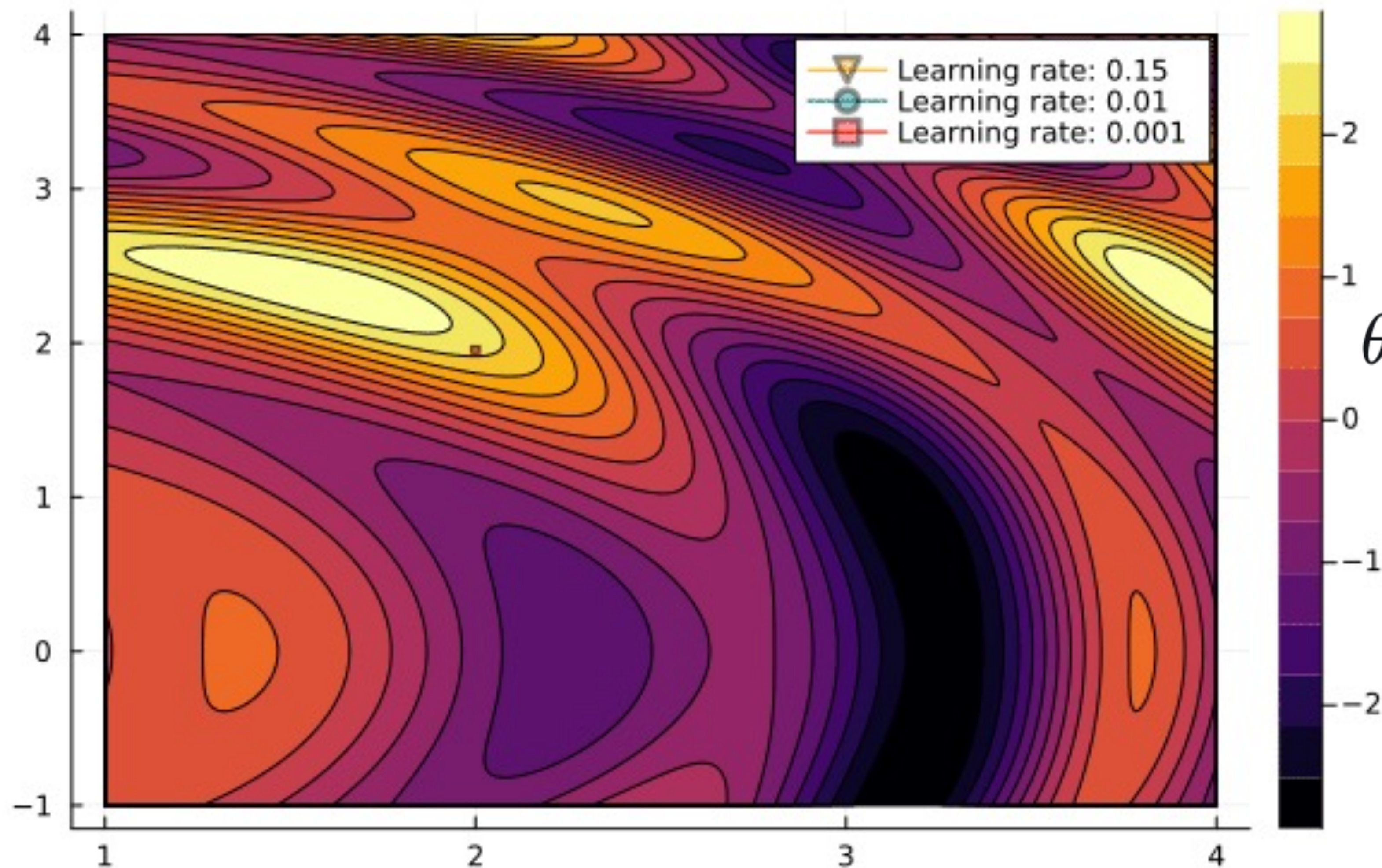
Local Minima can become a problem



Learning Rate

The first hyperparameter you should check if trouble arises

Iterations=1



$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta).$$

Vanilla Gradient Descent is slow

Speed it up by only looking at one data point at a time

Vanilla GD

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta).$$



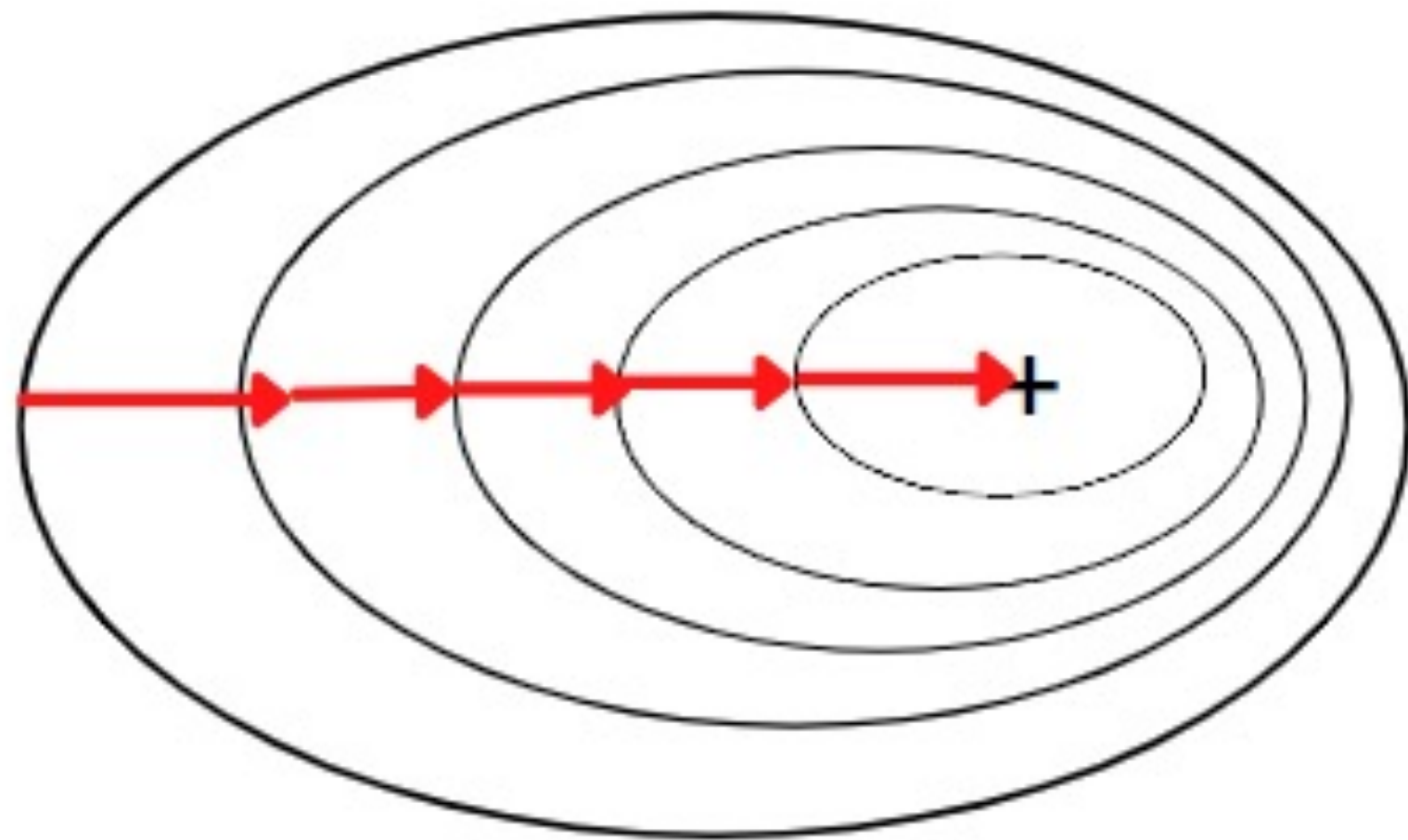
Stochastic GD

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)}).$$

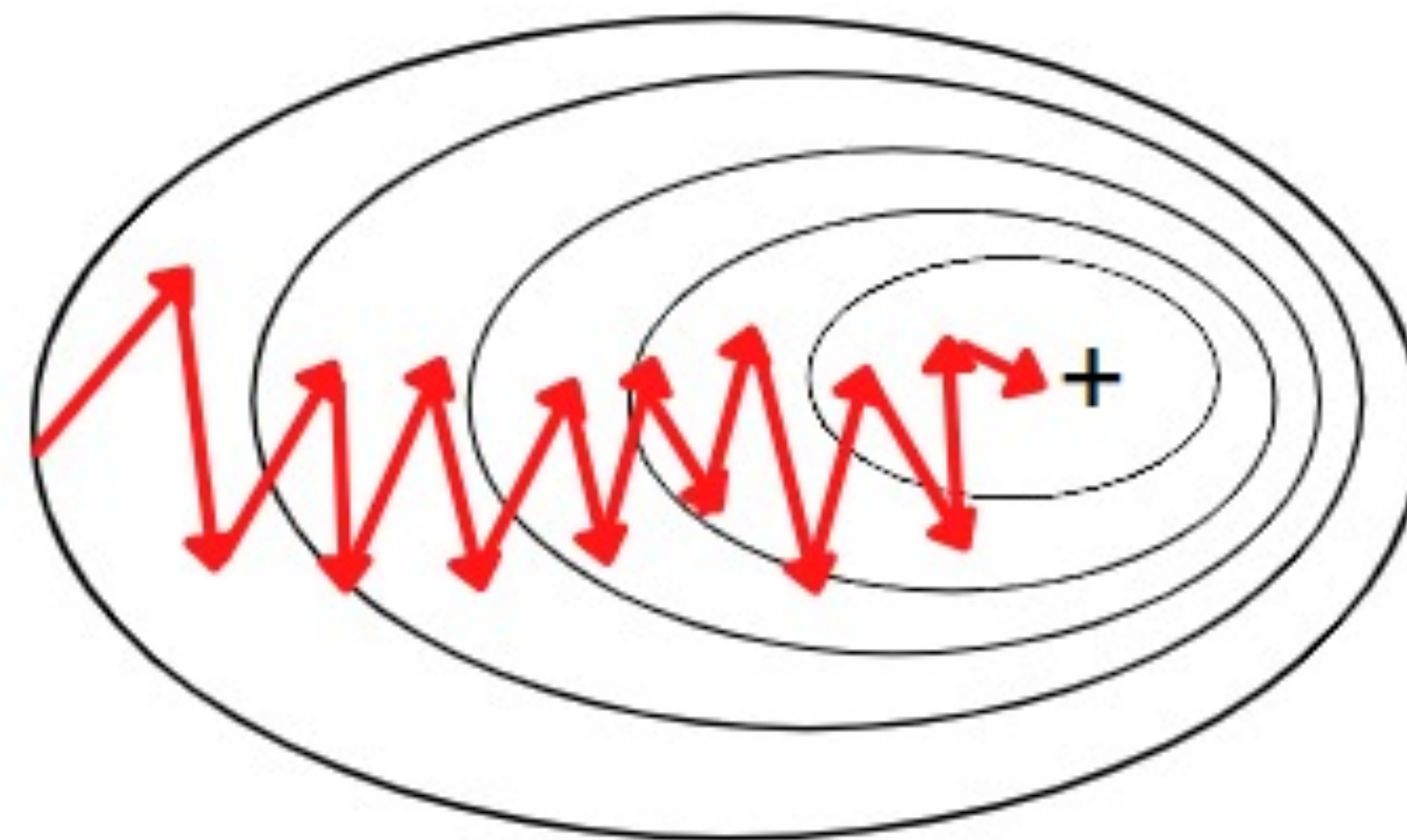
How can I imagine that?

Planned and cautious versus spontaneous and chaotic

Batch Gradient Descent



Stochastic Gradient Descent



Mini-batch GD: The best of both worlds

Only look at a subset of your data for each update step

Stochastic GD

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)}).$$



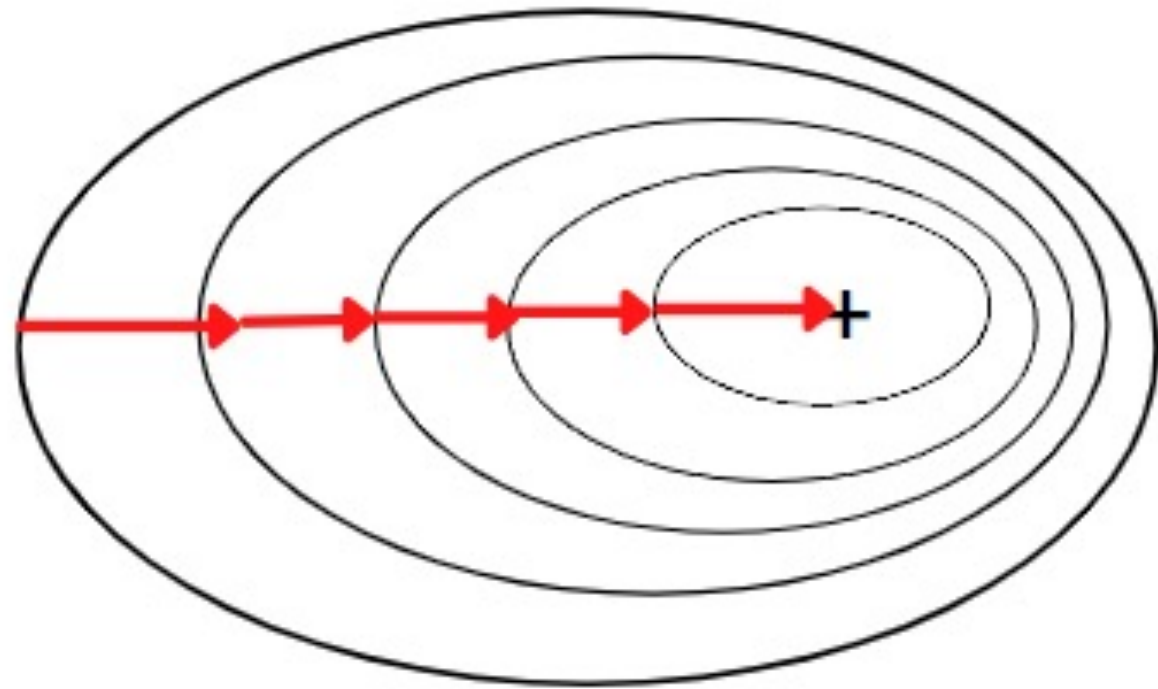
Mini-batch GD

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i:i+n)}; y^{(i:i+n)})$$

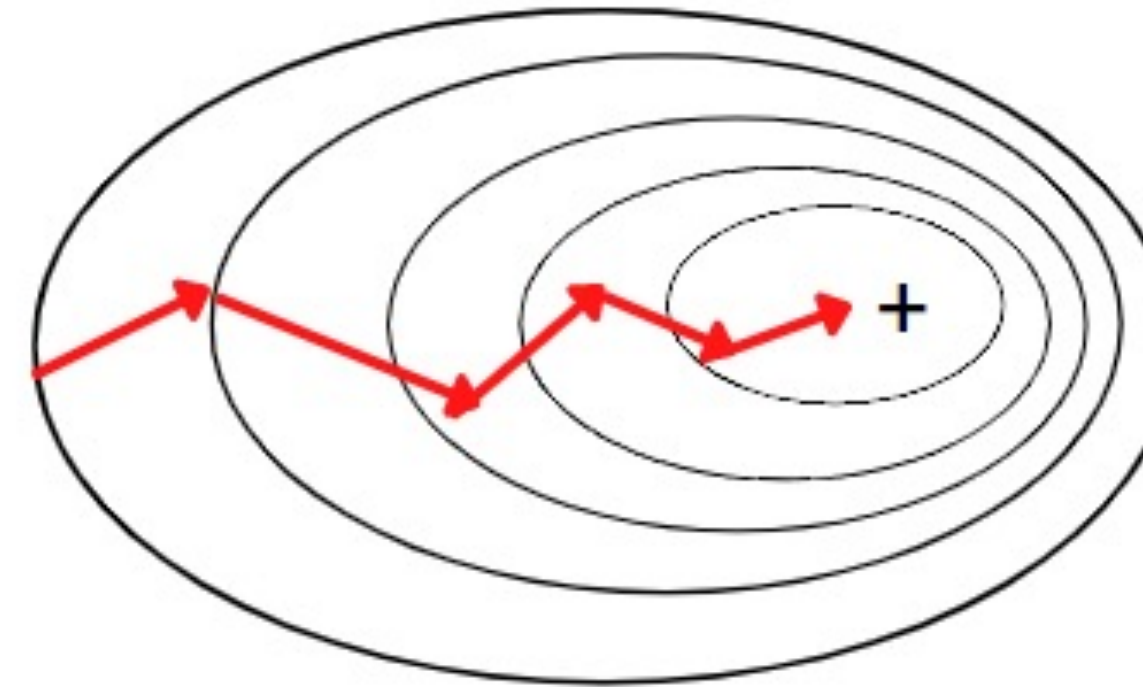
Mini-batch GD: The best of both worlds

Only look at a subset of your data for each update step

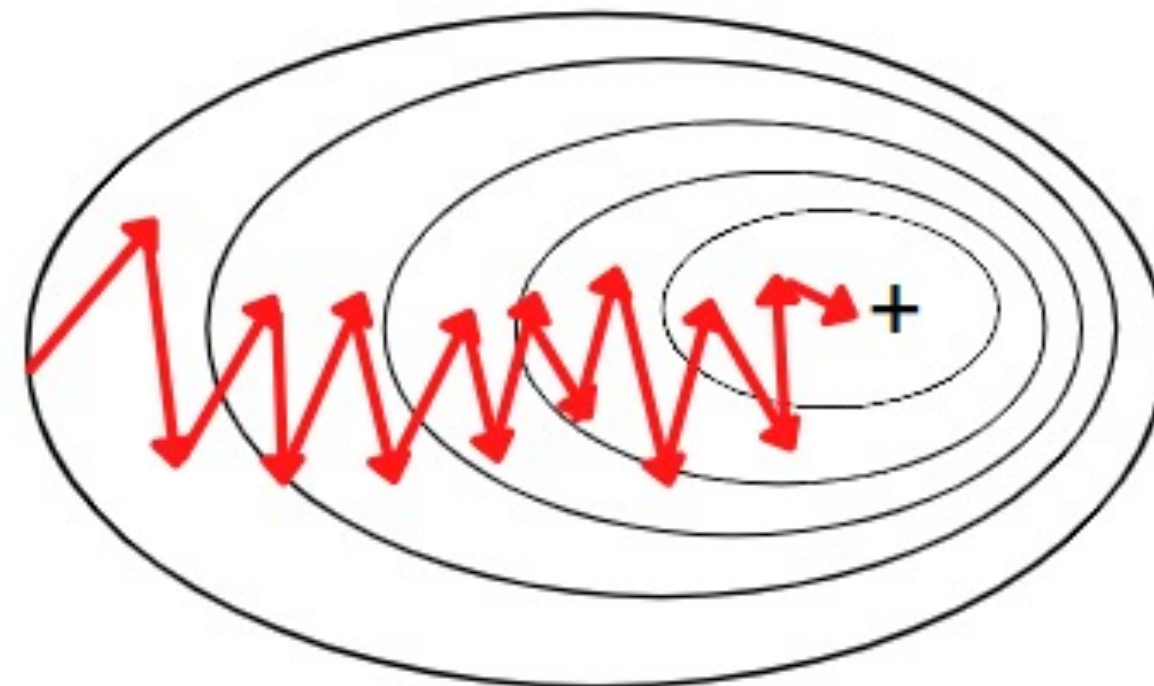
Batch Gradient Descent



Mini-Batch Gradient Descent



Stochastic Gradient Descent



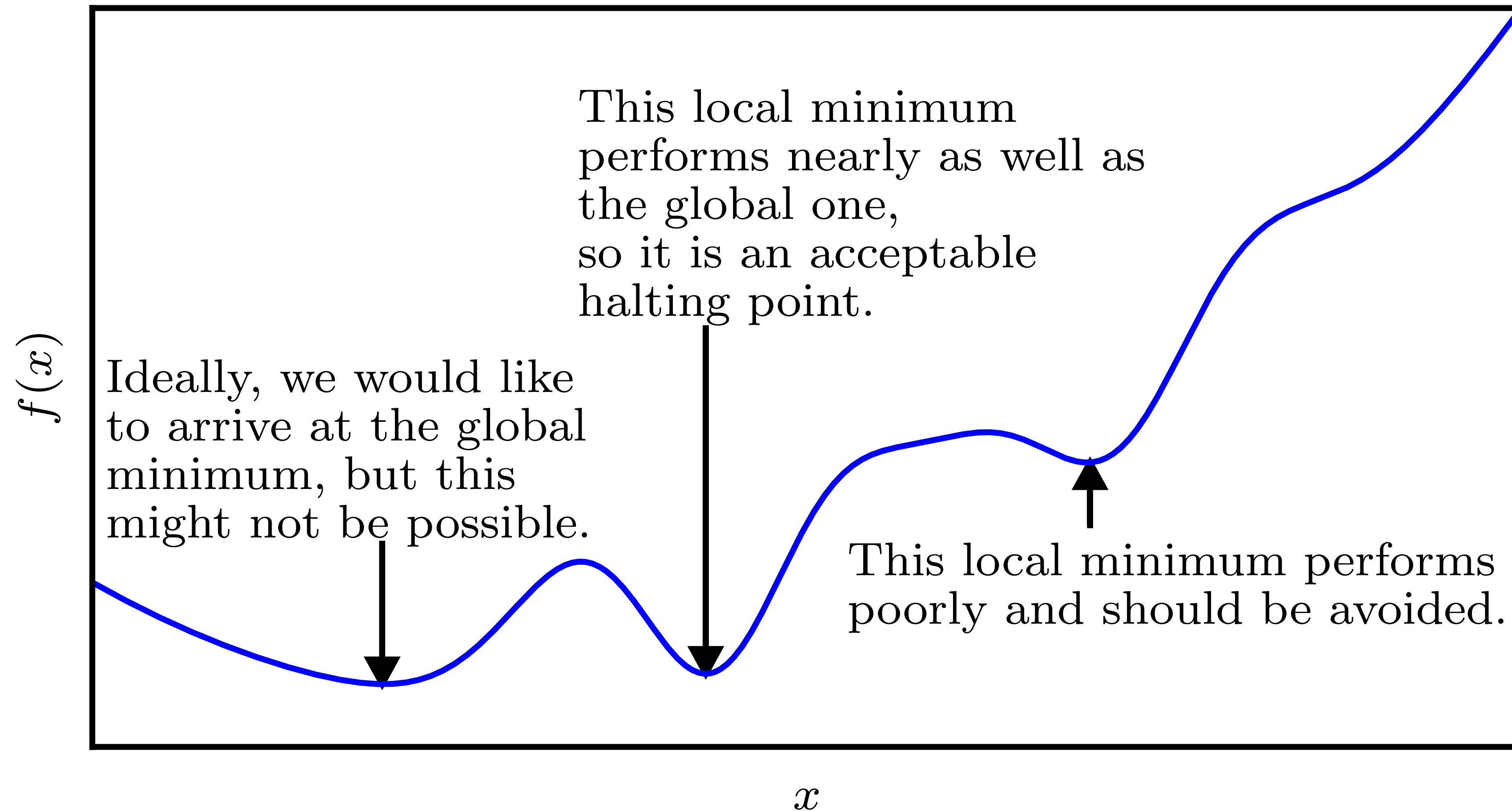
Which GD to choose?

