# Deep Learning Intuition

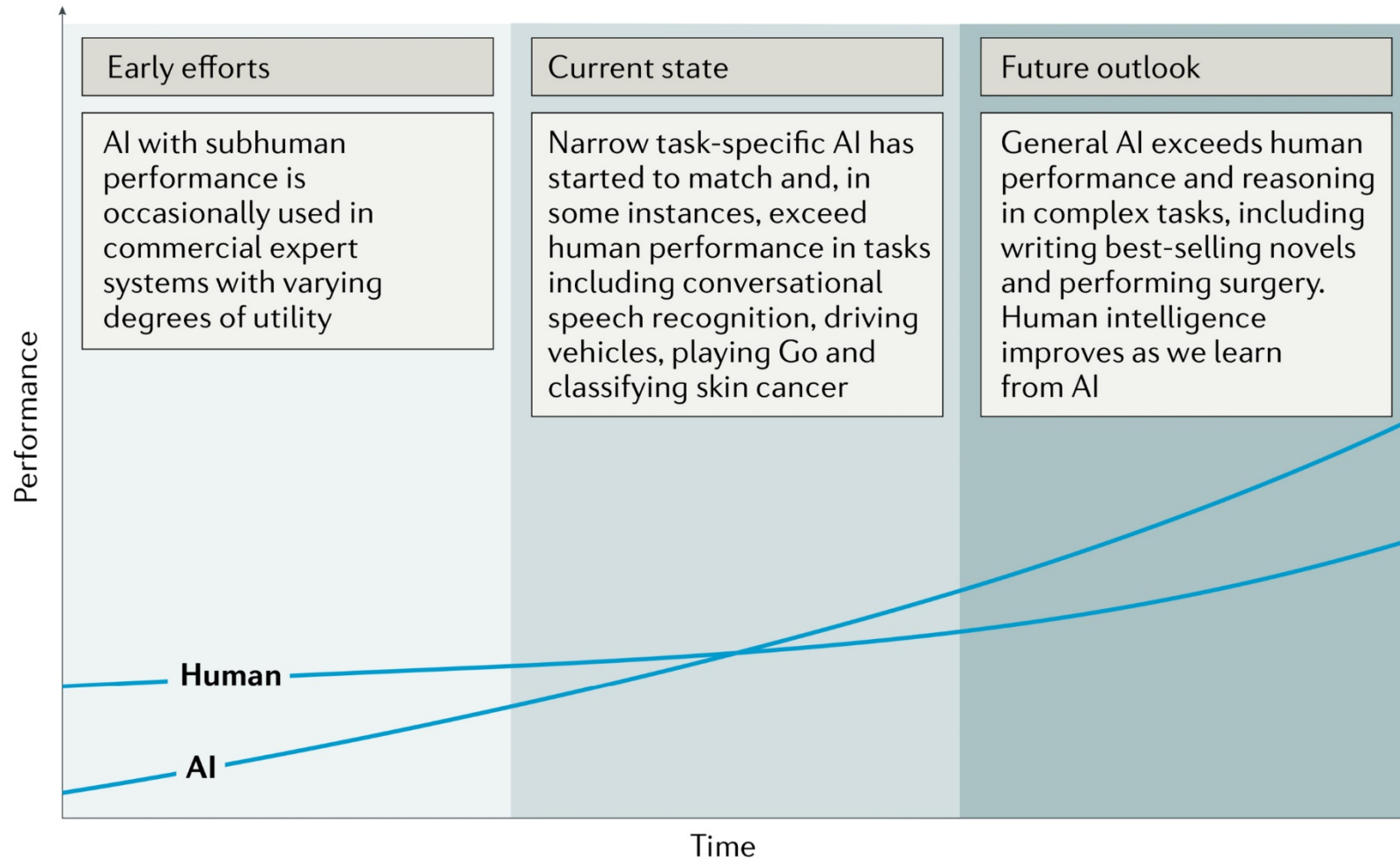*Ahmed Hosny*

HARVARD
MEDICAL SCHOOL

BWH
BRIGHAM AND
WOMEN'S HOSPITAL

DANA-FARBER
CANCER INSTITUTE

SAM Joint Imaging-Therapy Scientific Symposium (Certificate Series Session 2)
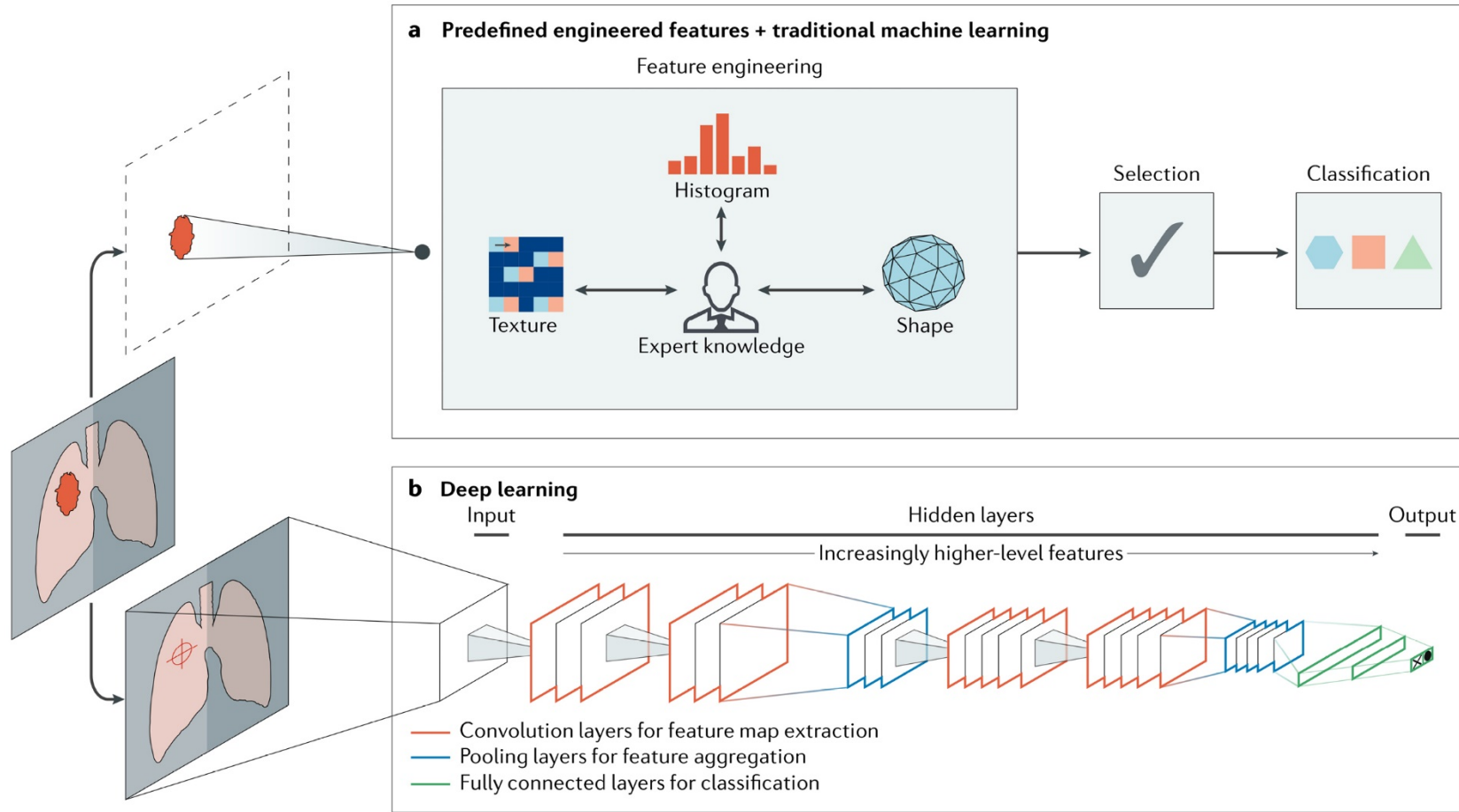Machine Learning for Radiomics - Wednesday, 8/1/2018 10:15 AM - 12:15 PM

# Deep Learning



Early efforts

AI with subhuman performance is occasionally used in commercial expert systems with varying degrees of utility

Current state

Narrow task-specific AI has started to match and, in some instances, exceed human performance in tasks including conversational speech recognition, driving vehicles, playing Go and classifying skin cancer

Future outlook

General AI exceeds human performance and reasoning in complex tasks, including writing best-selling novels and performing surgery. Human intelligence improves as we learn from AI

Performance

Human

AI

Time

*Ahmed Hosny , Chintan Parmar, John Quackenbush , Lawrence H Schwartz and Hugo JWL Aerts*

# Deep Learning



**a** Predefined engineered features + traditional machine learning

Feature engineering

Histogram

Texture — Expert knowledge — Shape

Selection

Classification

**b** Deep learning

Input | Hidden layers | Output

Increasingly higher-level features

- Convolution layers for feature map extraction
- Pooling layers for feature aggregation
- Fully connected layers for classification

*Ahmed Hosny , Chintan Parmar, John Quackenbush , Lawrence H Schwartz and Hugo JWL Aerts*

Artificial Intelligence in Radiology

# What is the intuition behind neural networks?

# How do neural networks learn?

# How to train neural networks?

# What is the intuition behind neural networks?

How do neural networks learn?