In practice, mini-batch is often a good choice

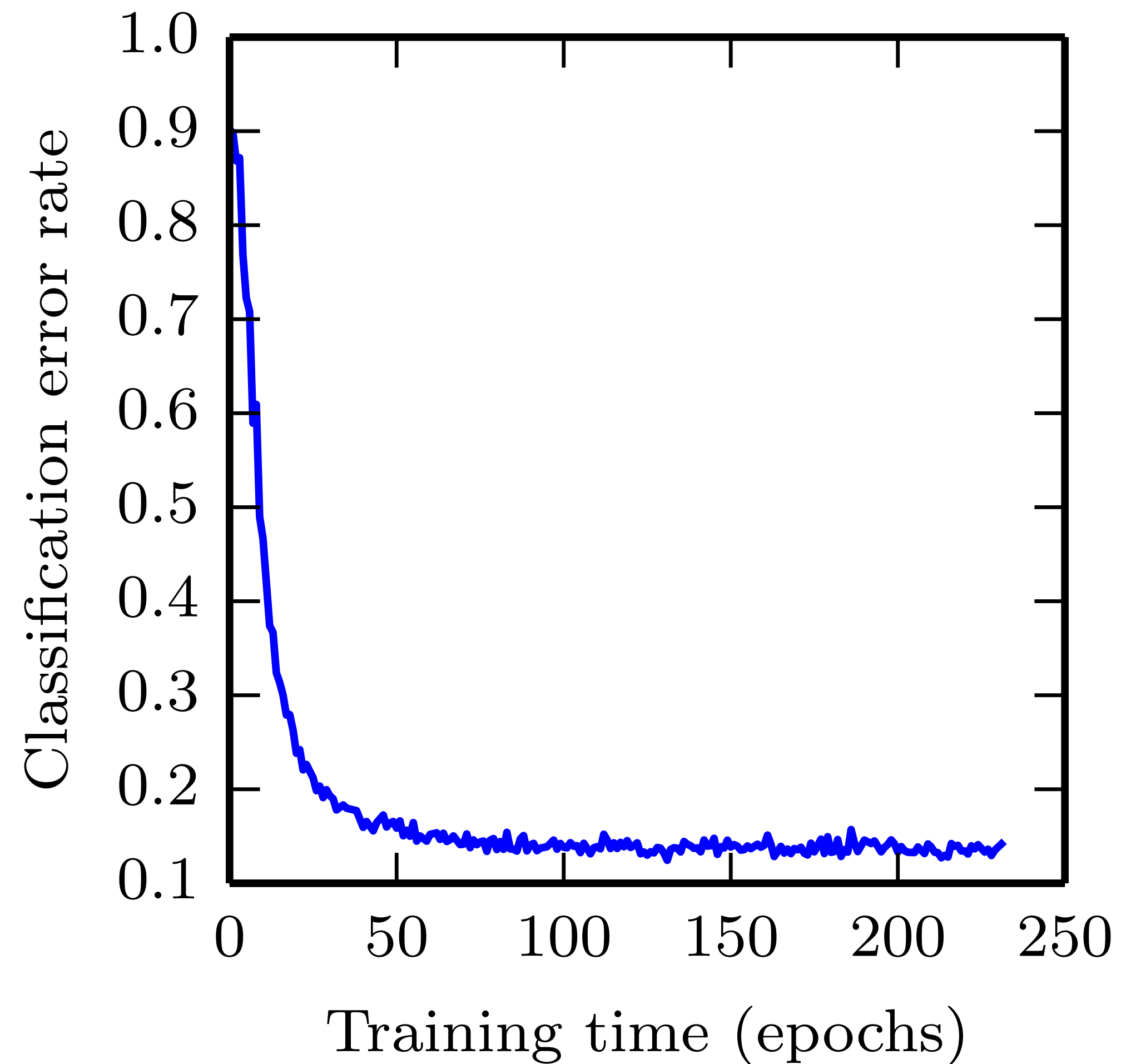
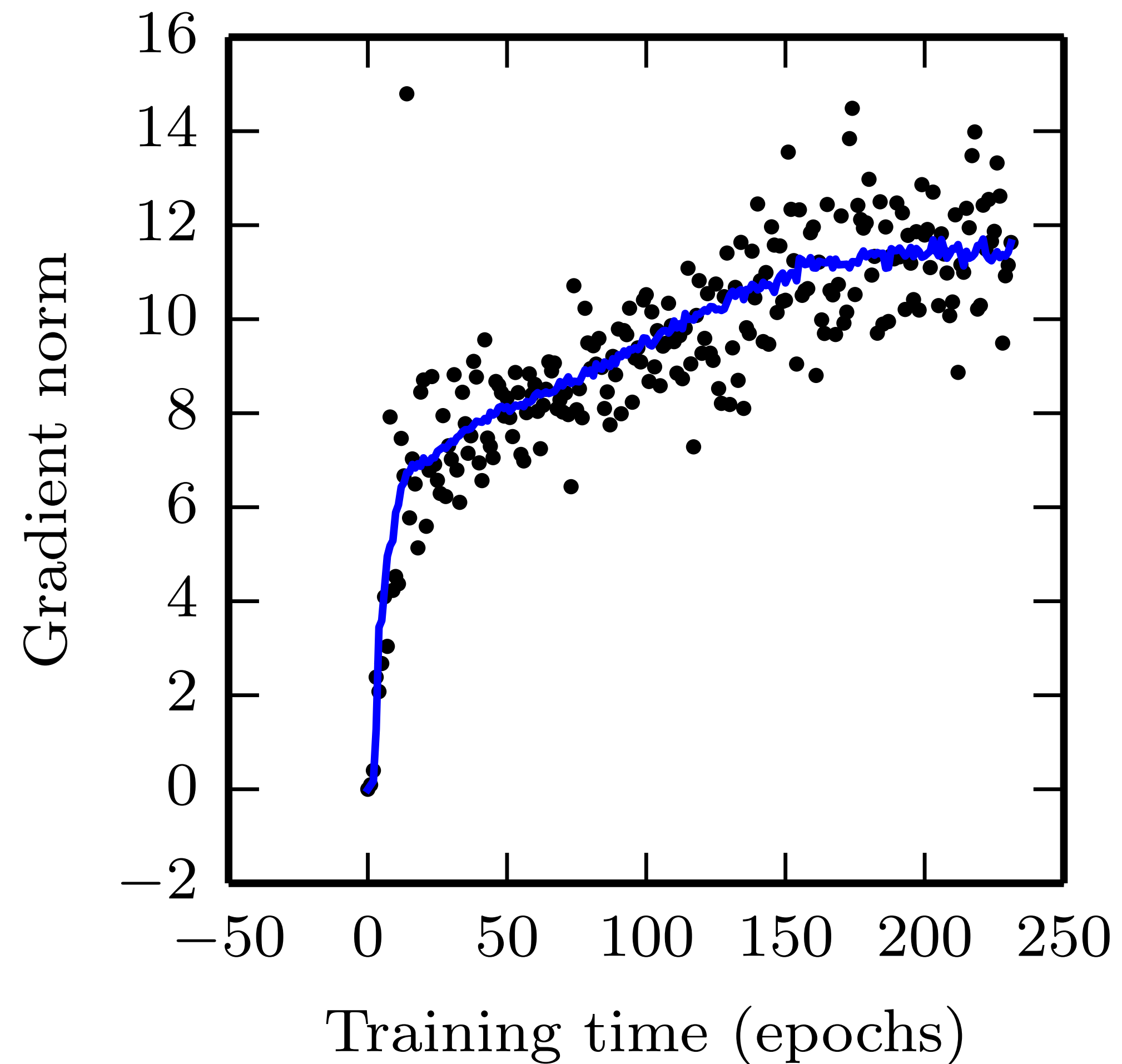
Method	Accuracy	Update Speed	Memory Usage	Online Learning
Batch gradient descent	Good	Slow	High	No
Stochastic gradient descent	Good (with annealing)	High	Low	Yes
Mini-batch gradient descent	Good	Medium	Medium	Yes

Real-life optimisation is hard

We do not expect to find the global minimum in most cases

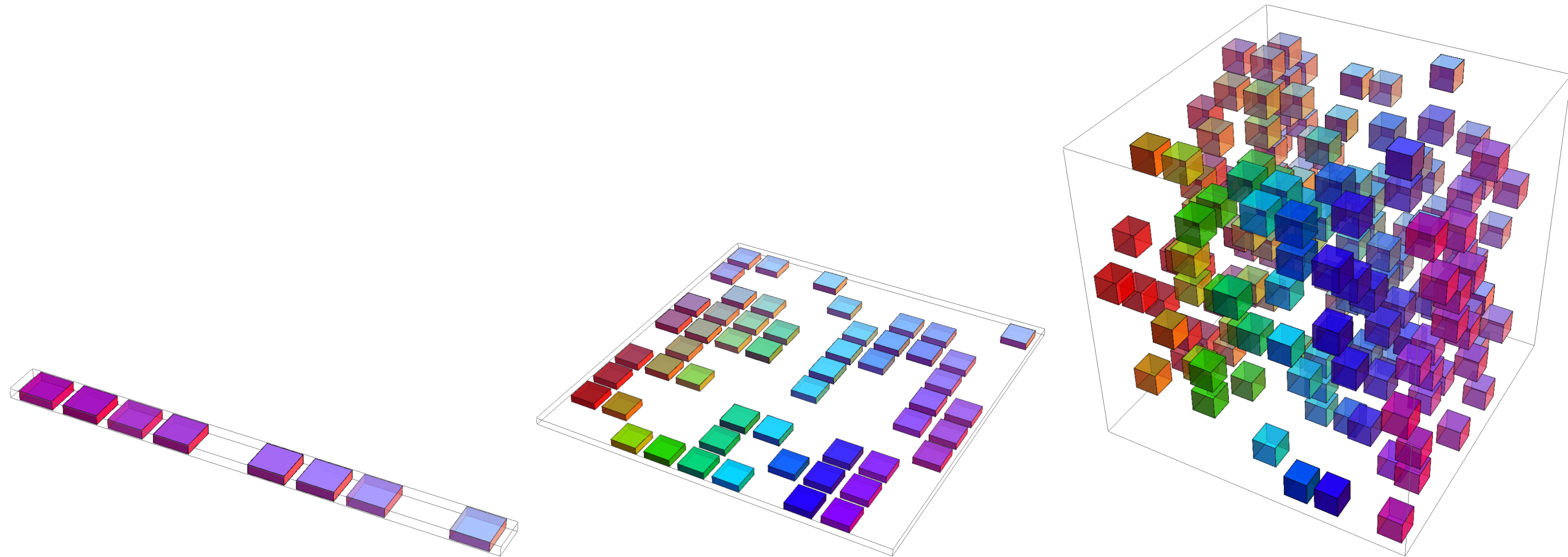


We usually don't even reach a local minimum



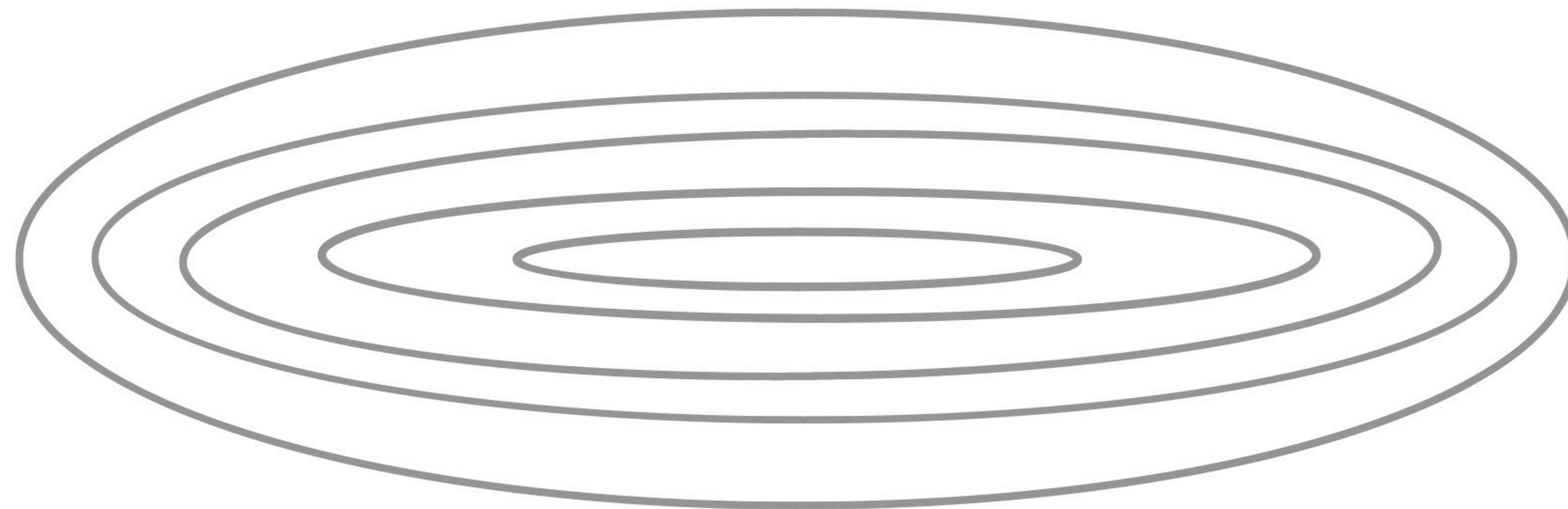
Curse of Dimensionality

The higher-dimensional the data, the more we need of it



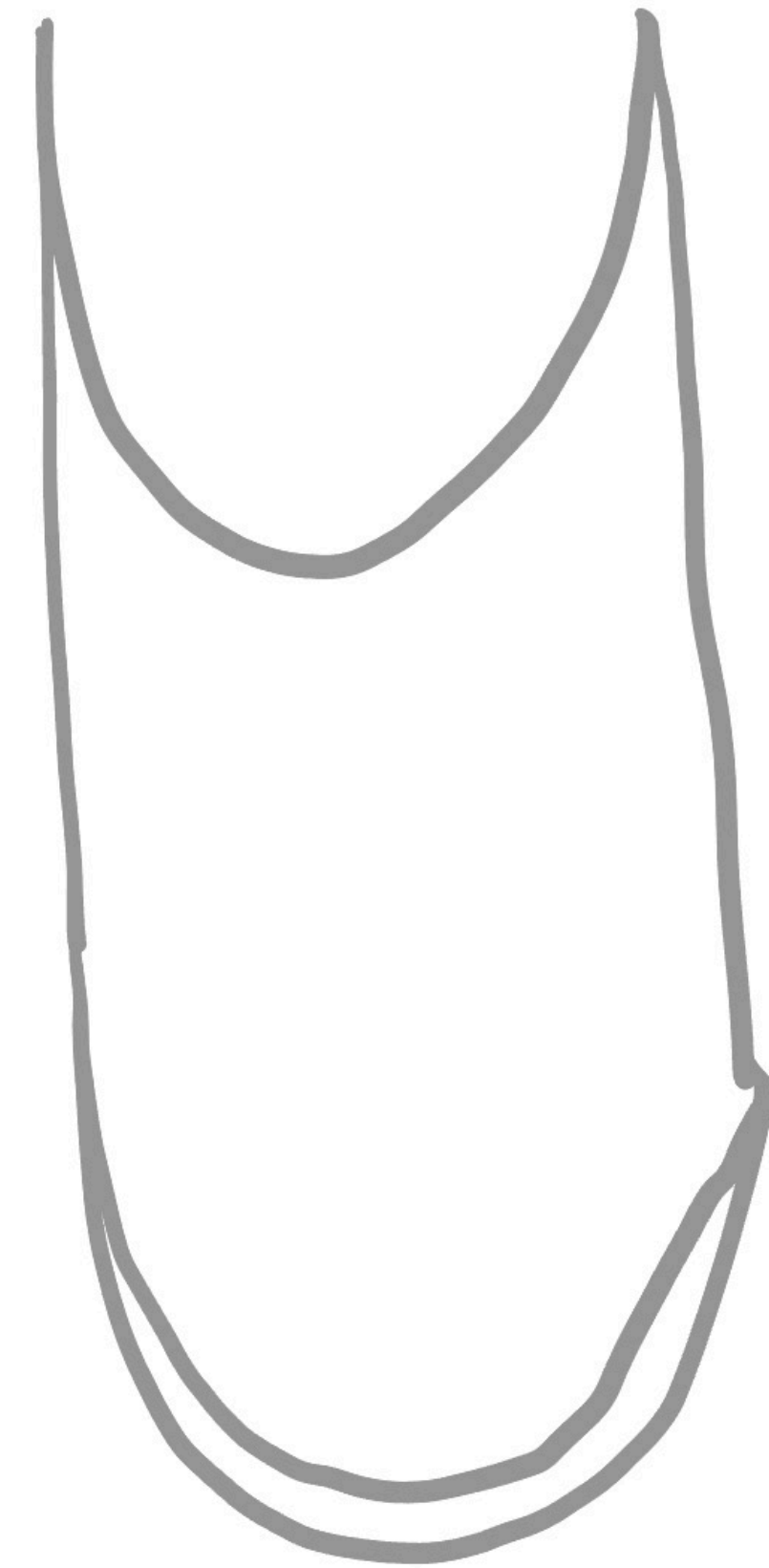
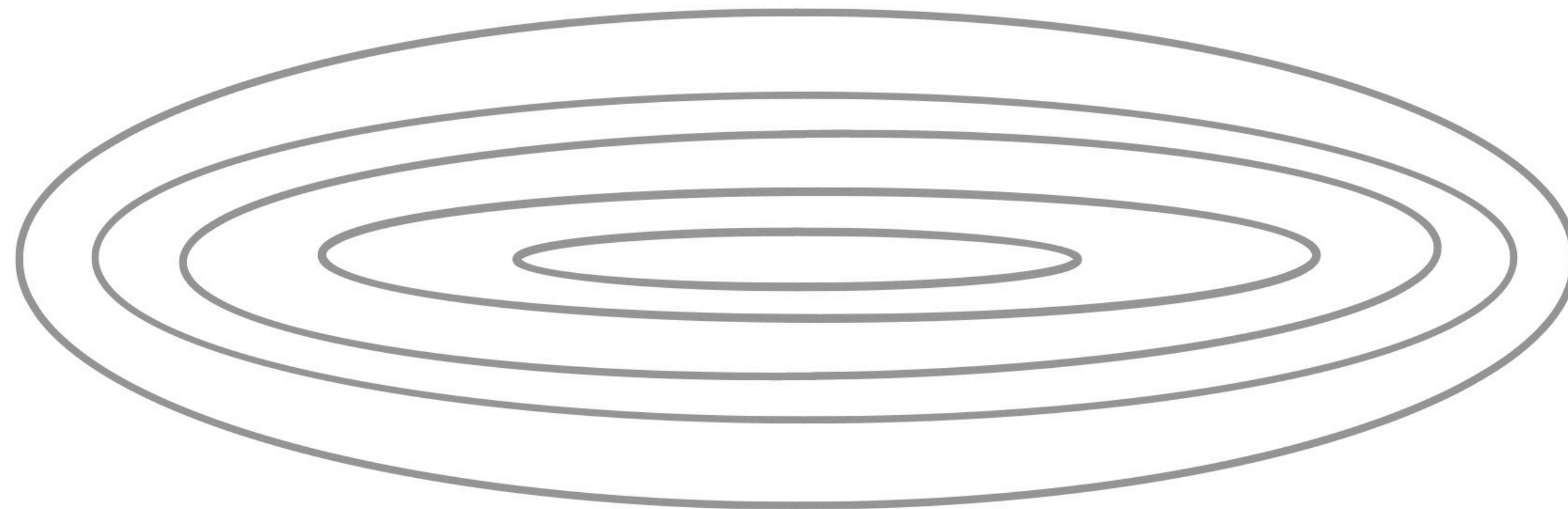
Momentum Optimisers

Have a memory of the past to overcome dire times



Momentum Optimisers

Have a memory of the past to overcome dire times



Momentum Optimisers

Have a memory of the past to overcome dire times

Vanilla GD

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta).$$

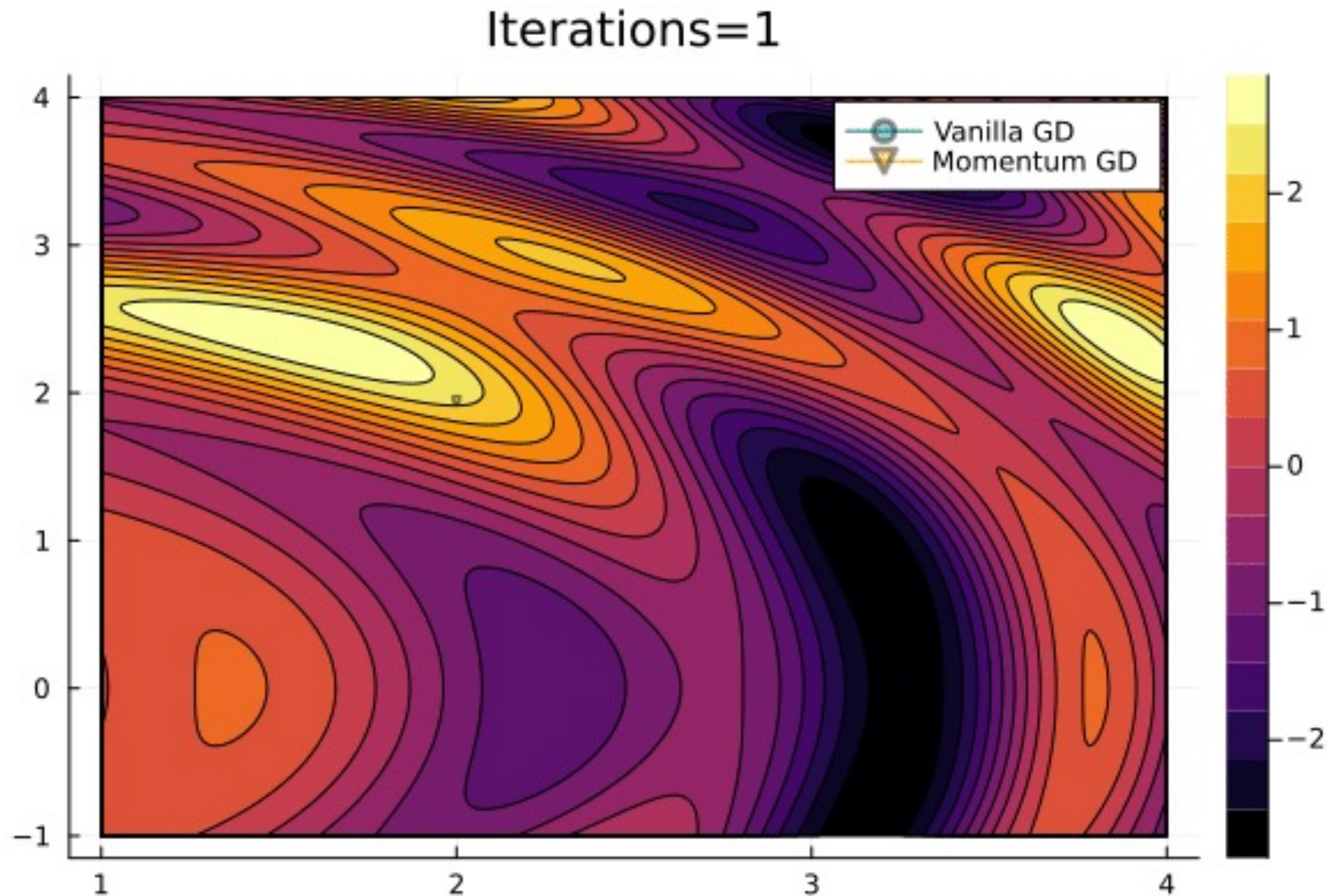


Momentum GD

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$
$$\theta = \theta - v_t$$

Momentum

Speeding things up when the gradient becomes shallow



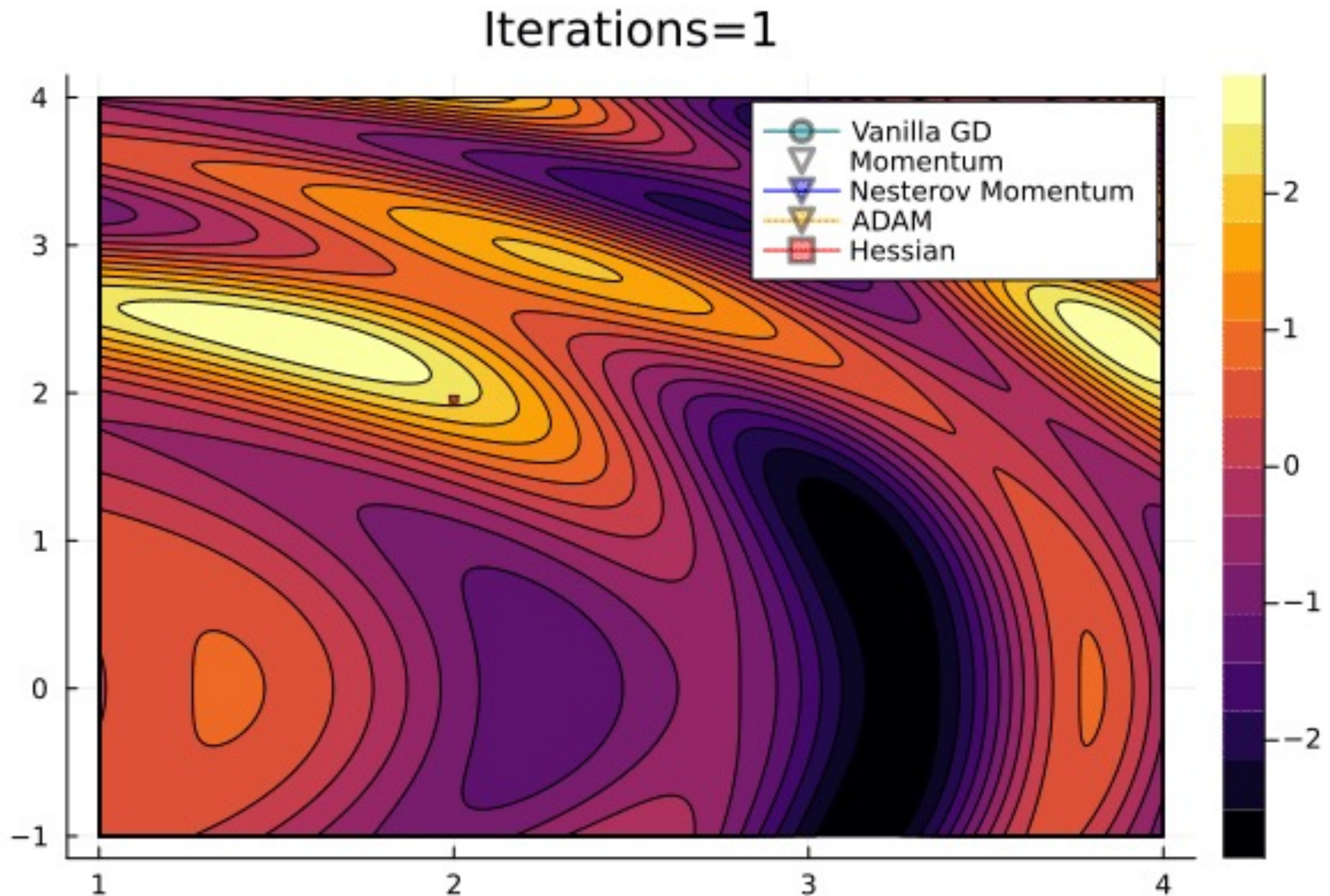
A Zoo of Momentum Methods...

ADAM is the standard default choice

Method	Update equation
SGD	$g_t = \nabla_{\theta_t} J(\theta_t)$ $\Delta\theta_t = -\eta \cdot g_t$ $\theta_t = \theta_t + \Delta\theta_t$
Momentum	$\Delta\theta_t = -\gamma v_{t-1} - \eta g_t$
NAG	$\Delta\theta_t = -\gamma v_{t-1} - \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$
Adagrad	$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$
Adadelta	$\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$
RMSprop	$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$
Adam	$\Delta\theta_t = -\frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$

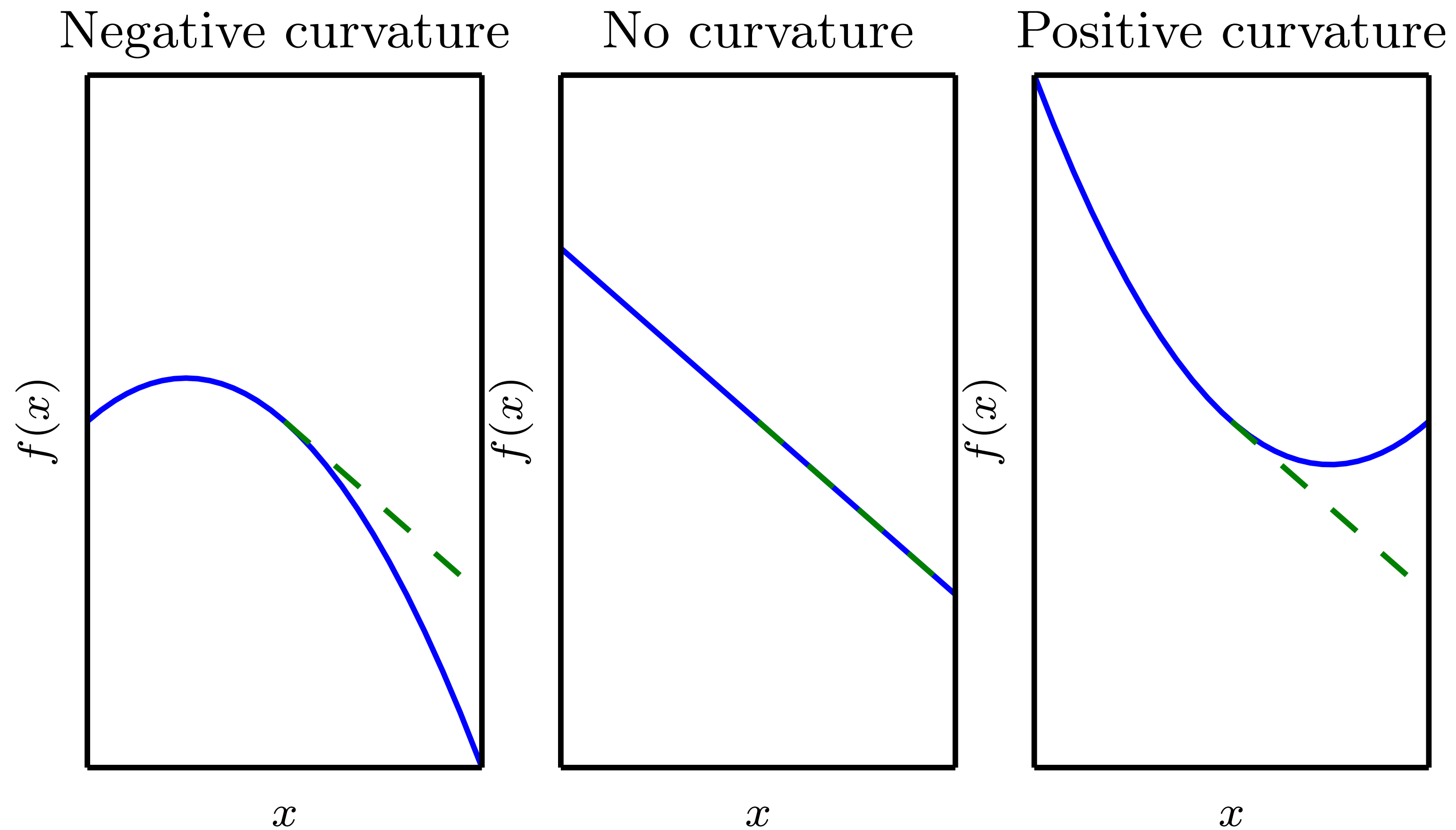
A Zoo of Momentum Methods...

ADAM is the standard default choice



Back to learning rate: Is there an optimum?

Curvature helps out



Newton's method: Using curvature

The Hessian tells you your optimal step size

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}. \quad (4.9)$$

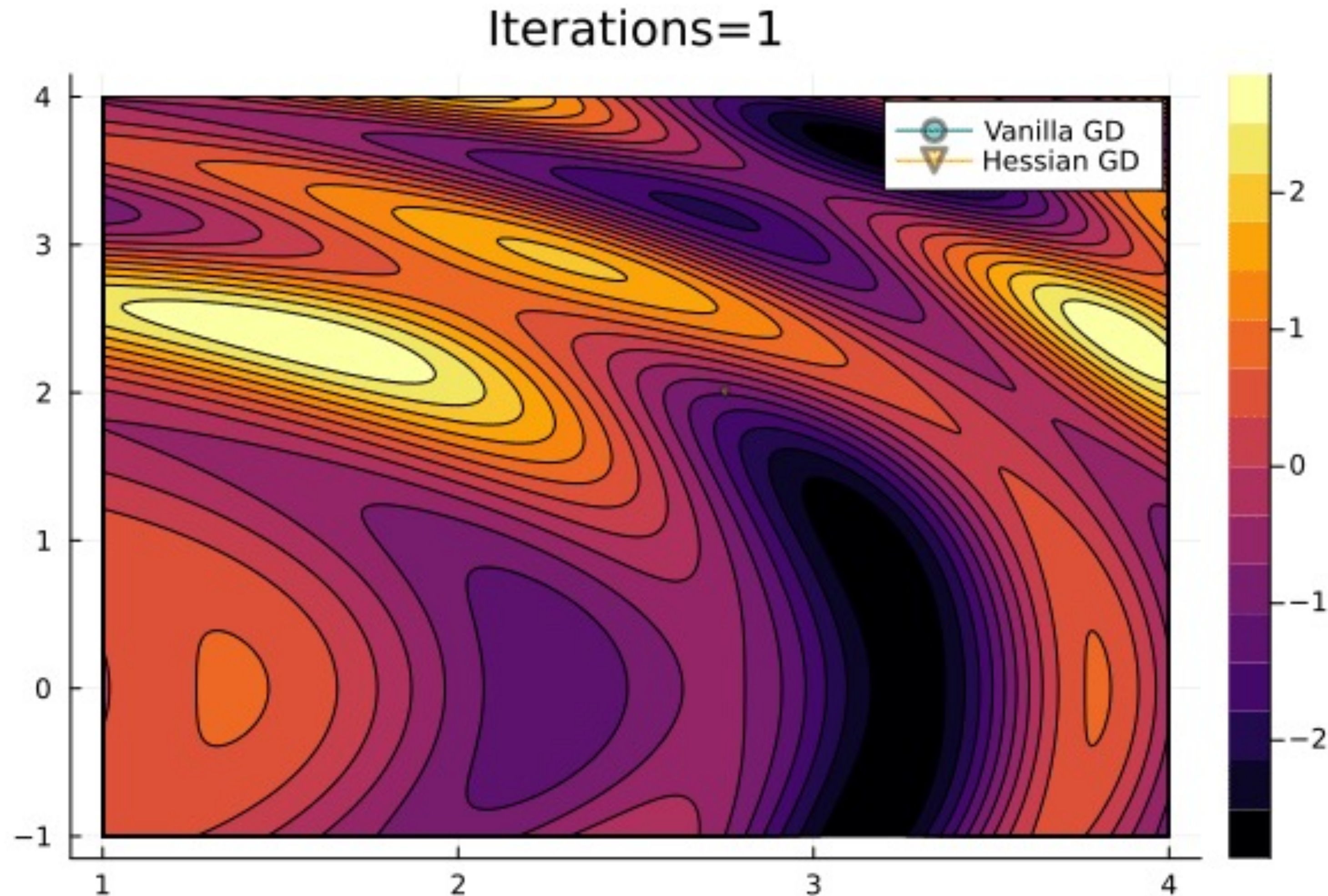
$$\epsilon^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}. \quad (4.10)$$

Big gradients speed you up

Big eigenvalues slow you down if you align with their eigenvectors

Hessian

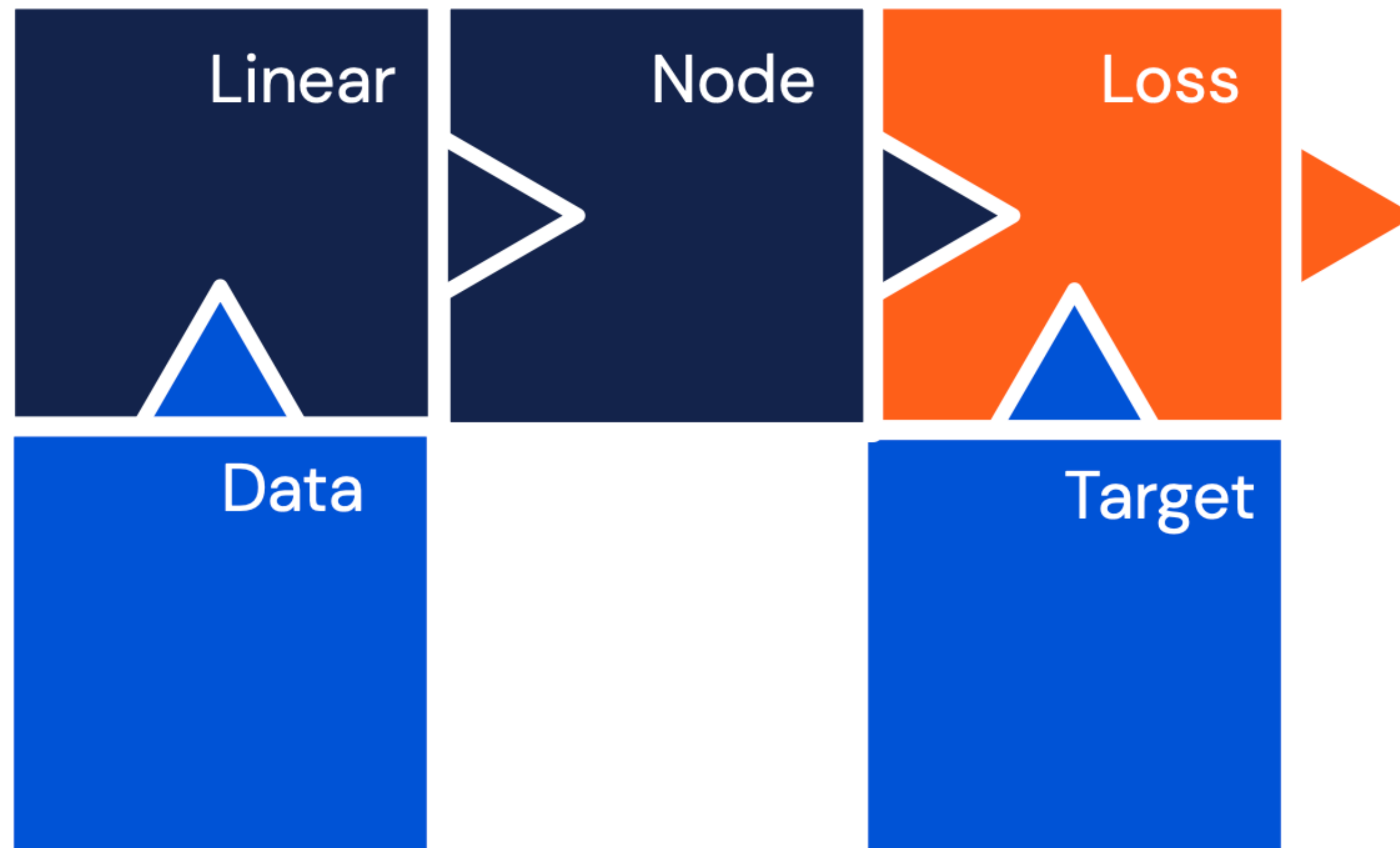
Each step is better, but evaluating each one is very expensive