How to train neural networks?

# Machine Learning: 4 Main Components

- Data
- Model/Representation
- Cost/Error/Loss of model
- Model Optimizer

# Logistic Regression

# Logistic Regression



$$y = wx + b$$

nodule concavity

nodule radius

malignant
benign

# Logistic Regression

$$y = wx + b$$

# Logistic Regression

$$y = wx + b$$

nodule concavity

nodule radius

malignant

benign

$log\ loss = 3*\ \bigcirc\ + 23*\bullet + 14*\ \bigcirc\ + 11*\bullet$

# Logistic Regression

$$y = wx + b$$



log loss = 3* ⬤ + 23* • + 14* ⬤ + 11* •

# Logistic Regression

# Logistic Regression

# Logistic Regression



$$y = wx + b$$

nodule concavity

- malignant
- benign

nodule radius

log loss

$w$

log loss

$b$

$$log\ loss = 2*\bigcirc + 24*\bullet + 8*\bigcirc + 17*\bullet$$

# Logistic Regression



$$y = wx + b$$

nodule concavity

nodule radius

● malignant
● benign

log loss

$w$

log loss

$b$

$$log\ loss = 26^* ● + 25^* ●$$

# Dealing with Edge Conditions

# Dealing with Edge Conditions

# Dealing with Edge Conditions

# Dealing with Edge Conditions

# Dealing with Edge Conditions

# Dealing with Edge Conditions

# Quadrant Questioning

# Quadrant Questioning



- Line A, over or under?

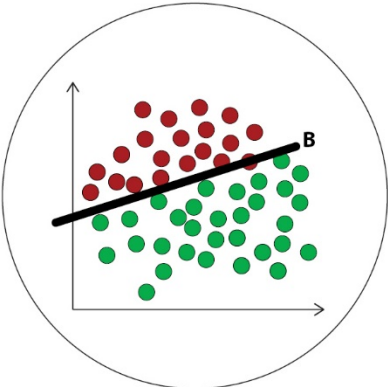# Quadrant Questioning



- Line A, over or under?
- Line B, over or under?

# Quadrant Questioning



- Line A, over or under?
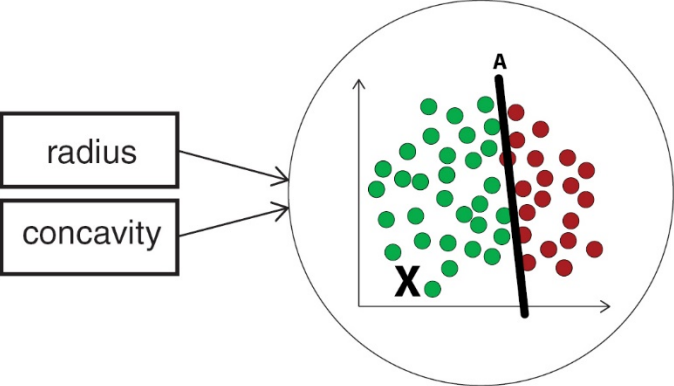- Line B, over or under?
- Both YES?

# Graph Representation

Line A, over or under?
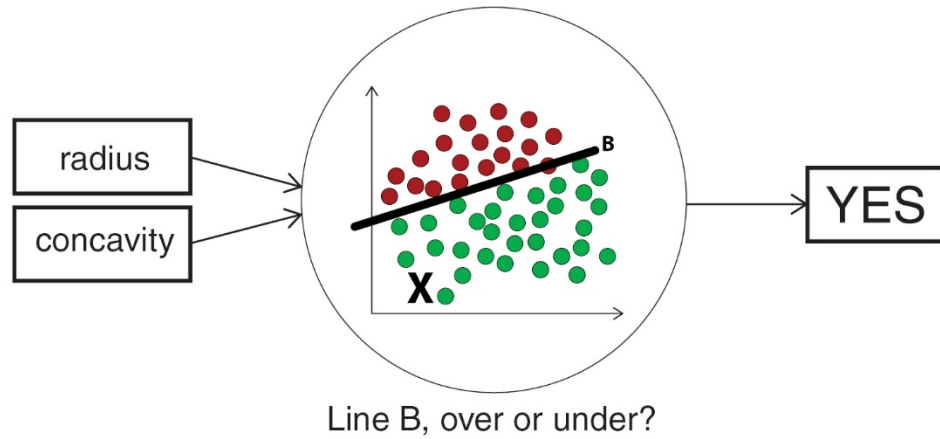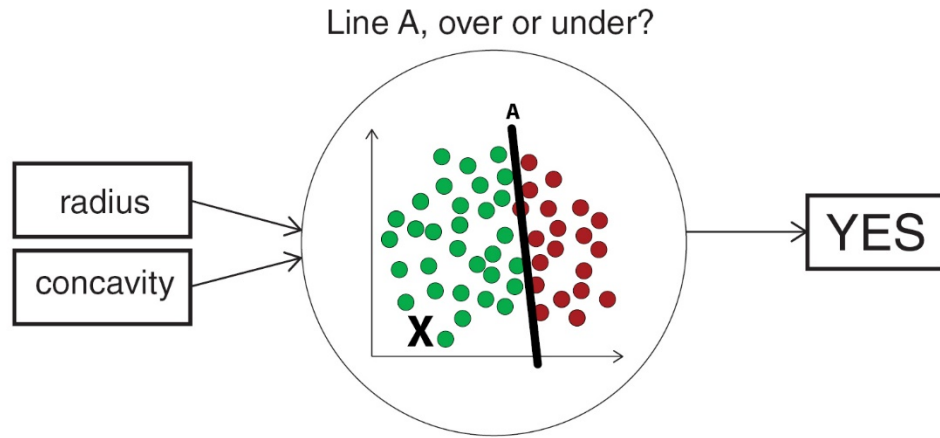