4 Deep Learning

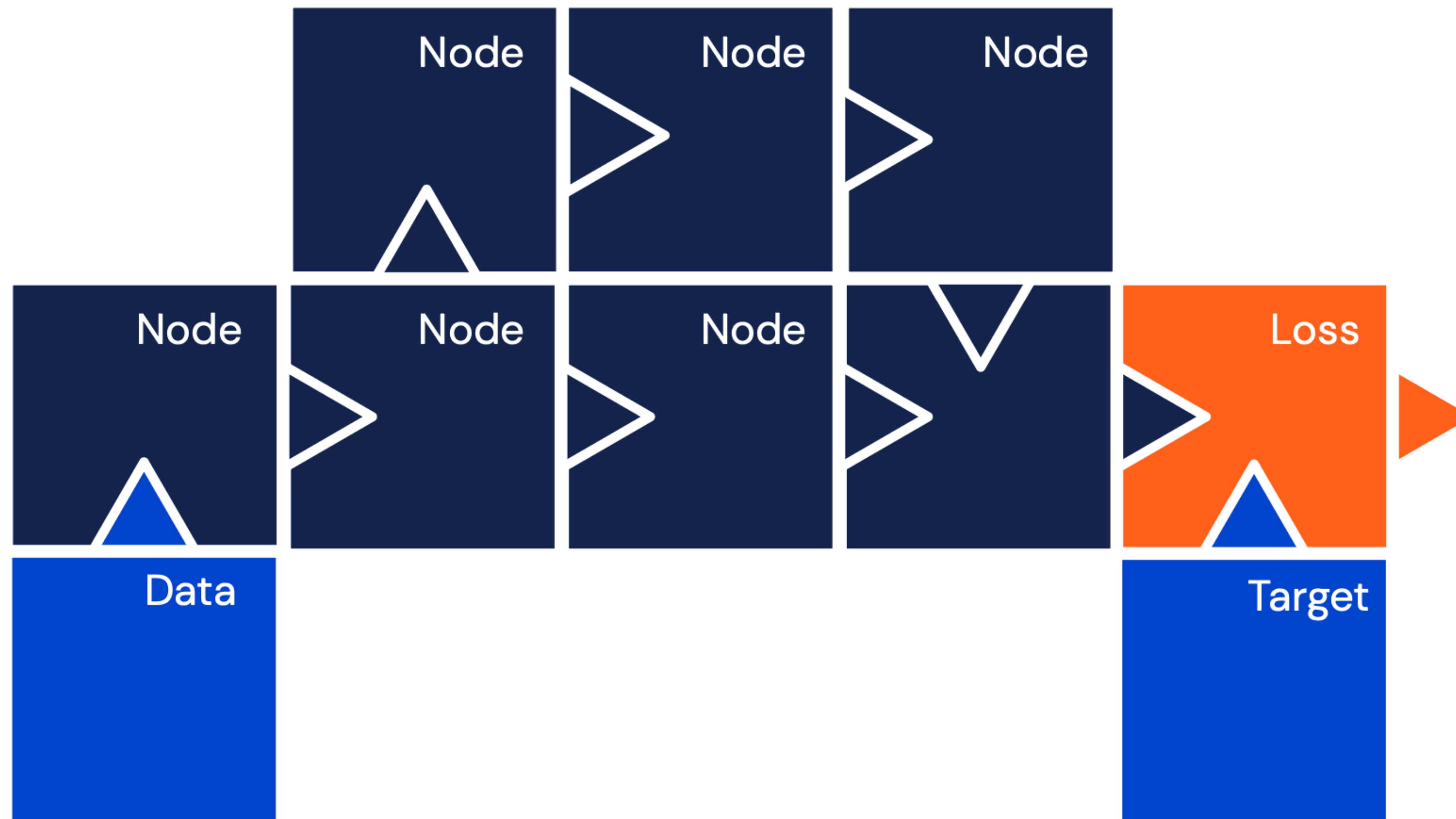
Deep Learning as LEGO for adults

How to build your best model



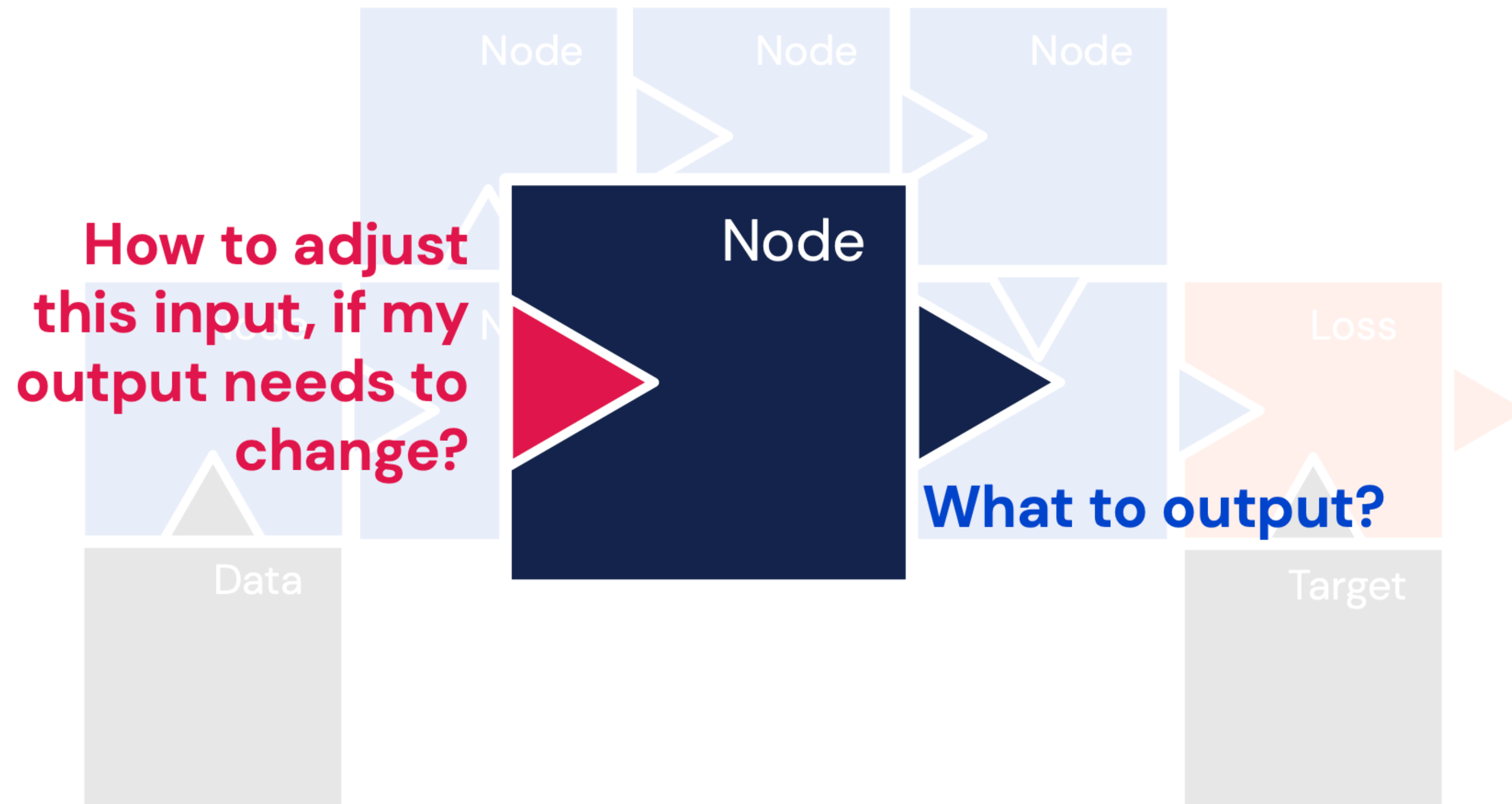
Deep Learning as LEGO for adults

How to build your best model



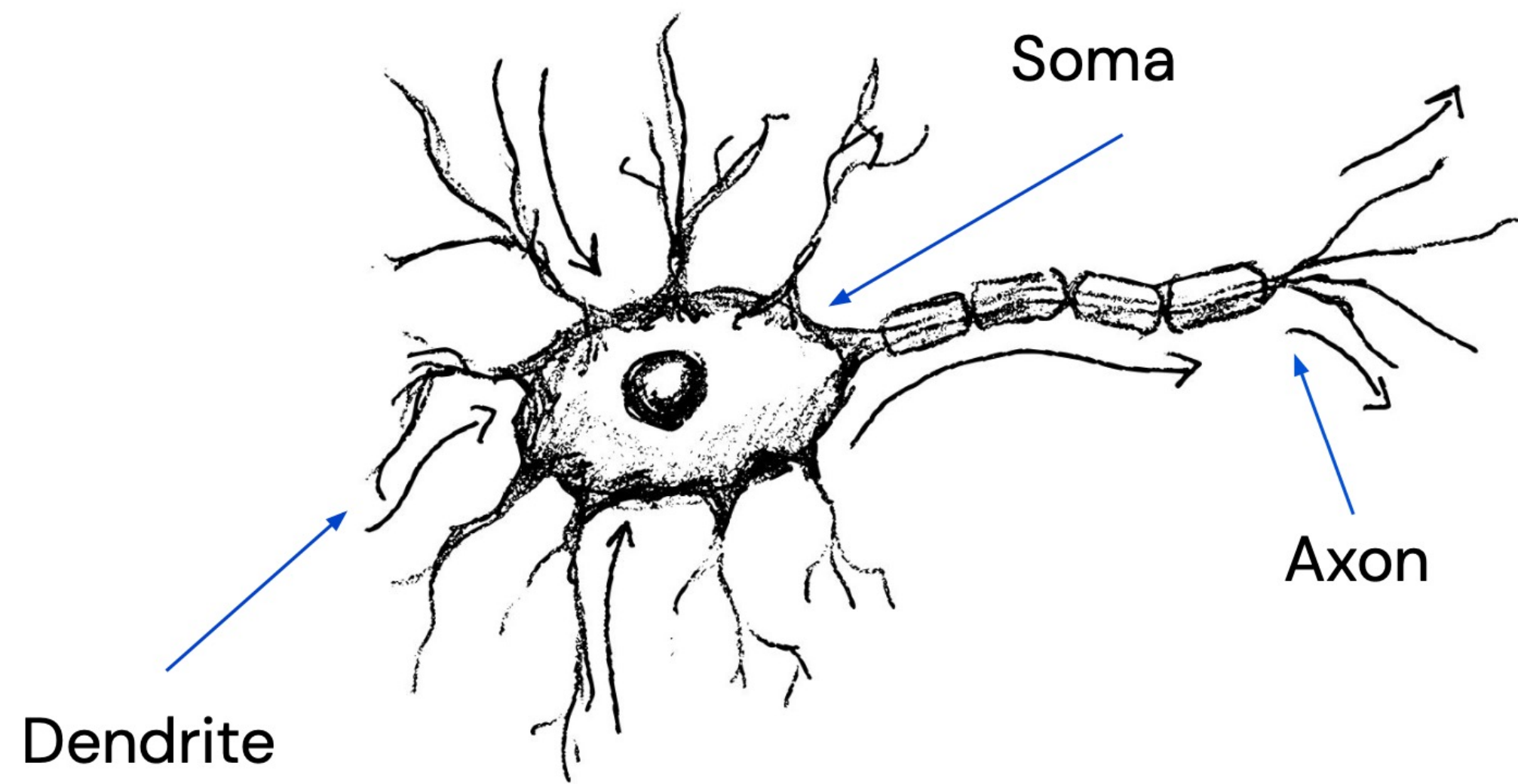
A

B

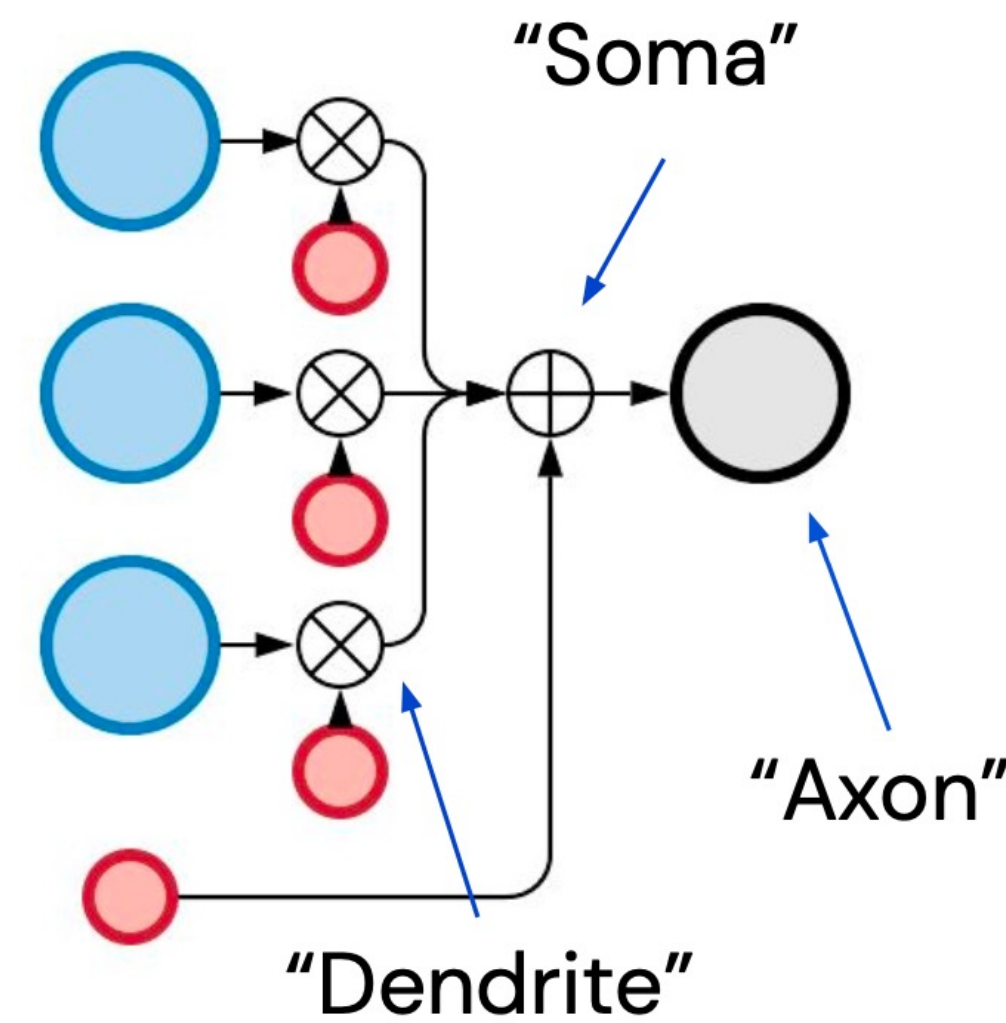


Artificial Neurons

Not biologically plausible, but still very useful



≠

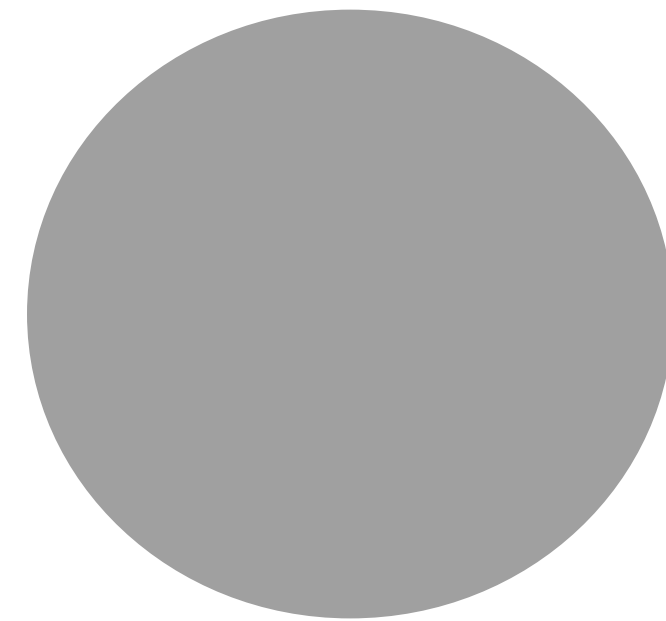
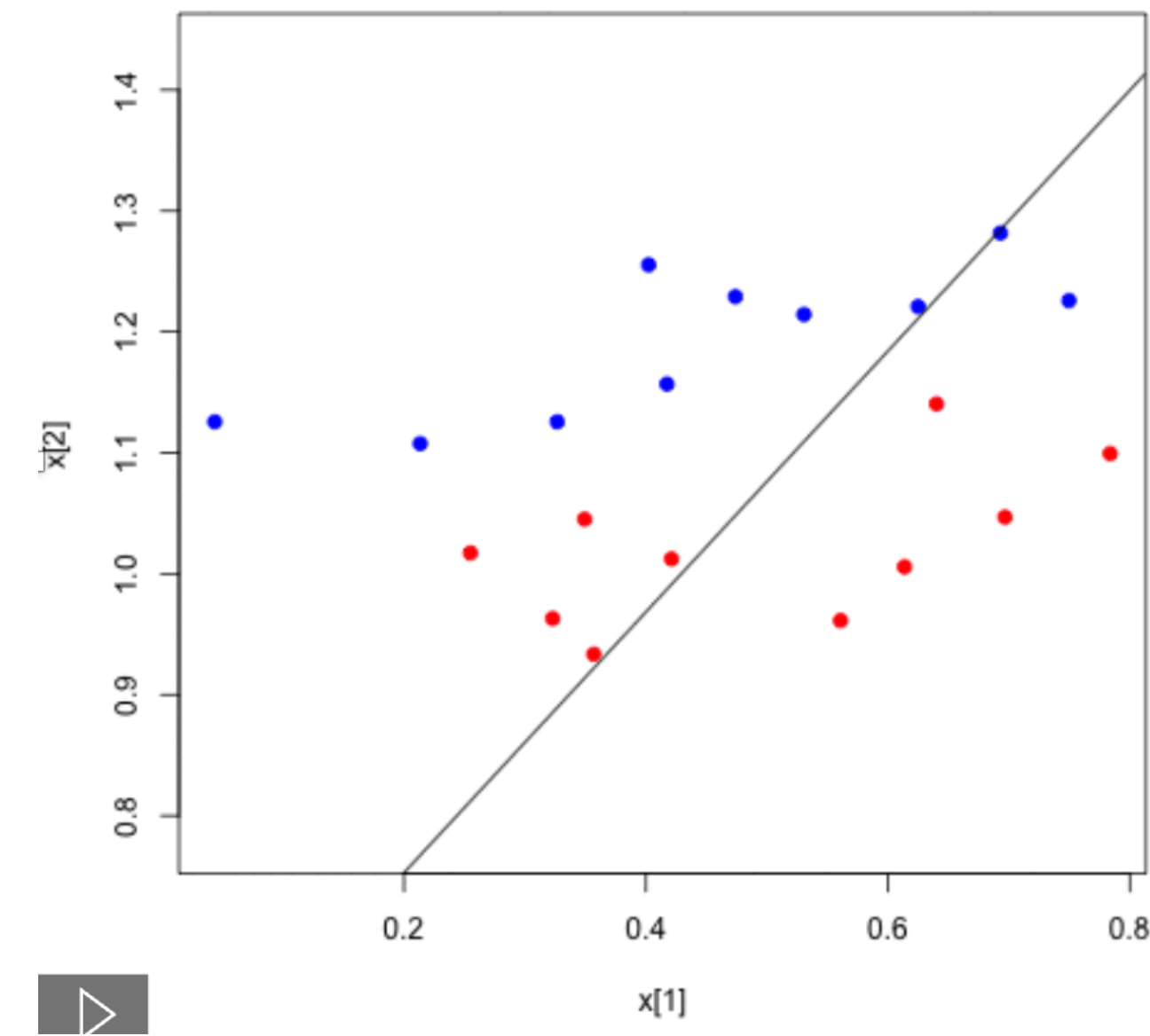


$$\sum_{i=1}^d \mathbf{w}_i \mathbf{x}_i + b$$
$$\sum_{i=0}^d \mathbf{w}_i \mathbf{x}_i \quad \mathbf{x}_0 := 1$$

The Perceptron (Rosenblatt, 1958)

Not biologically plausible, but still very useful

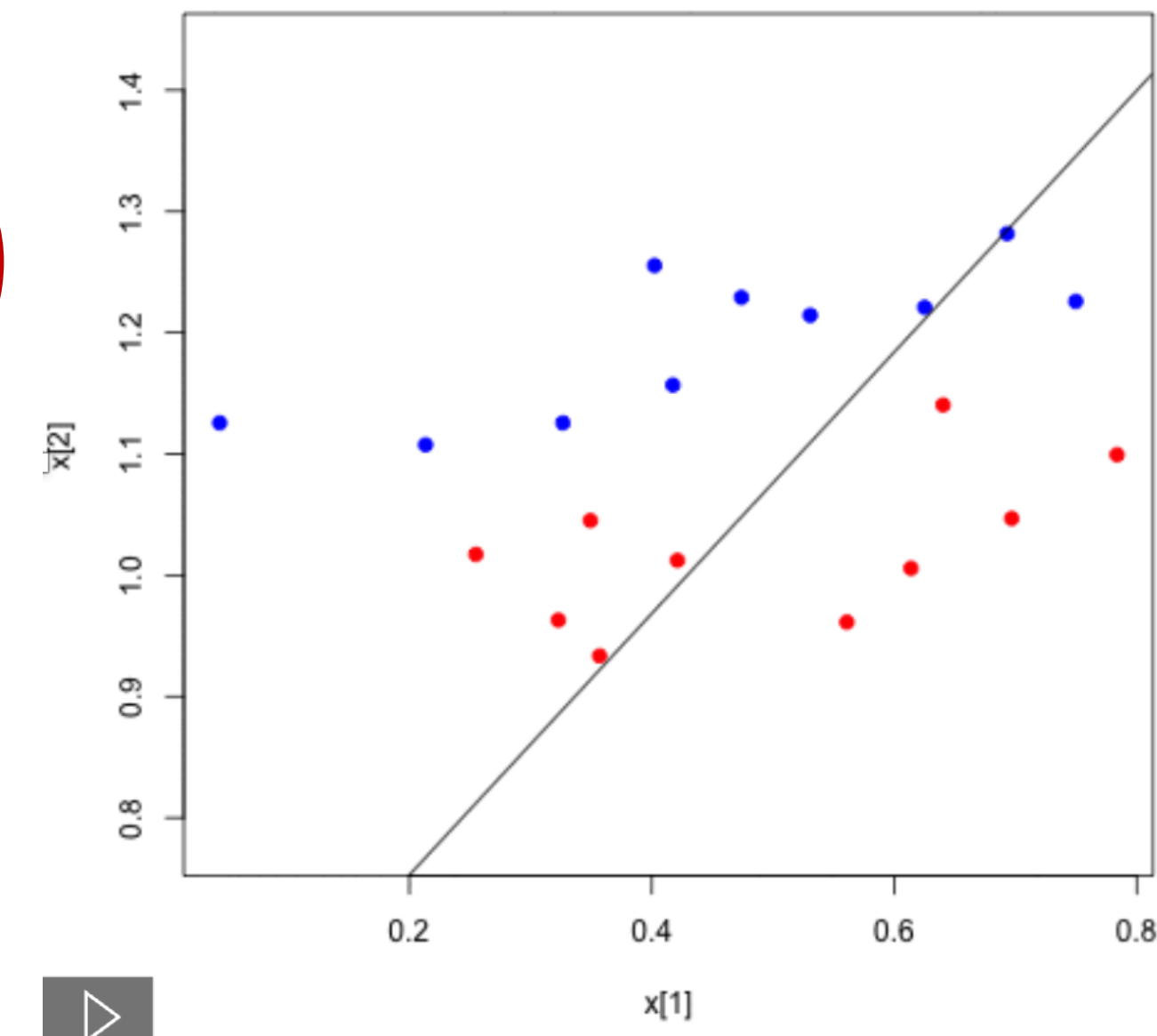
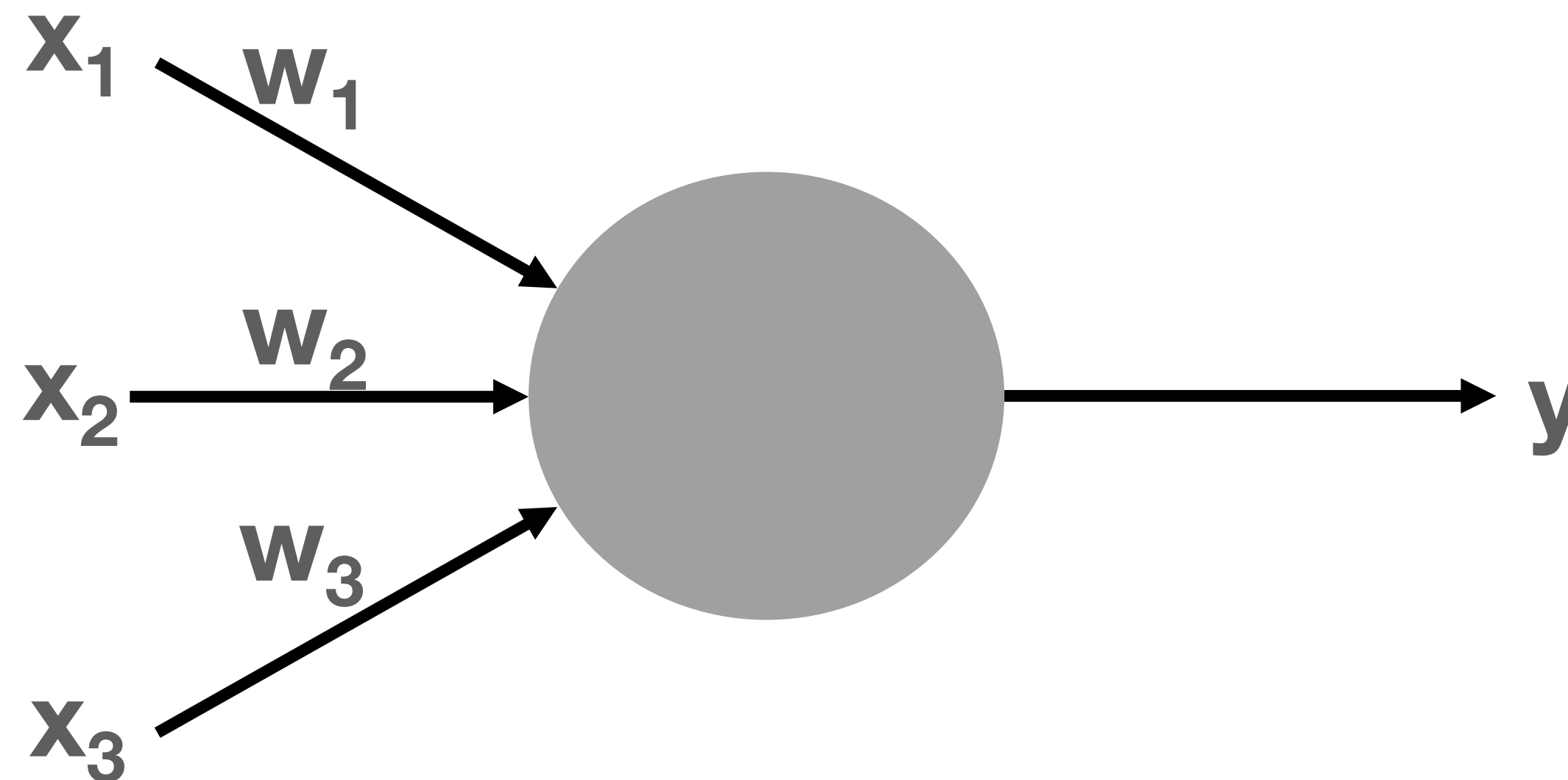
$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0, \\ 0 & \text{otherwise} \end{cases}$$



The Perceptron (Rosenblatt, 1958)

Not biologically plausible, but still very useful

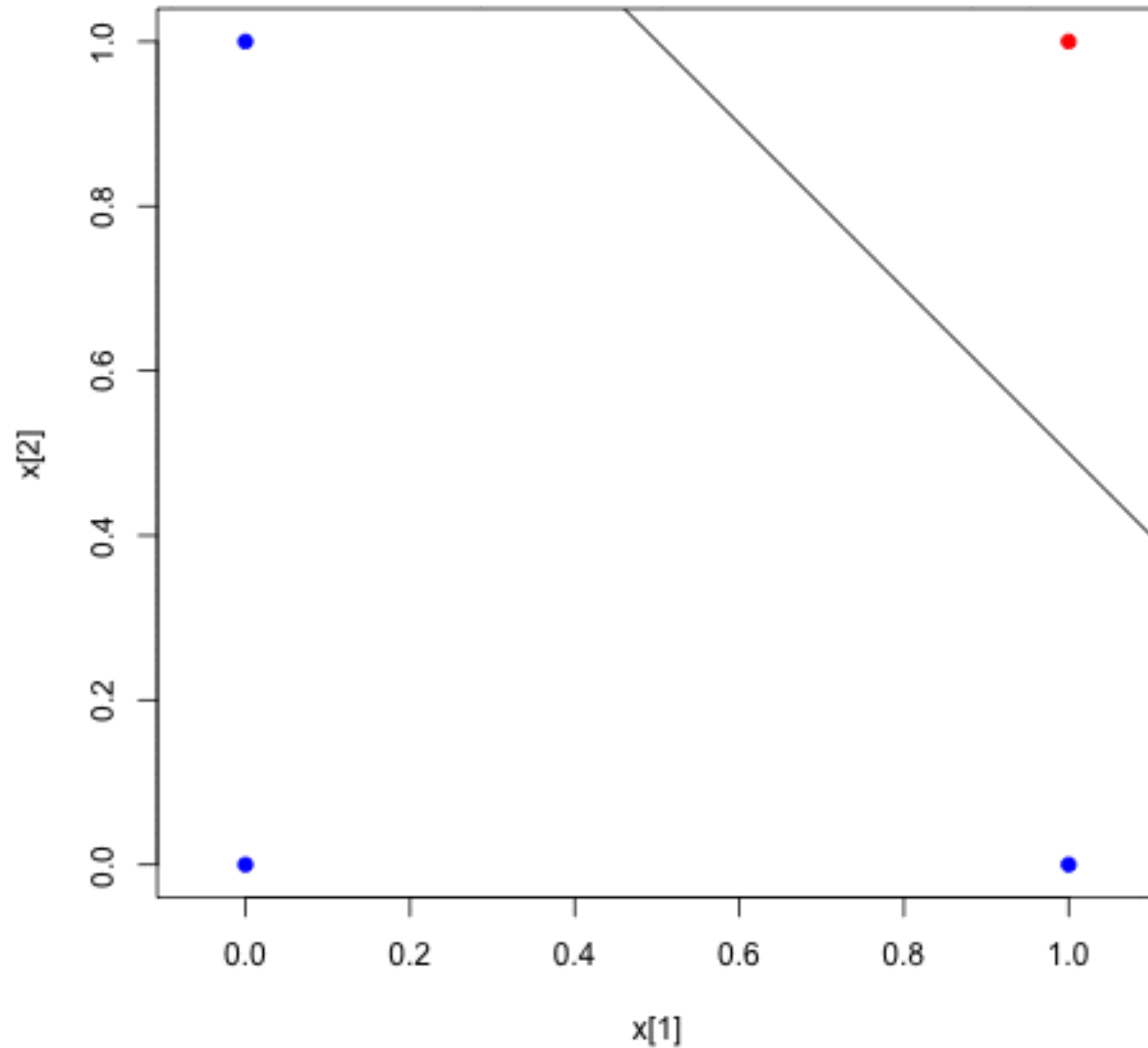
$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0, \\ 0 & \text{otherwise} \end{cases}$$