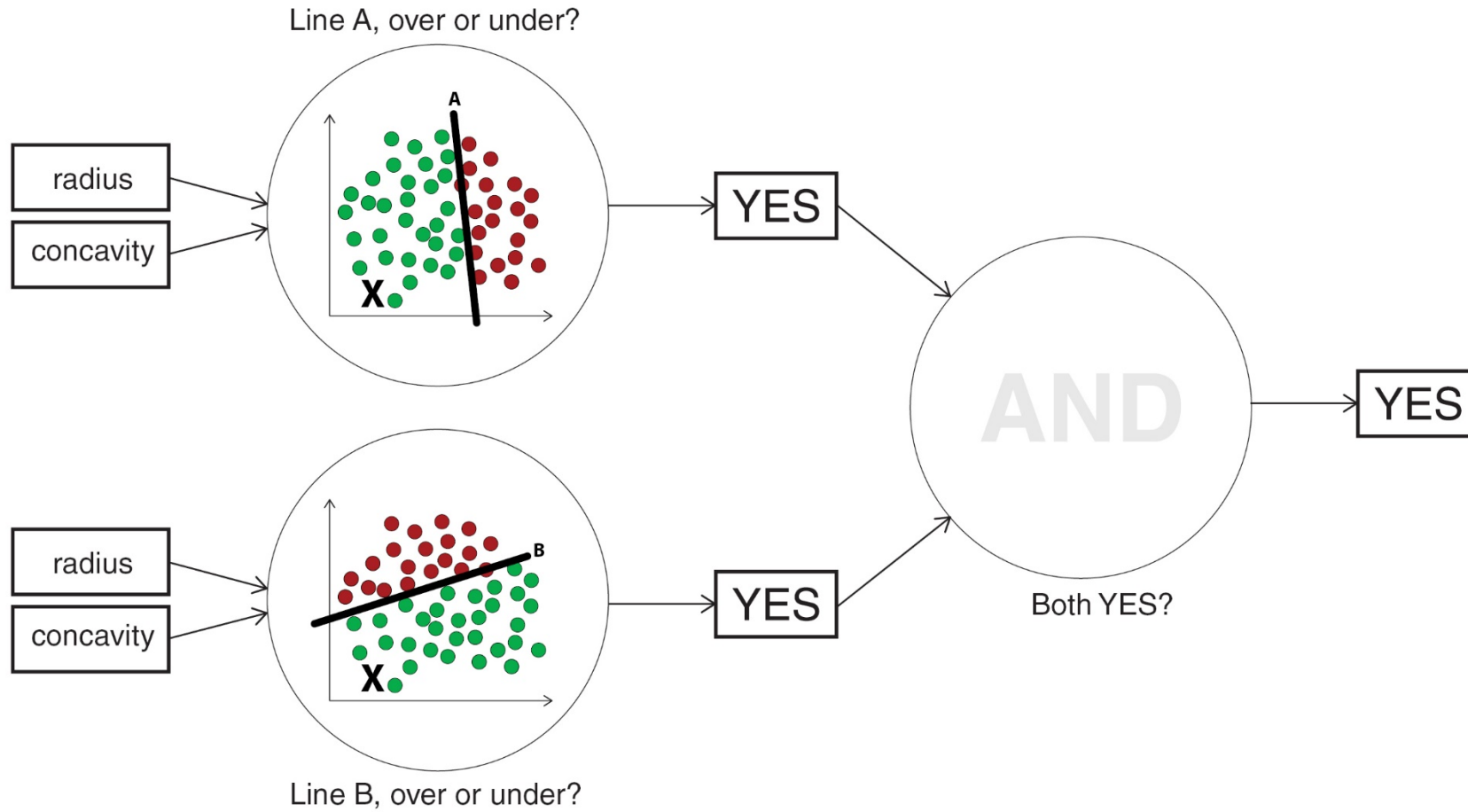