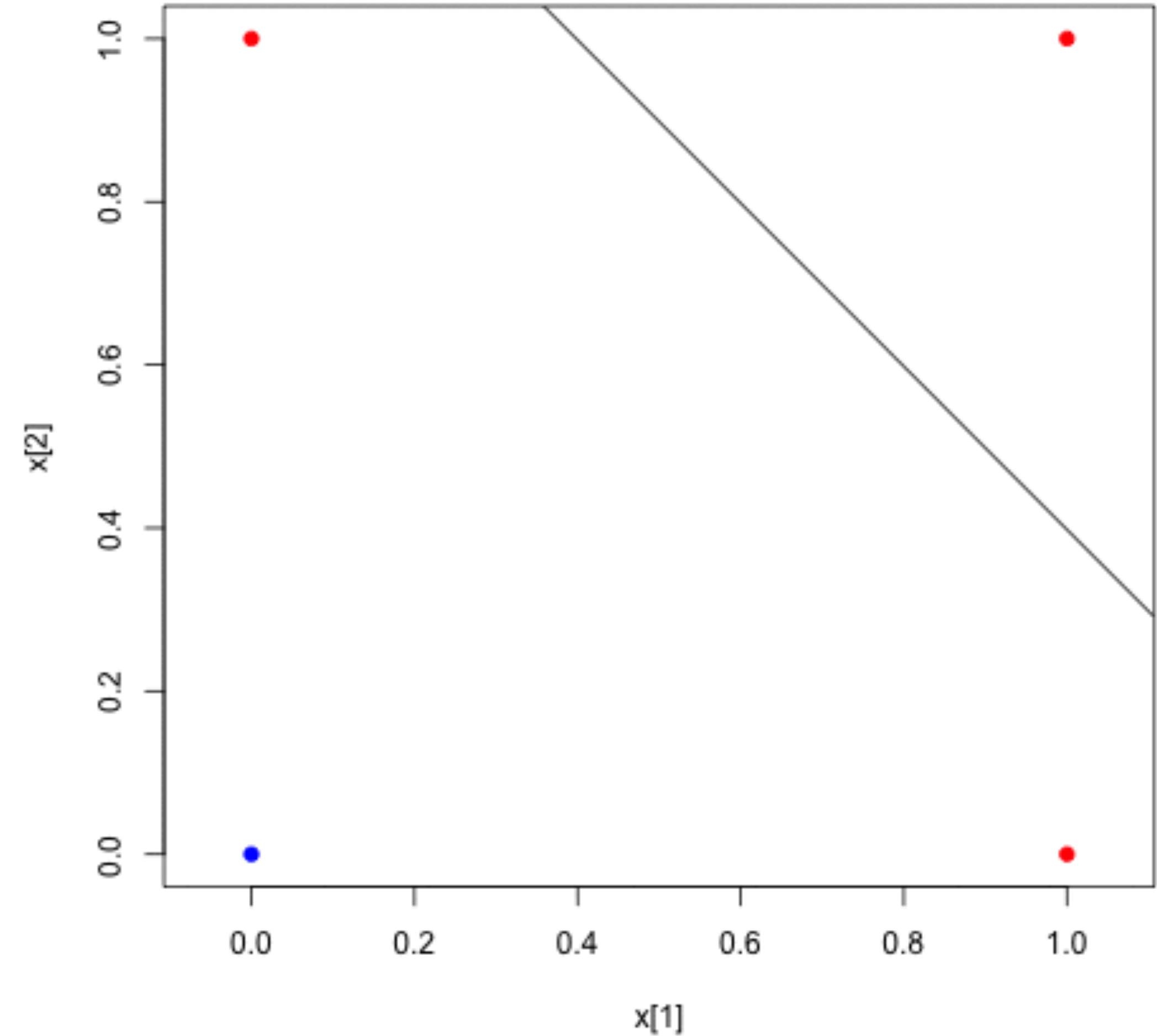
The perceptron can solve AND and OR

Linear problems can be solved

epoch 1 error 0.000

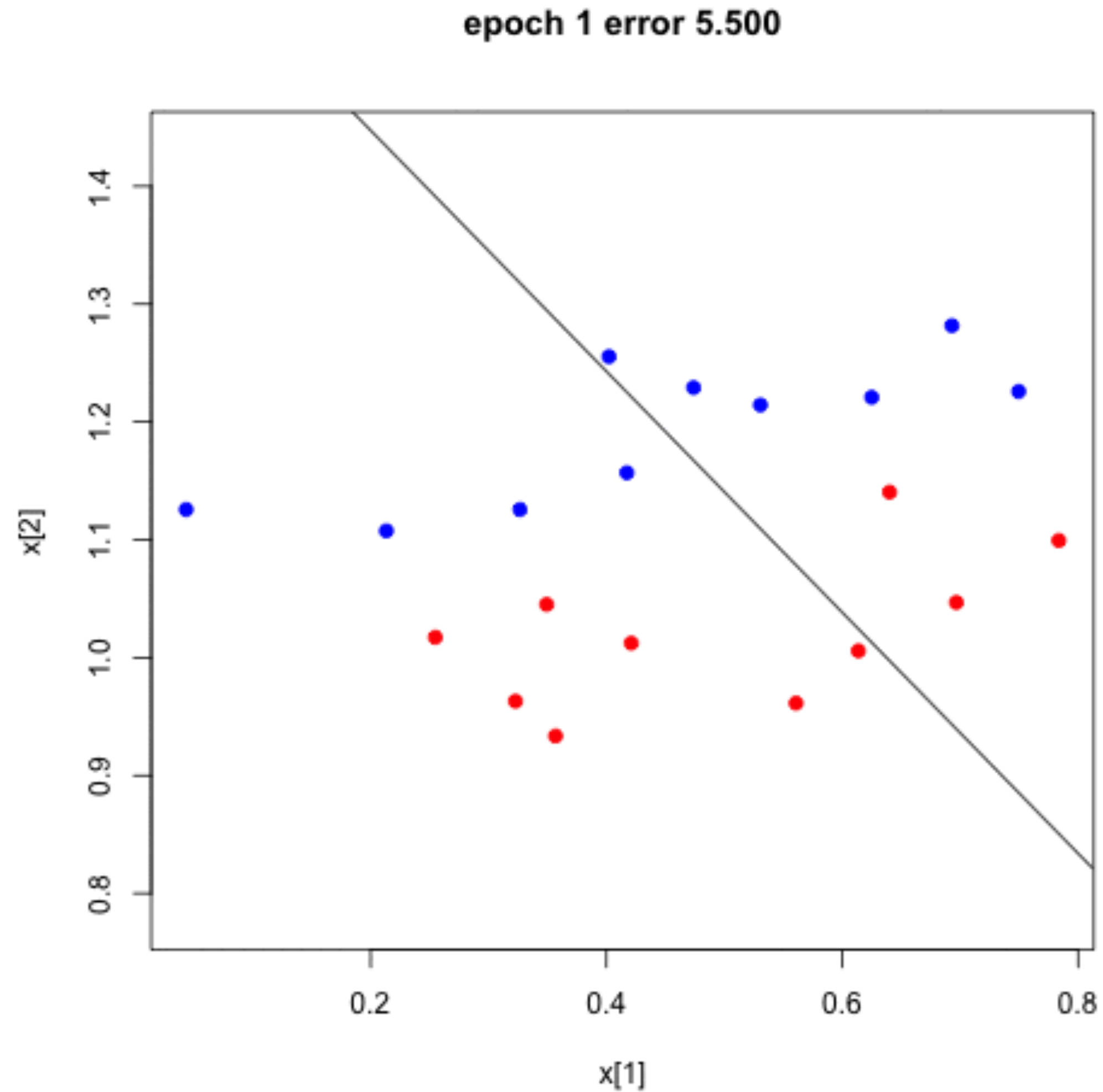


epoch 1 error 1.000



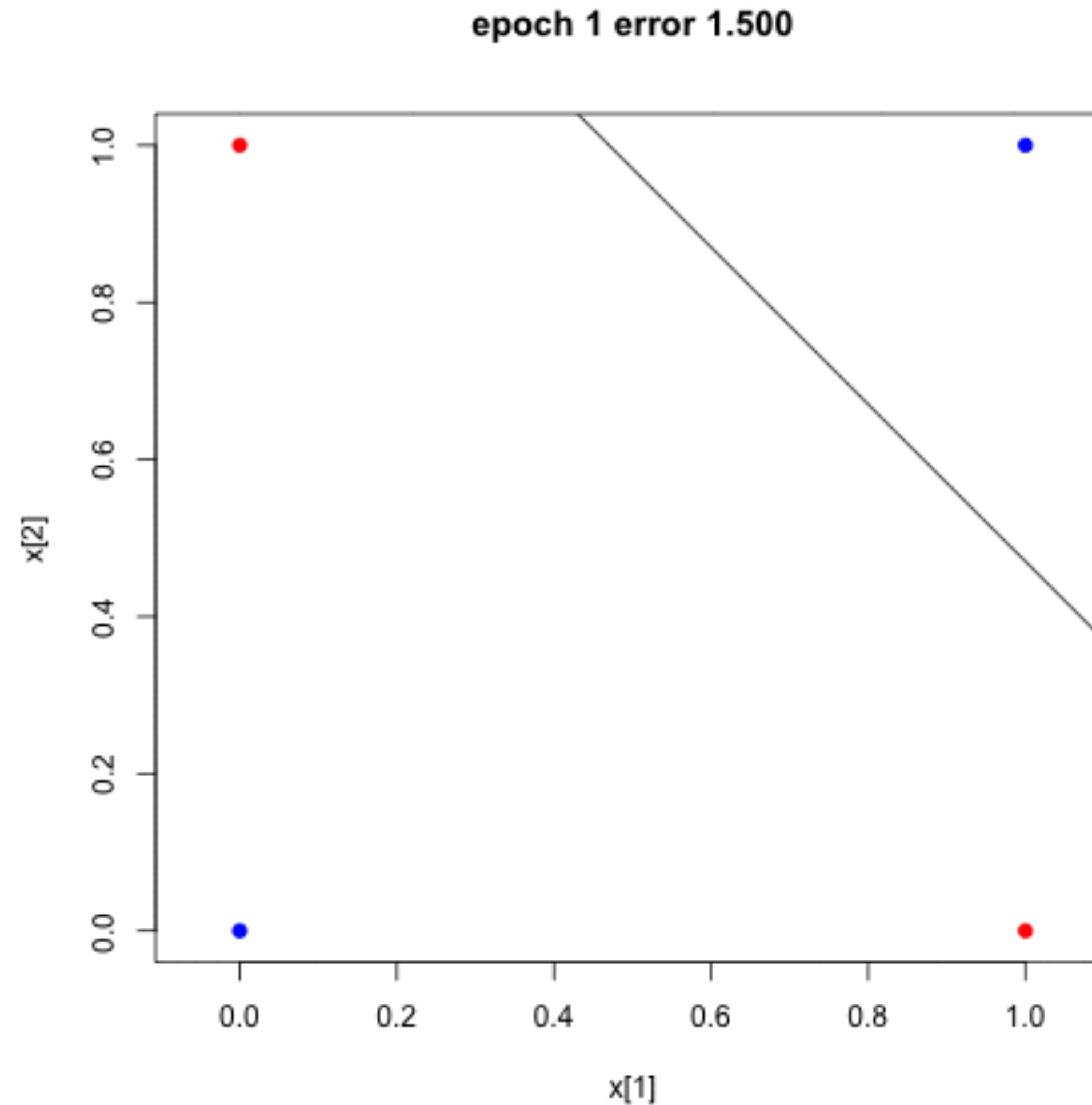
Perceptron on linearly separable data

Linear problems can be solved



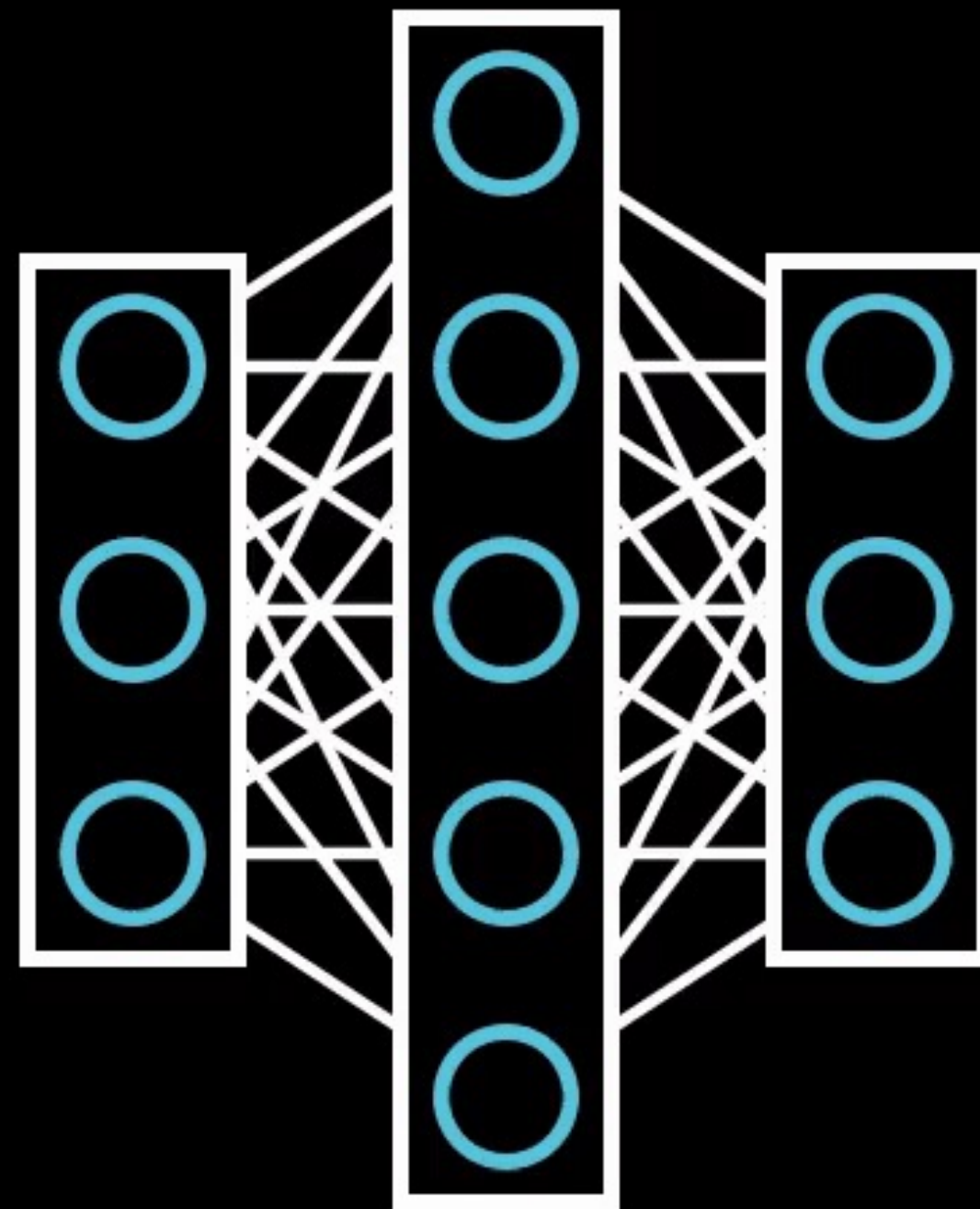
The perceptron cannot solve XOR

Minsky & Papert, 1969: Dawn of the first AI winter



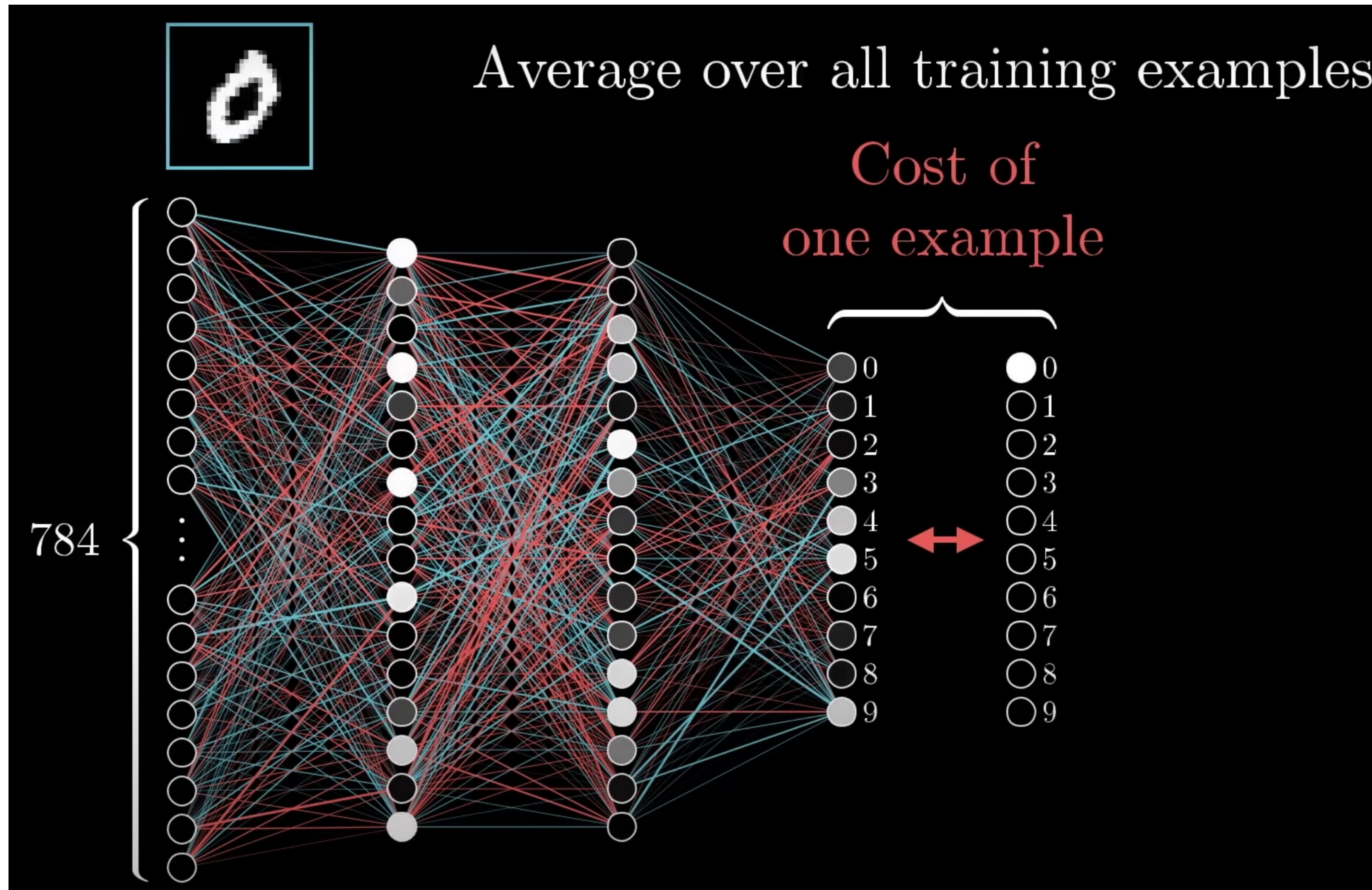
The secret sauce (1/2): Multilayer Perceptrons

Backpropagation allowed the use of deeper networks



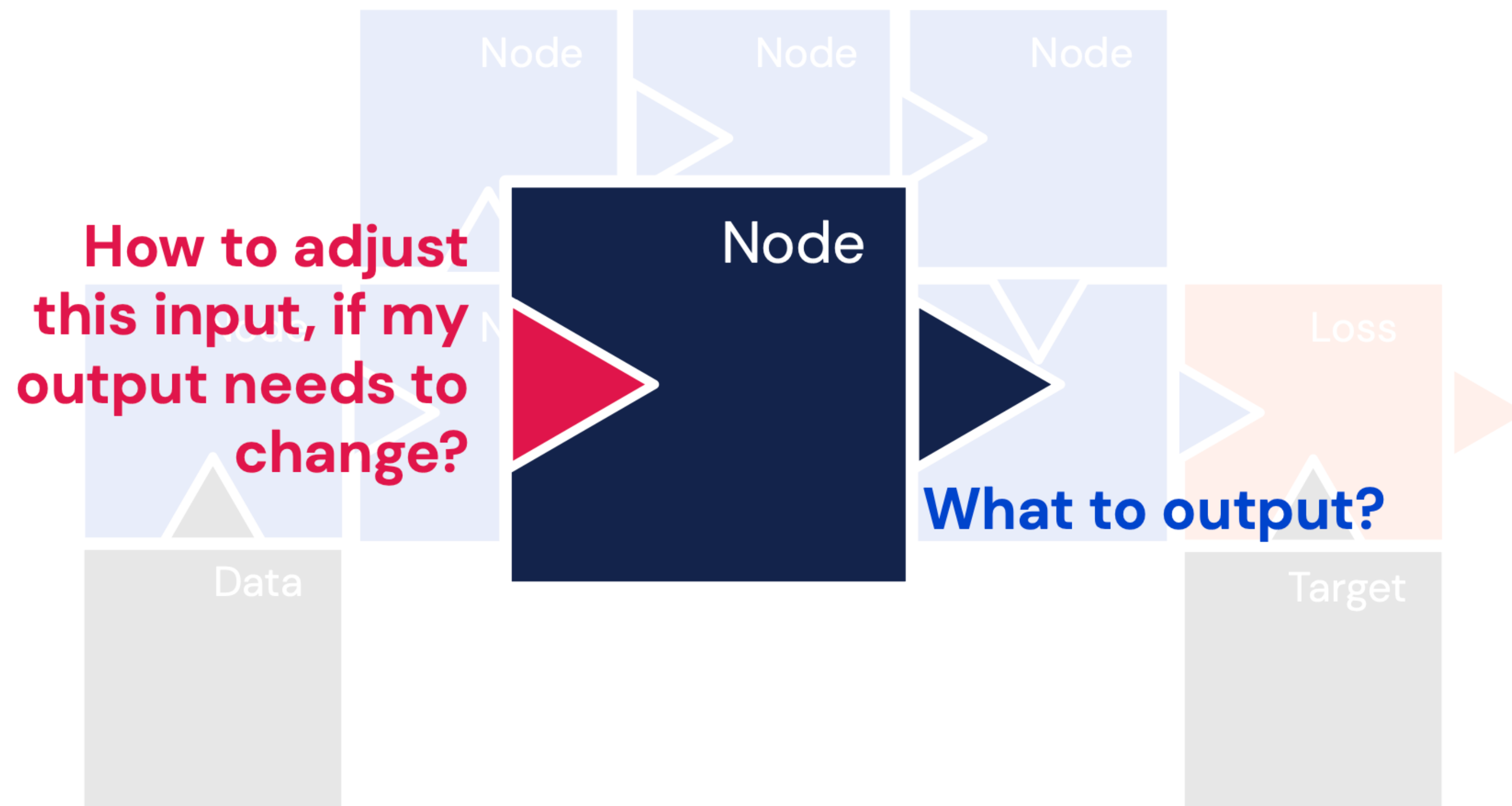
The secret sauce (1/2): Multilayer Perceptrons

Backpropagation allowed the use of deeper networks



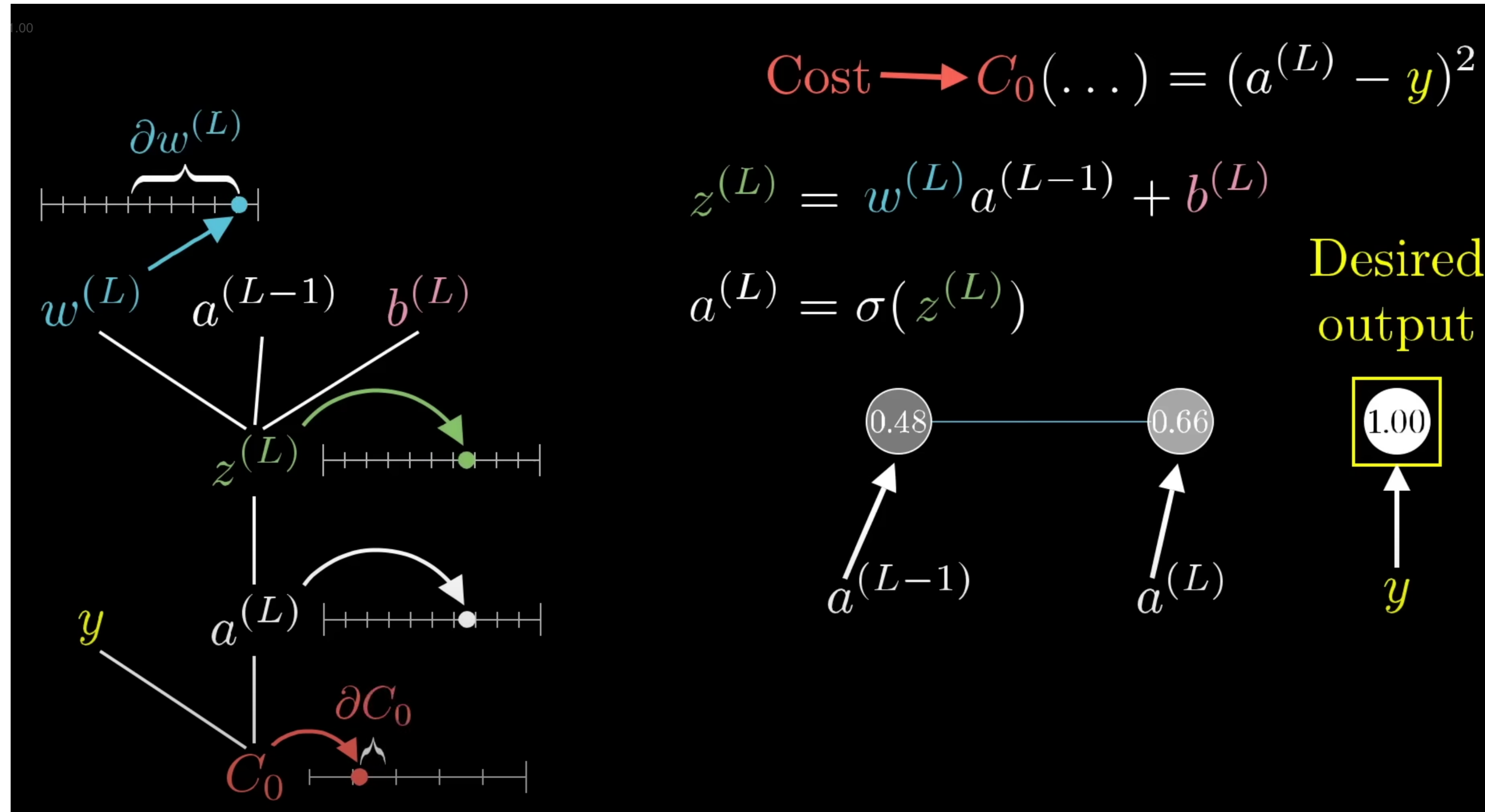
How do we tune hidden neurons

Backpropagation allowed the use of deeper networks



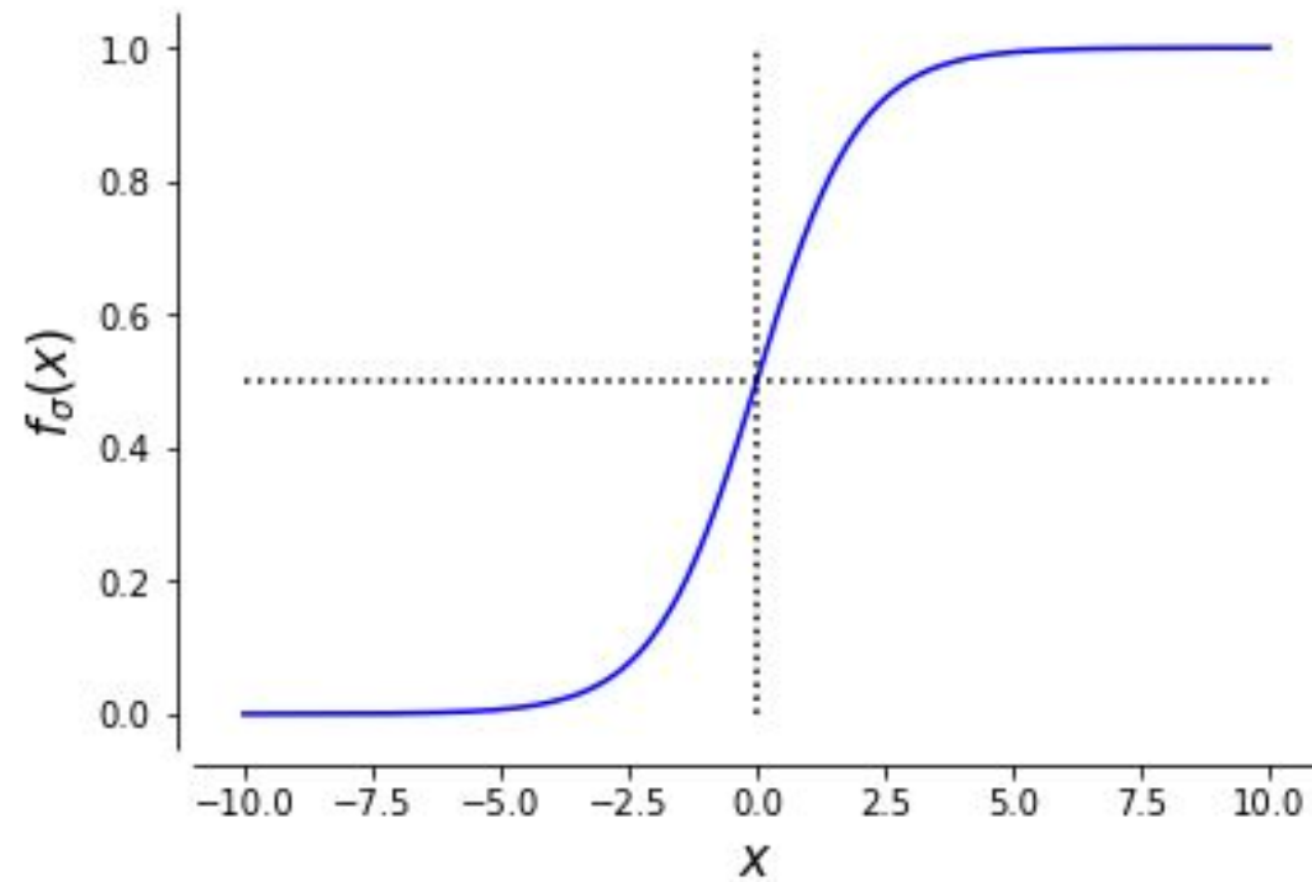
The secret sauce (1/2): Multilayer Perceptrons

Backpropagation allowed the use of deeper networks



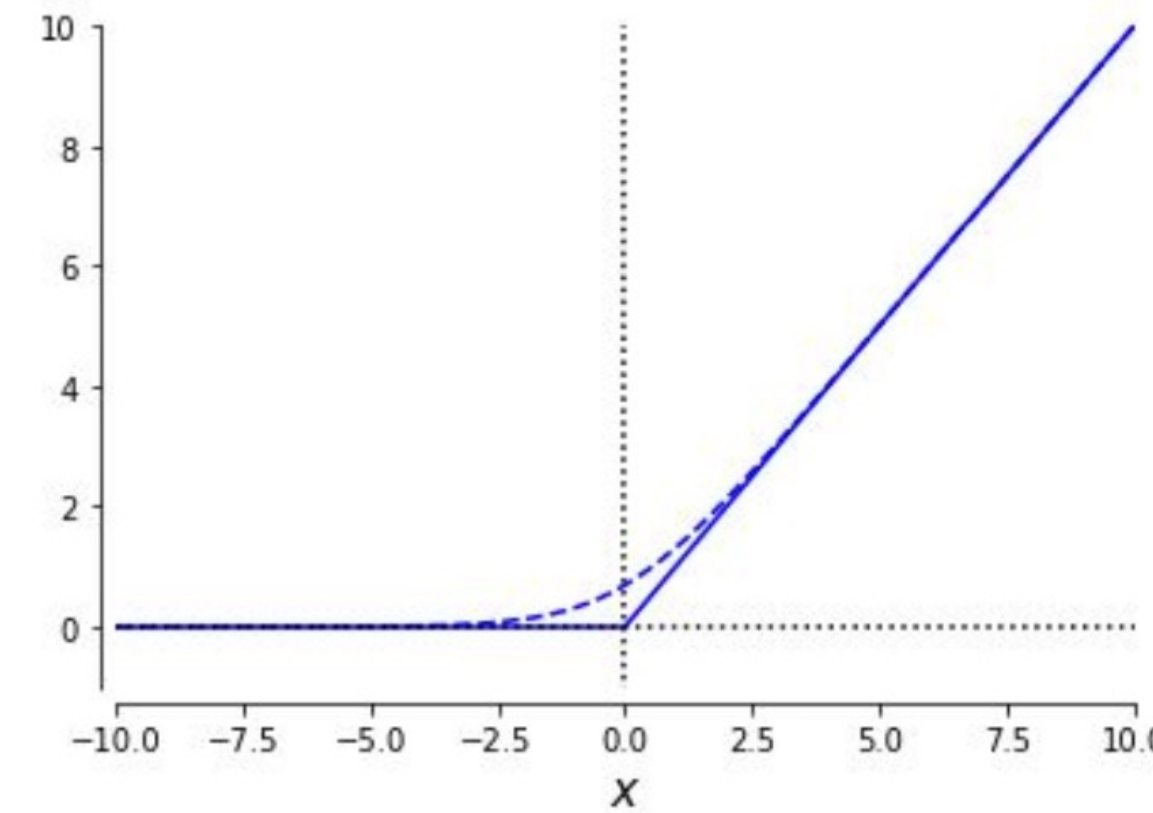
The secret sauce (2/2): Activation functions

Non-linearities allow us to solve non-linear problems



$$f_{\sigma}(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{x}}}$$

$$f_{\sigma}(\mathbf{x}) = \frac{e^{\mathbf{x}}}{e^{\mathbf{x}} + 1}$$



$$f_{\text{relu}}(\mathbf{x}) = \max(0, \mathbf{x})$$

$$f_{\text{sp}}(\mathbf{x}) = \log(1 + e^{\mathbf{x}})$$

Activation functions are often called **non-linearities**.
Activation functions are applied **point-wise**.

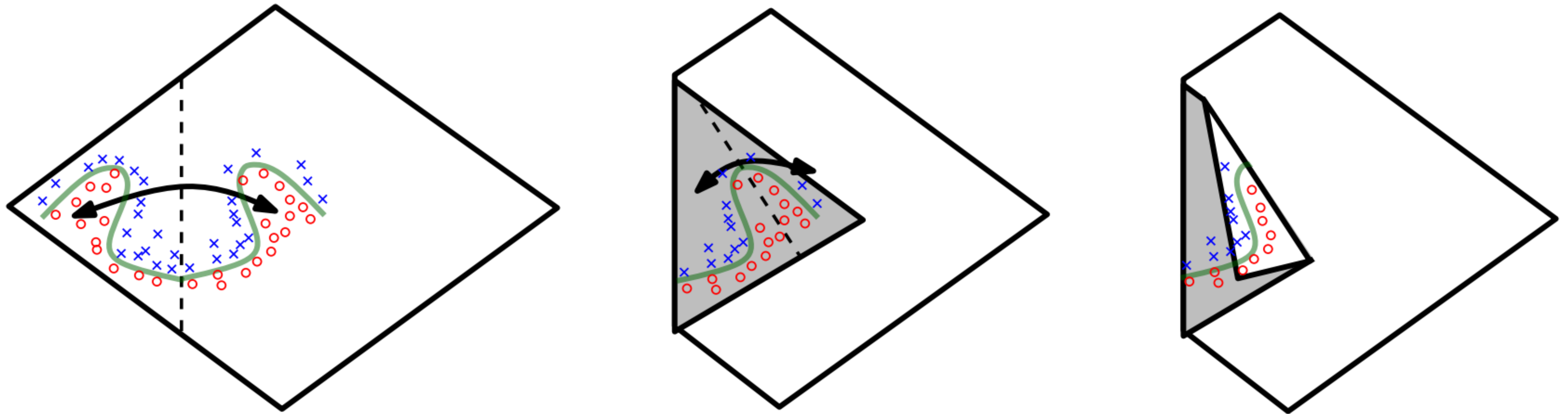
One of the **most commonly used activation functions**.
Made **math analysis of networks much simpler**.

Solve XOR with just 2 hidden neurons

- Hidden layers bend and deform input space**
- Last linear model does linear classification**
- Non-linear transformations are key for deep learning!**
- Our network became a feature extractor!**

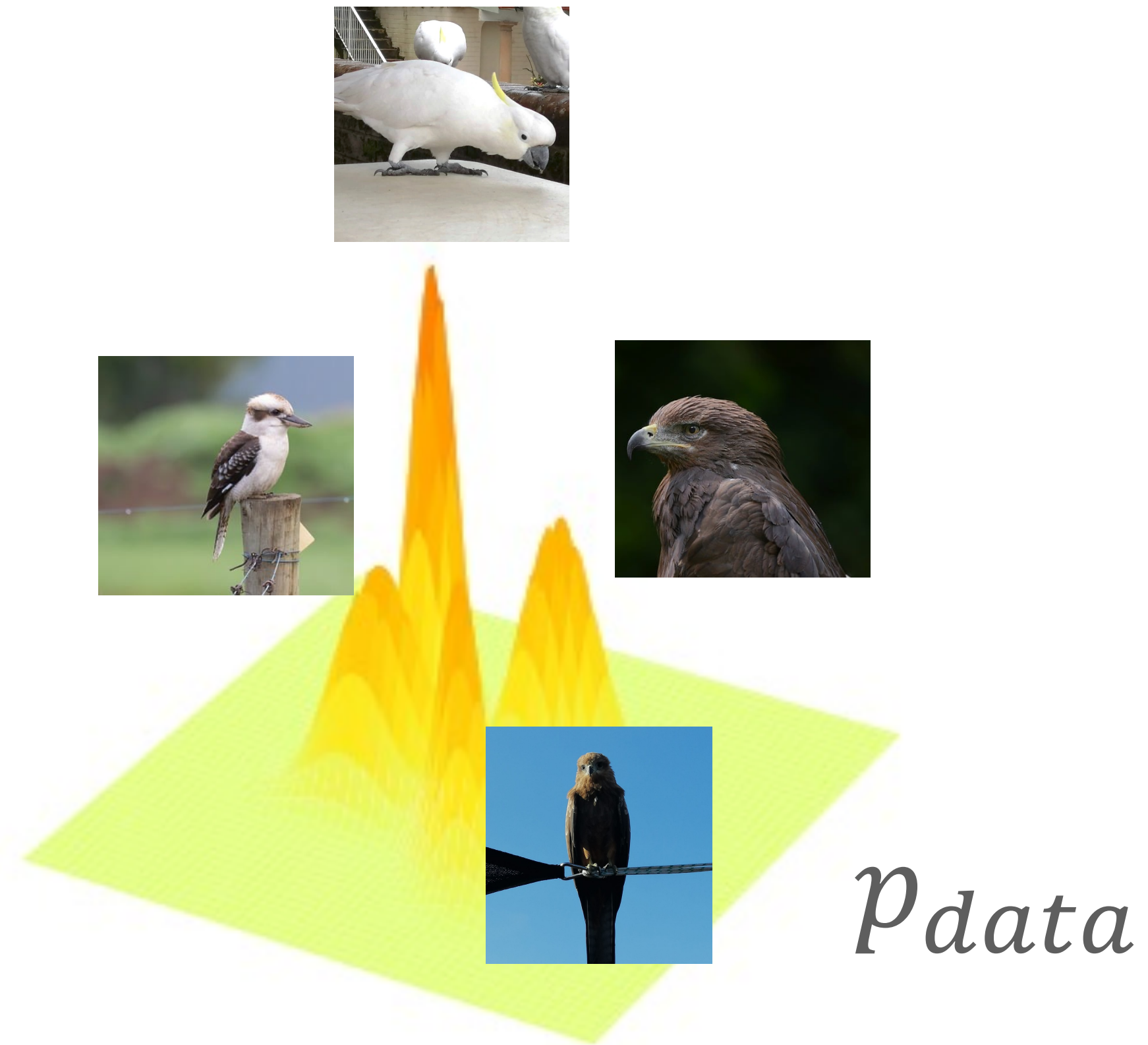
What is deep learning doing?

Transforming data non-linearly into a better representation



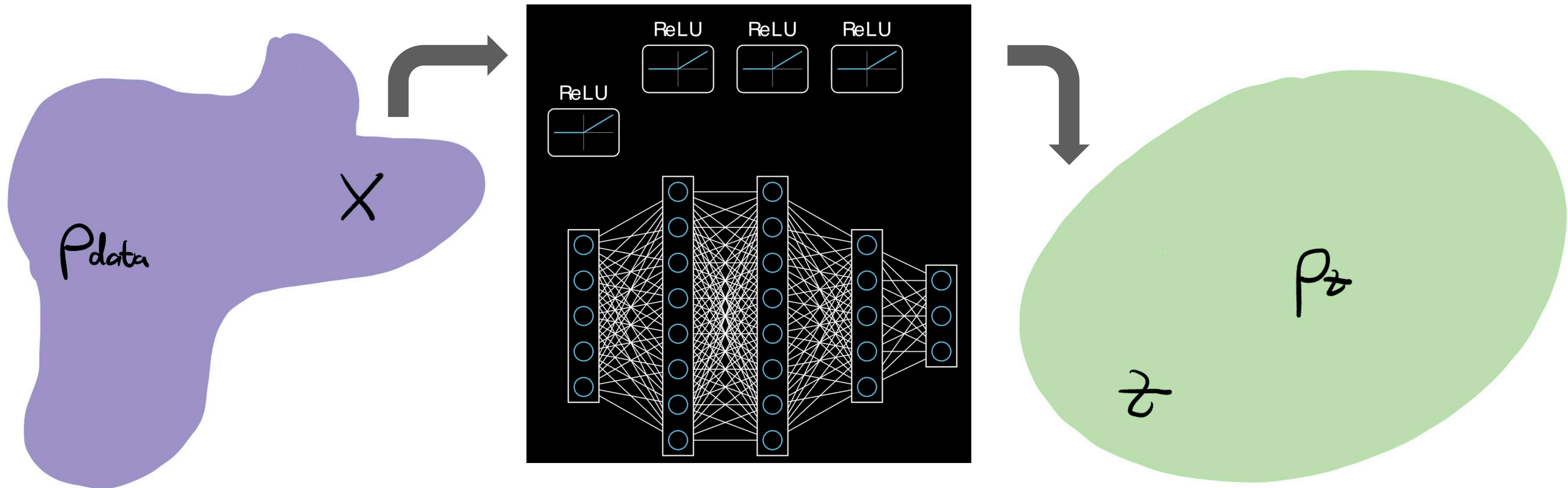
What is deep learning doing?

Transforming data non-linearly into a better representation



What is deep learning doing?

Transforming data non-linearly into a better representation



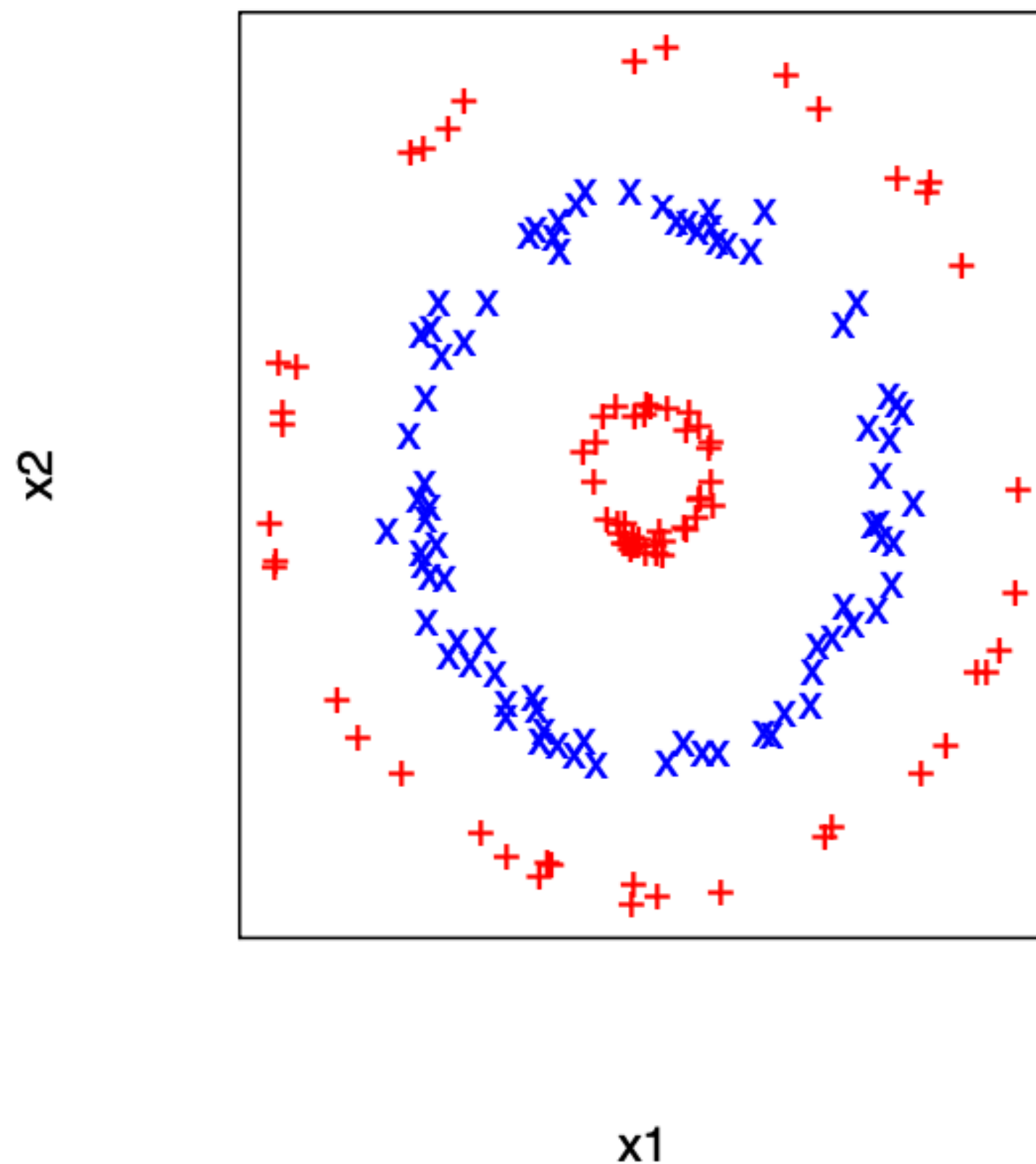
Let's try it ourselves

See the power of hidden layers

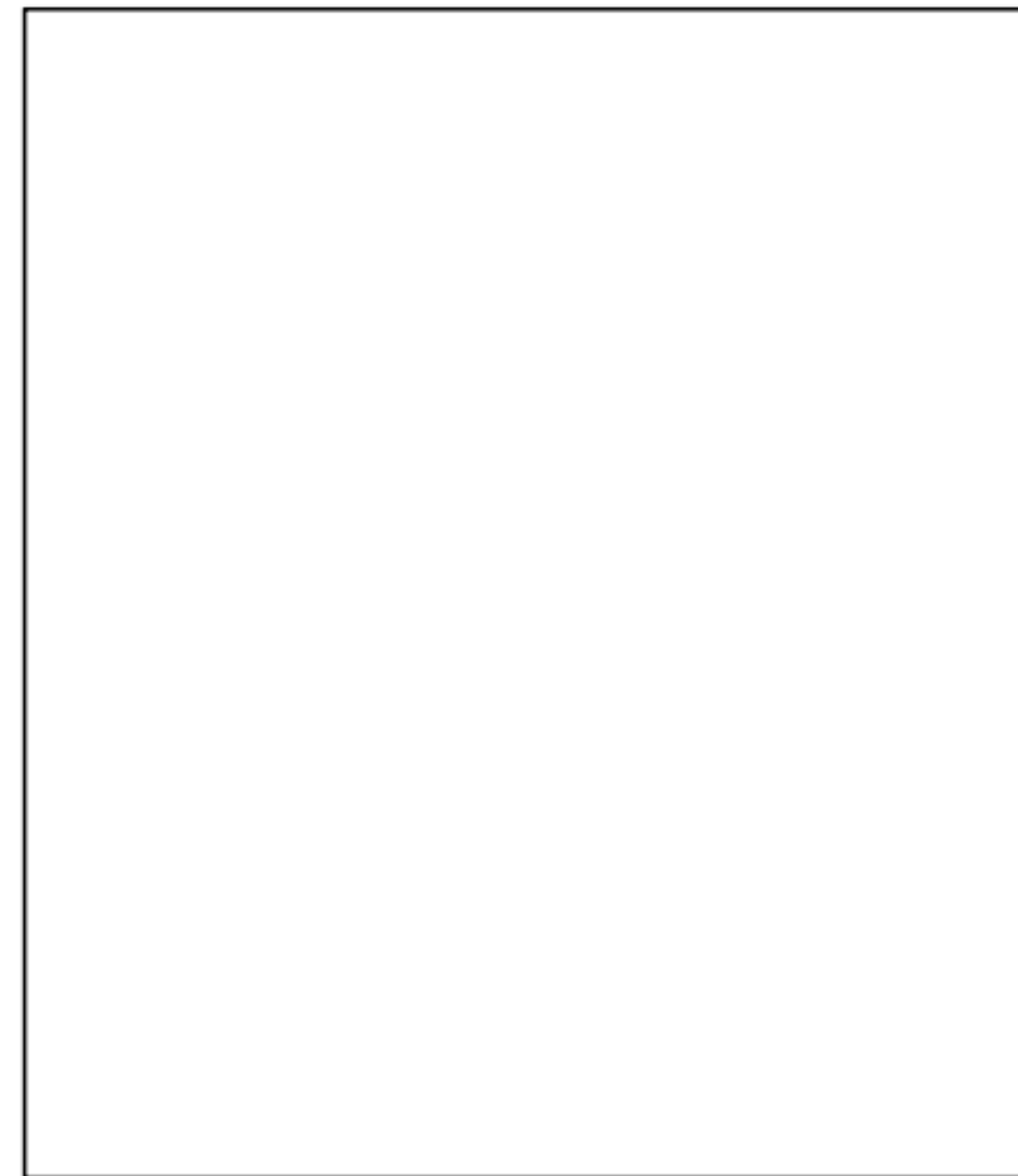
<http://playground.tensorflow.org/>

Feature/Representation engineering is hard

The Deep Learning way: Let the network figure it out for you



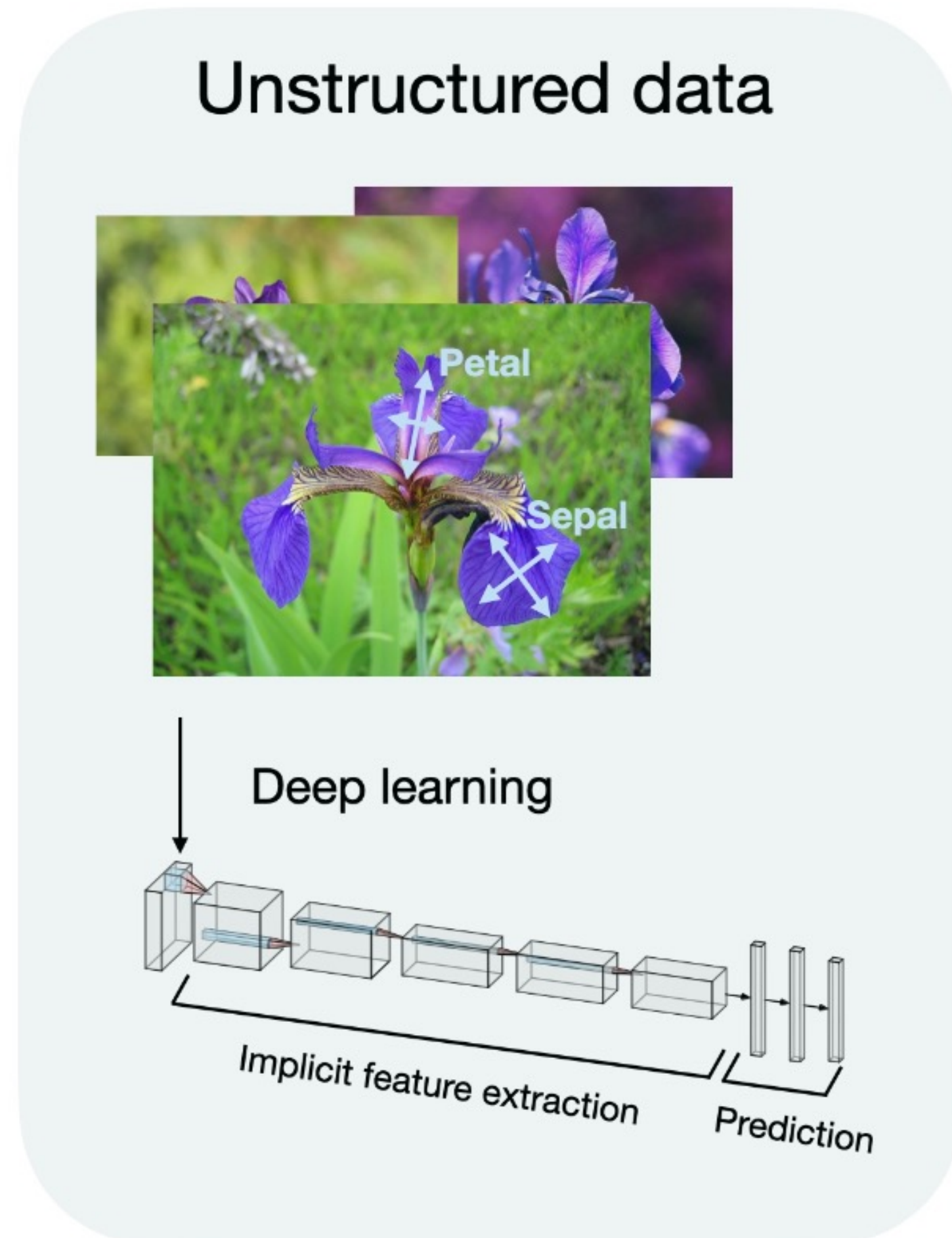
?



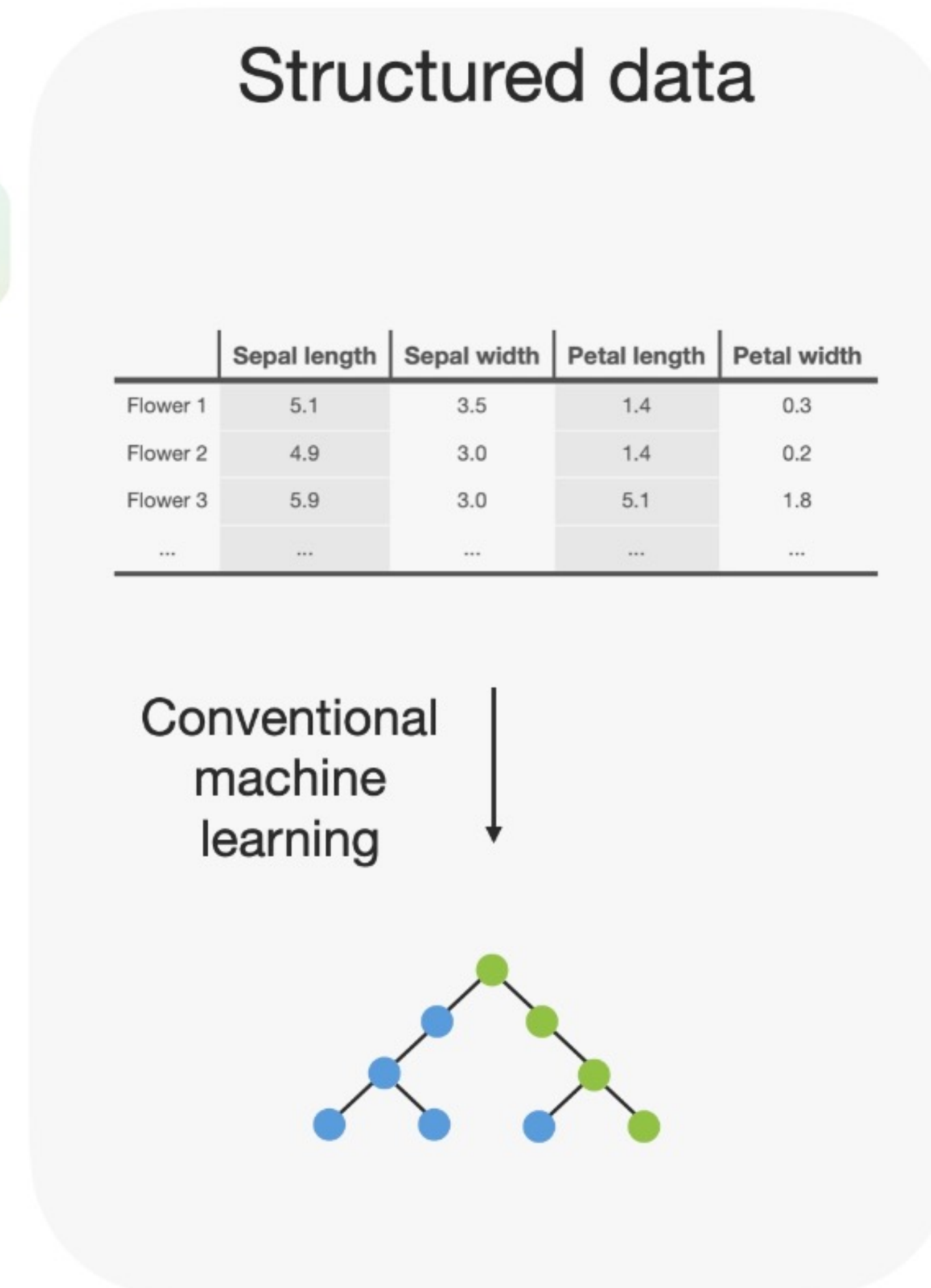
?