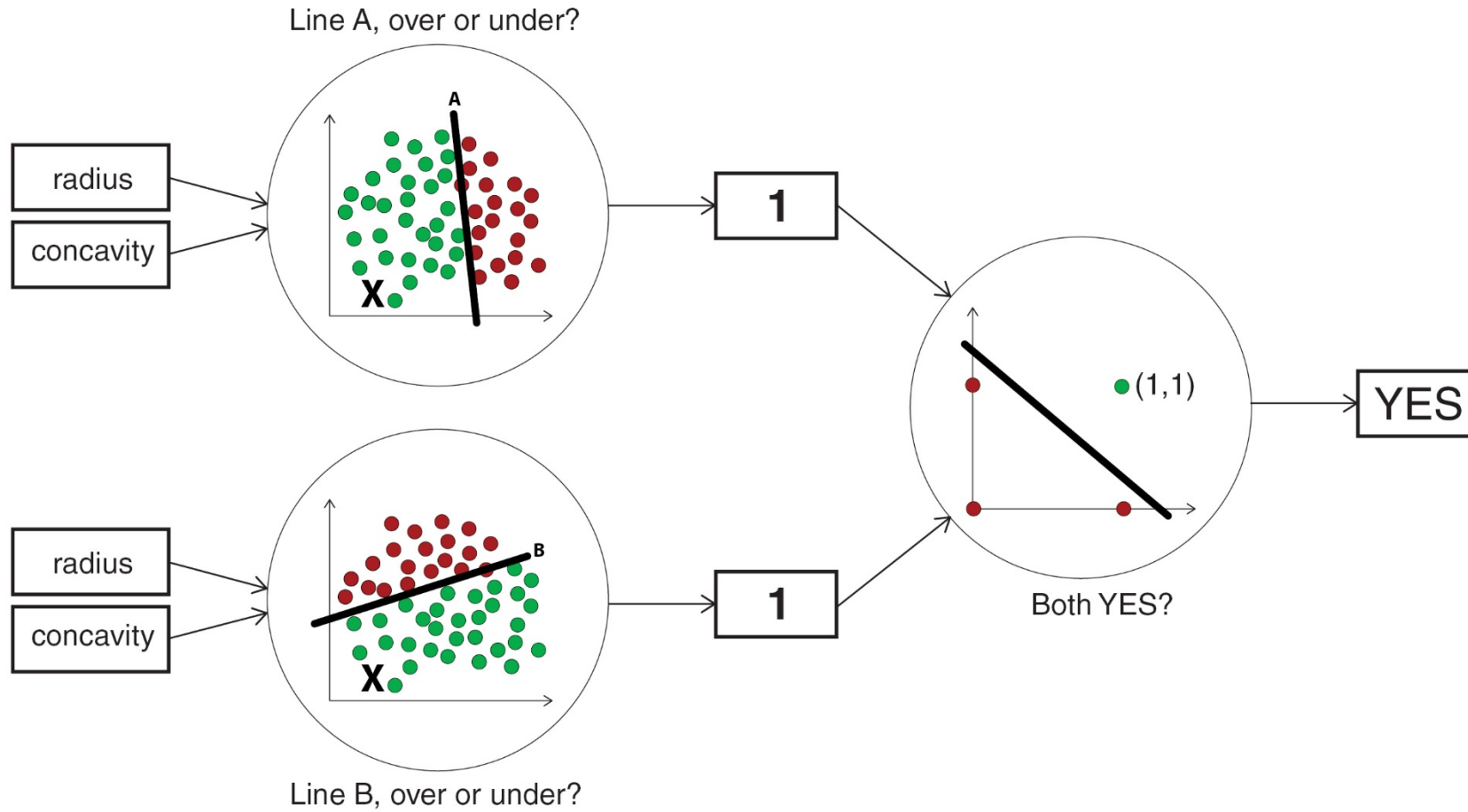Line B, over or under?

# Graph Representation
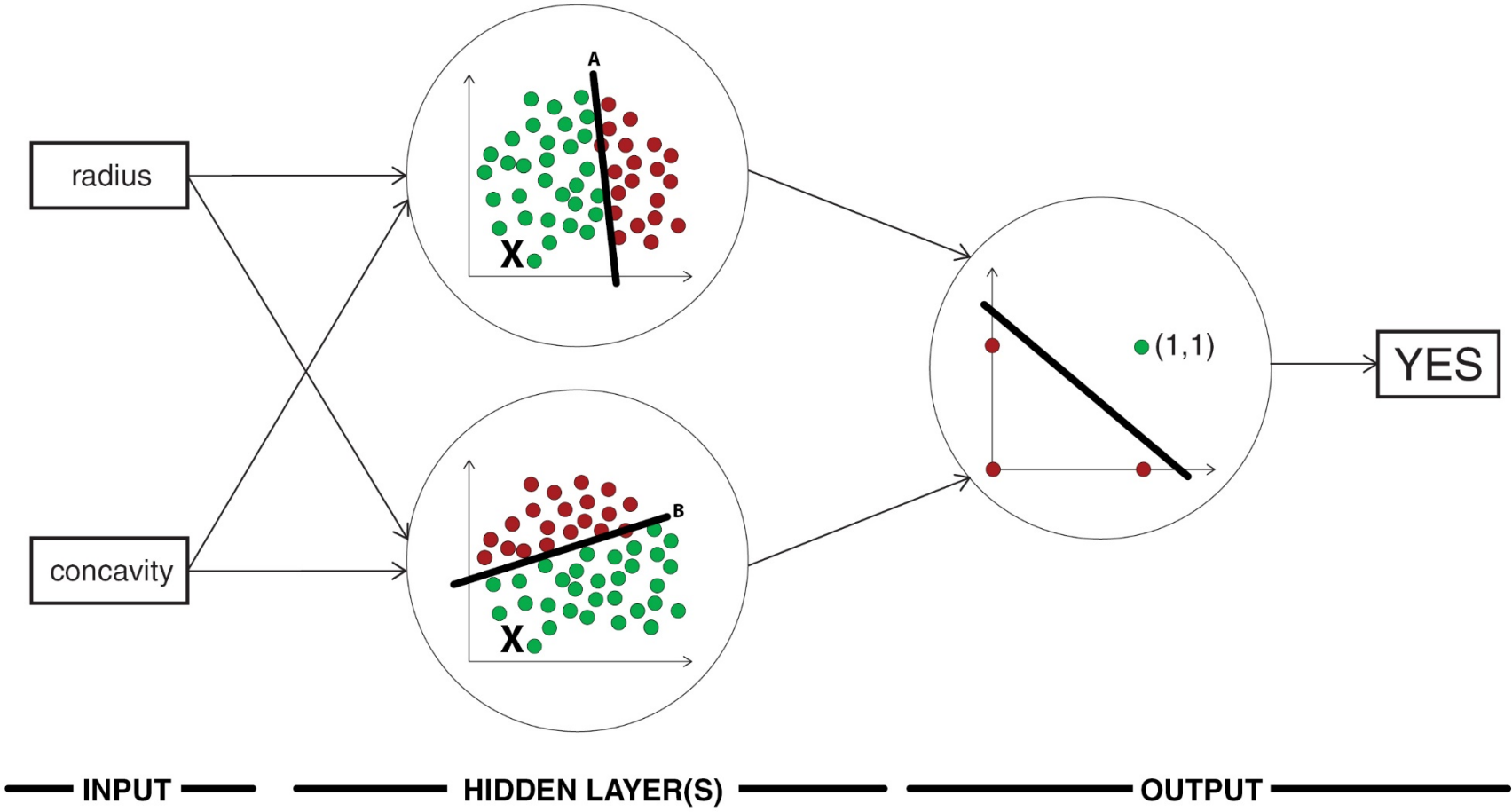
Line A, over or under?

radius

concavity



Line B, over or under?

# Graph Representation

Line A, over or under?

radius

concavity

YES

X

A

radius

concavity

YES

X

B

Line B, over or under?

# Graph Representation



Line A, over or under?

radius

concavity

YES

AND

YES

Both YES?

radius

concavity

YES

Line B, over or under?

# Graph Representation

Line A, over or under?

radius

concavity

**1**

Line B, over or under?

radius

concavity

**1**

(1,1)

Both YES?