Deep Learning performs feature extraction

The deeper the network, the more complex the patterns can be

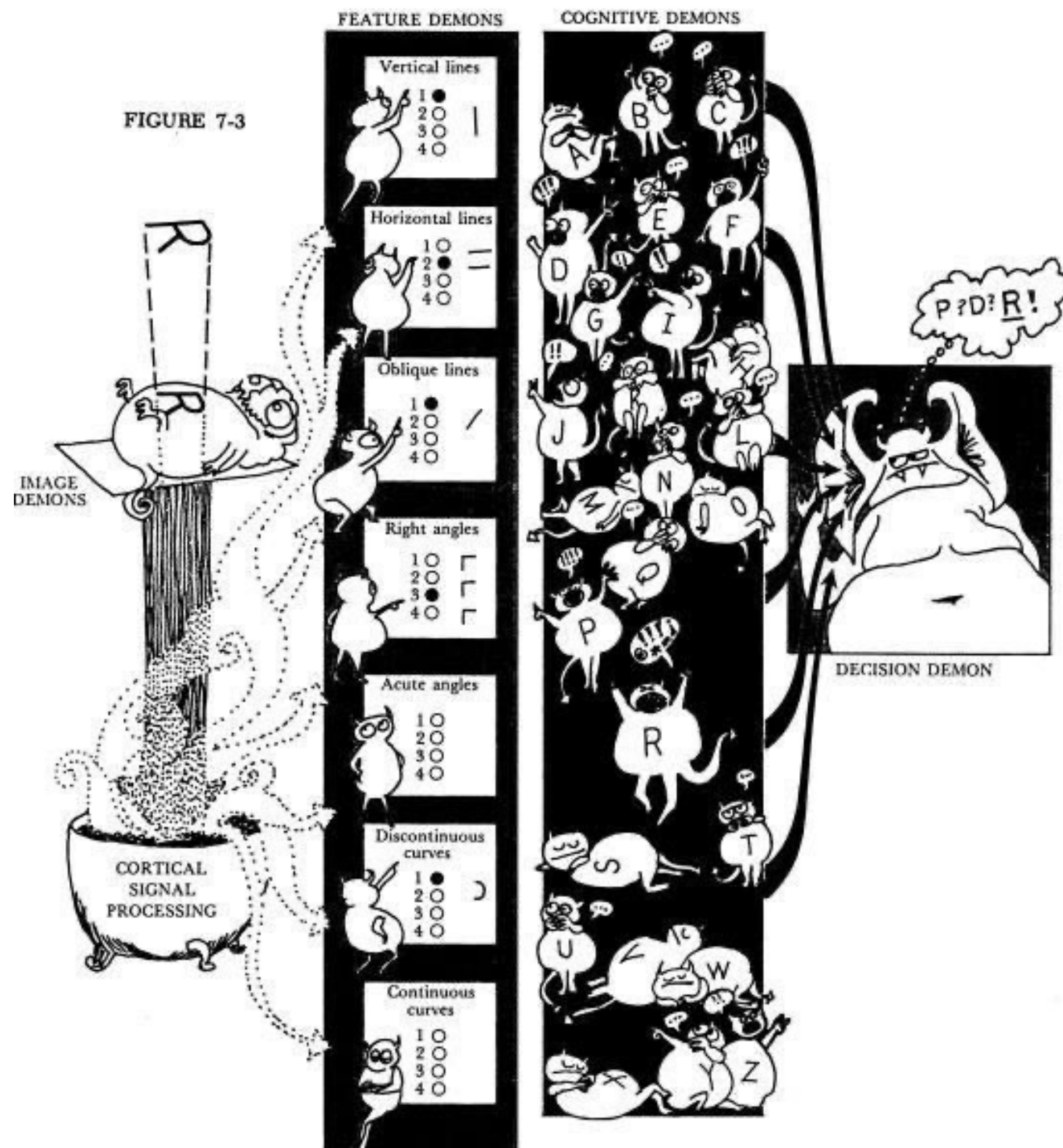


“Manual”
feature extraction



The demons in the pandemonium

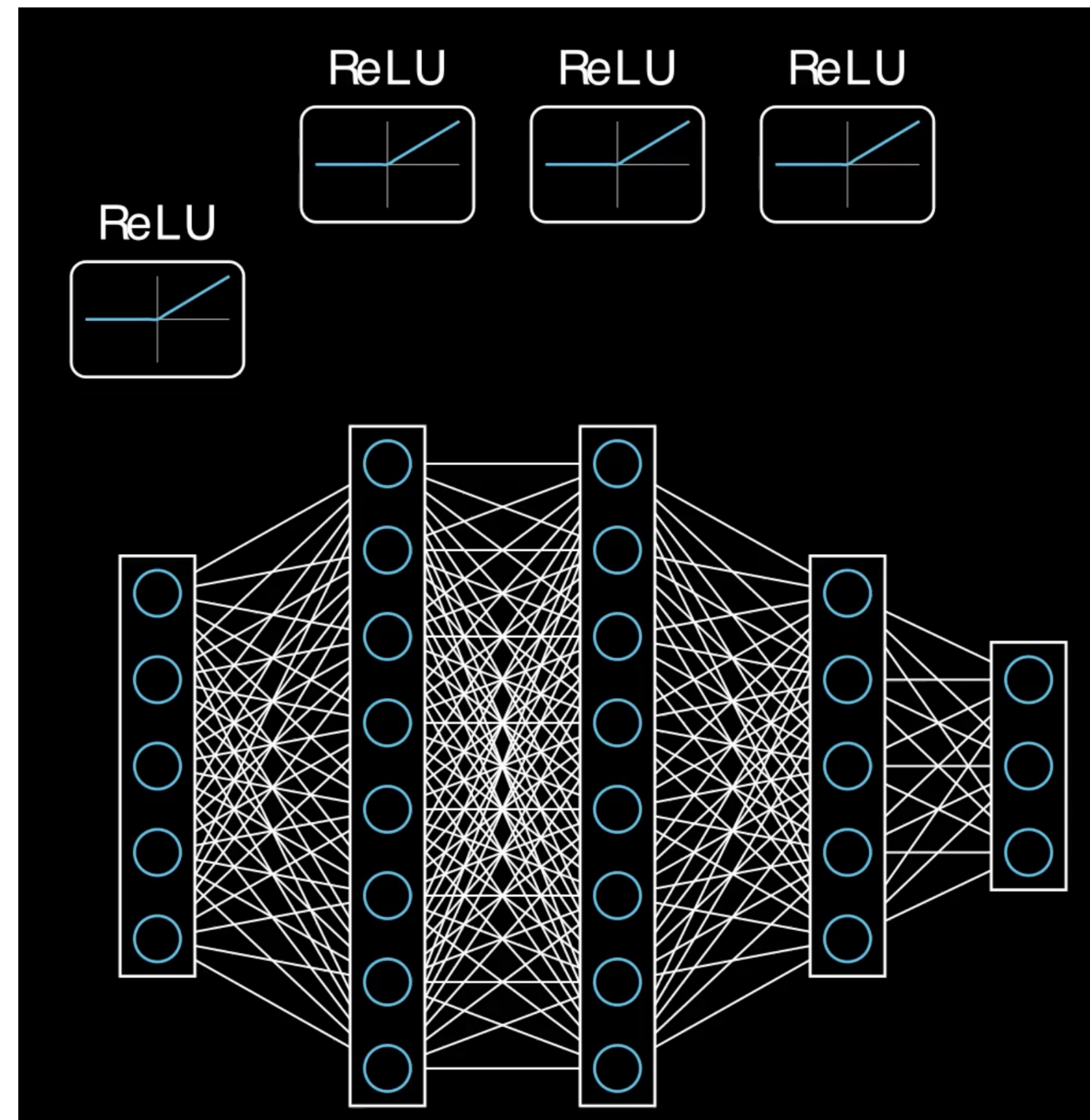
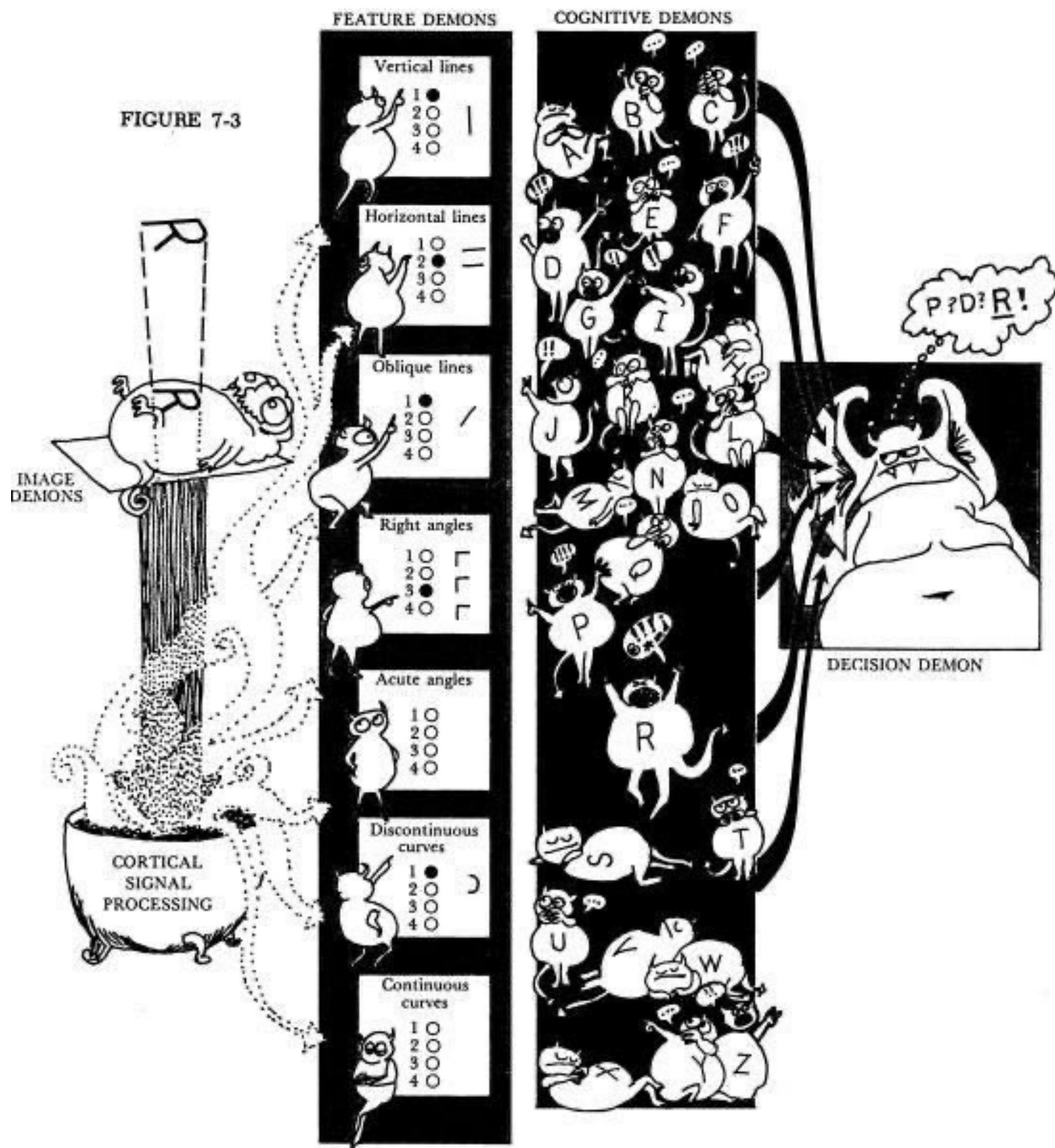
Selfridge, 1959: Pattern recognition triggers downstream cognition



The demons in the pandemonium

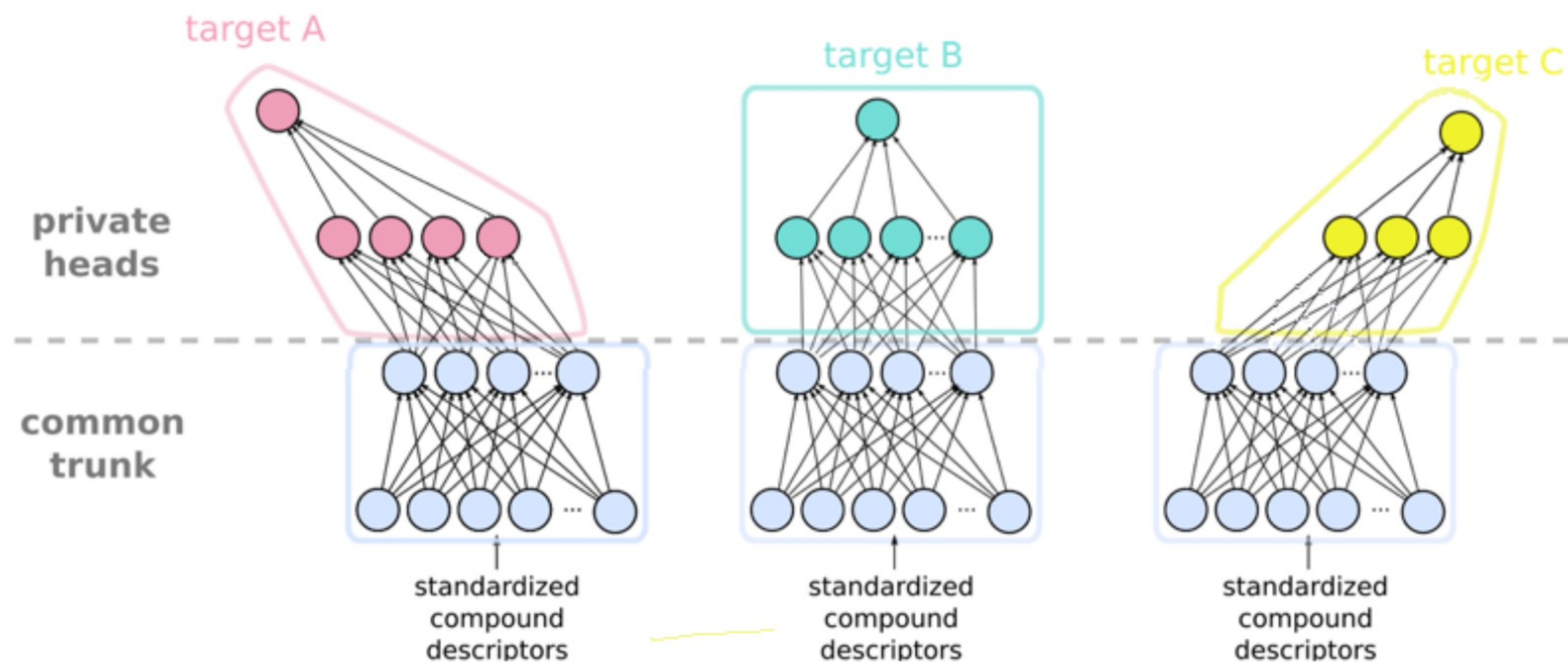
Trunk versus head

266 7. Pattern recognition and attention



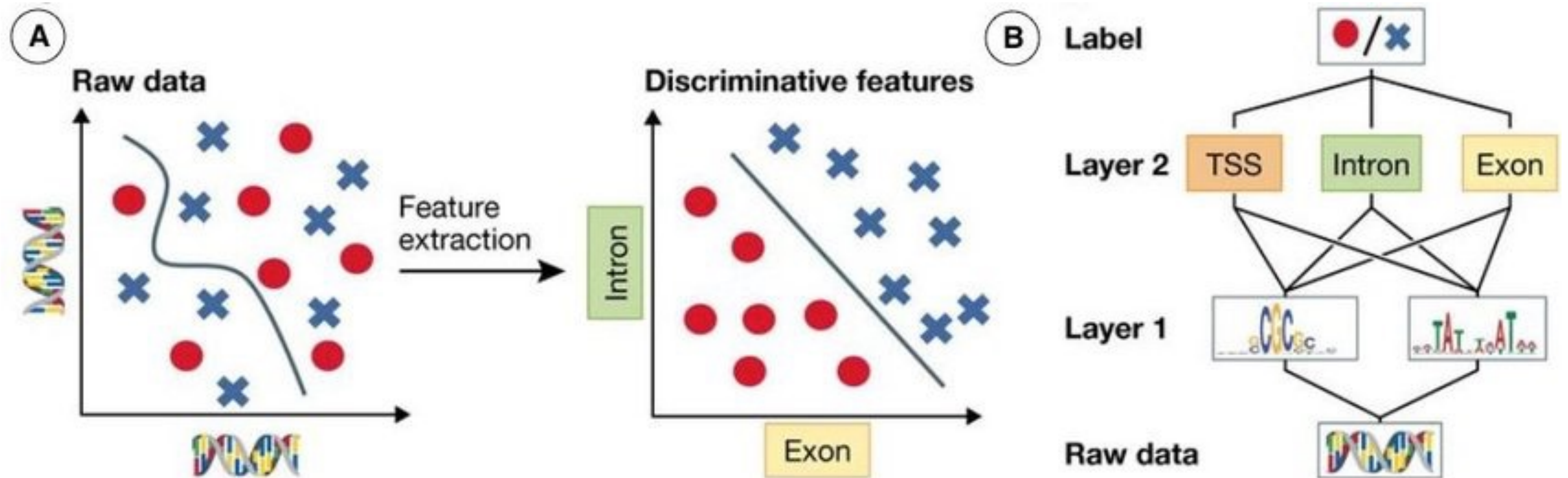
Deep Learning as feature extractors

Trunk versus head



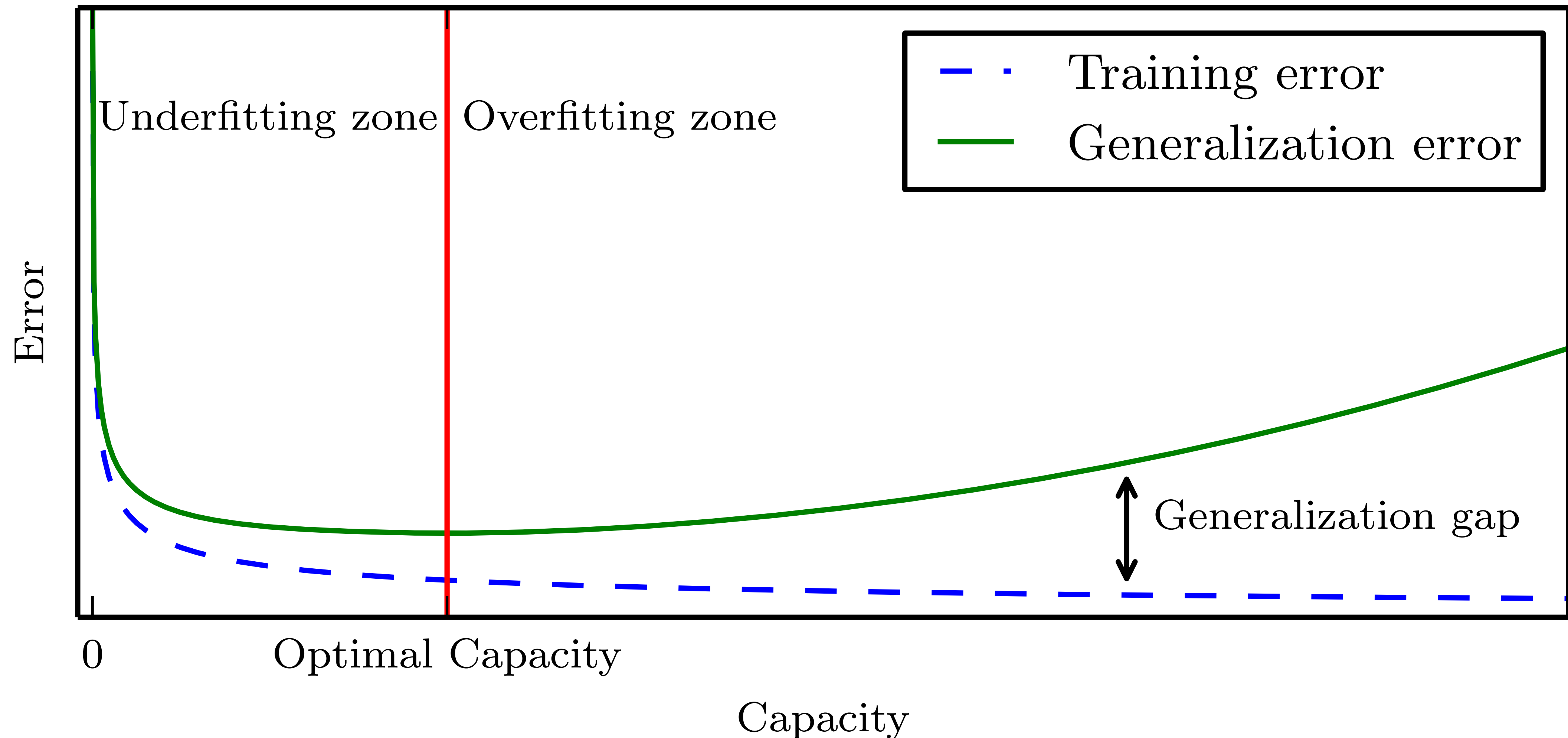
Deep Learning as feature extractors

Trunk versus head



Regularisation: How do we avoid overfitting?

Machine learning models like to overfit your data



Regularisation: How do we avoid overfitting?

Machine learning models like to overfit your data

Explicit Regularisation

-> Add reg. term to loss function:

- $\sum_{w \in \theta} \|w\|_1$
- $\sum_{w \in \theta} \|w\|_2$
- $\sum_{w \in \theta} \|w\|_\infty$
- $\sum_x p(x) \log(x)$

Implicit Regularisation

-> Various hacks to improve optimisation

- data augmentation
- dropout
- early stopping
- label smoothing
- Batch/Layer Normalisation



Takeaway



Neural networks perform **implicit feature extraction** and are optimized via **gradient descent**.

Outlook for Part II

1. **Images: Convolutional Neural Networks**
2. **Sequences: RNNs + Transformers**
3. **Current developments**