YES

# A Neural Network



INPUT — HIDDEN LAYER(S) — OUTPUT

# A Neural Network



radius

concavity

$W1^{\text{radius}}$

$W2^{\text{radius}}$

$W1^{\text{concavity}}$

$W2^{\text{concavity}}$

$W1^{\text{radius}} * radius$
$+$
$W1^{\text{concavity}} * concavity$
$+$
$b$

$W2^{\text{radius}} * radius$
$+$
$W2^{\text{concavity}} * concavity$
$+$
$b$

$f(h)$

$f(h)$
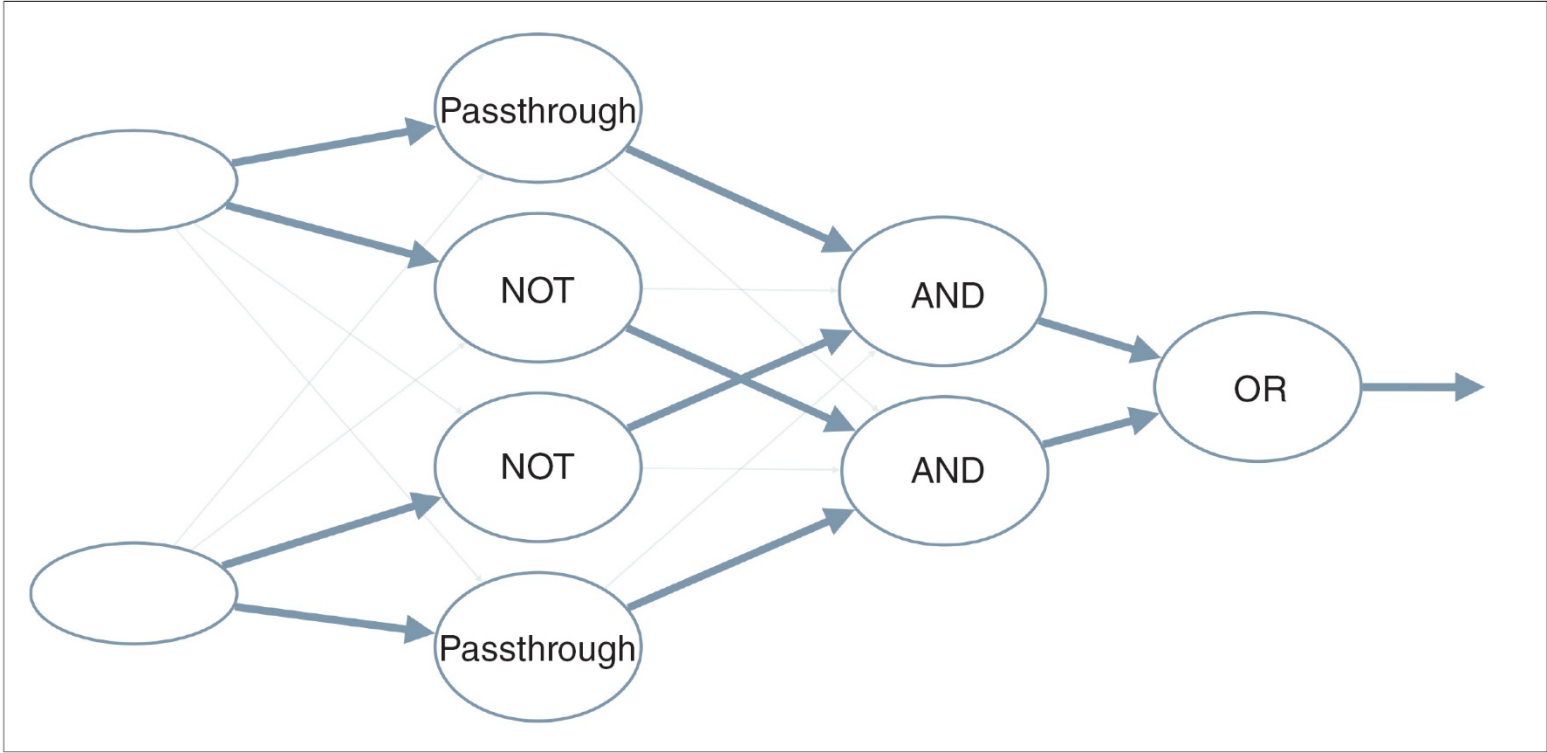
$(1,1)$

YES

**INPUT**  **HIDDEN LAYER(S)**  **OUTPUT**

# XOR Perceptron
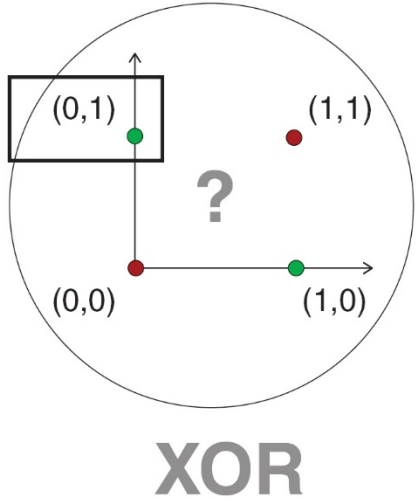


AND

# XOR Perceptron



AND          OR

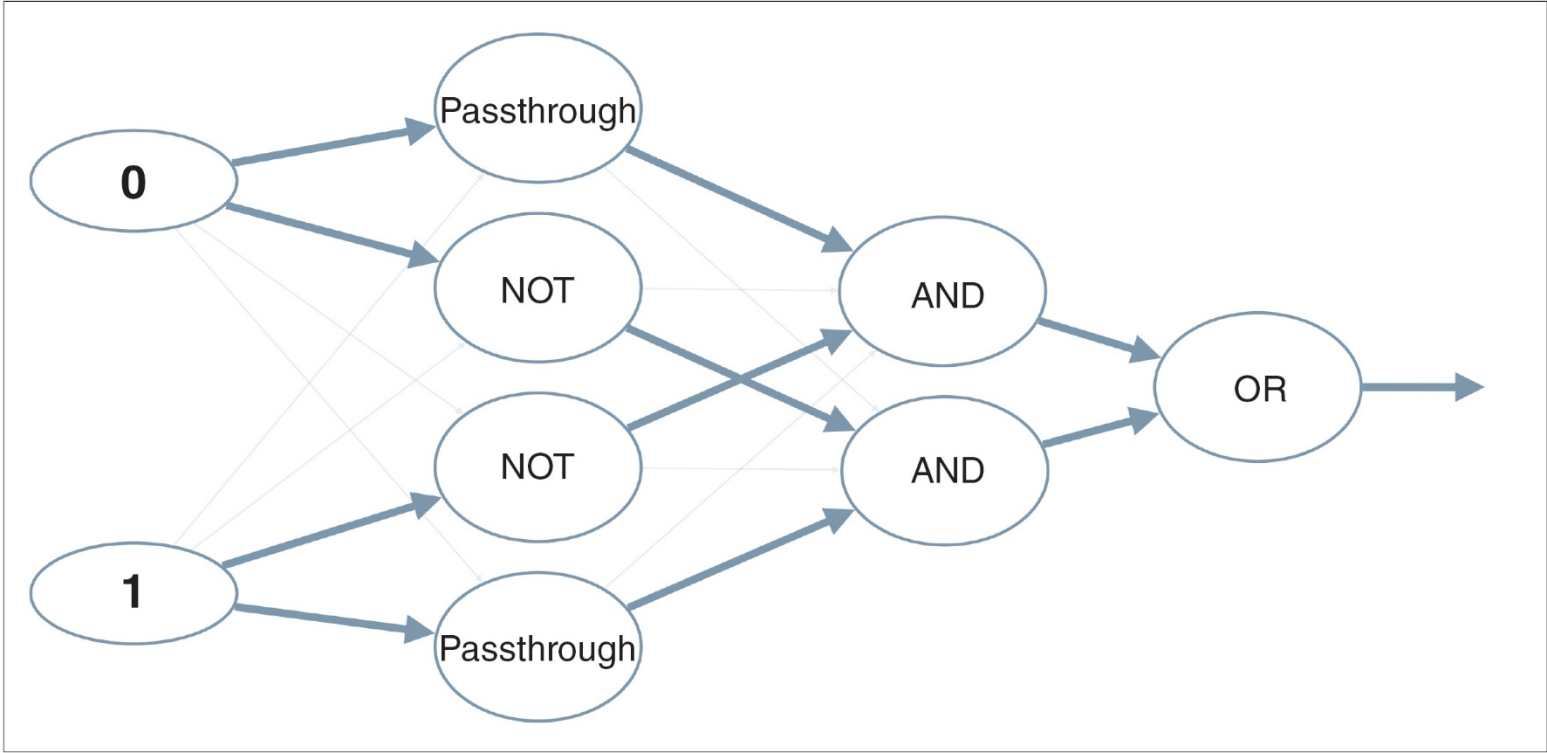# XOR Perceptron



AND        OR        XOR
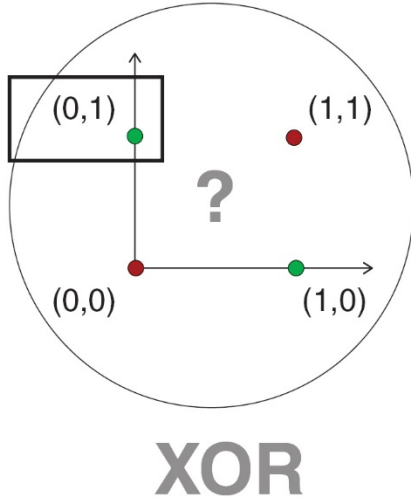
# XOR Perceptron
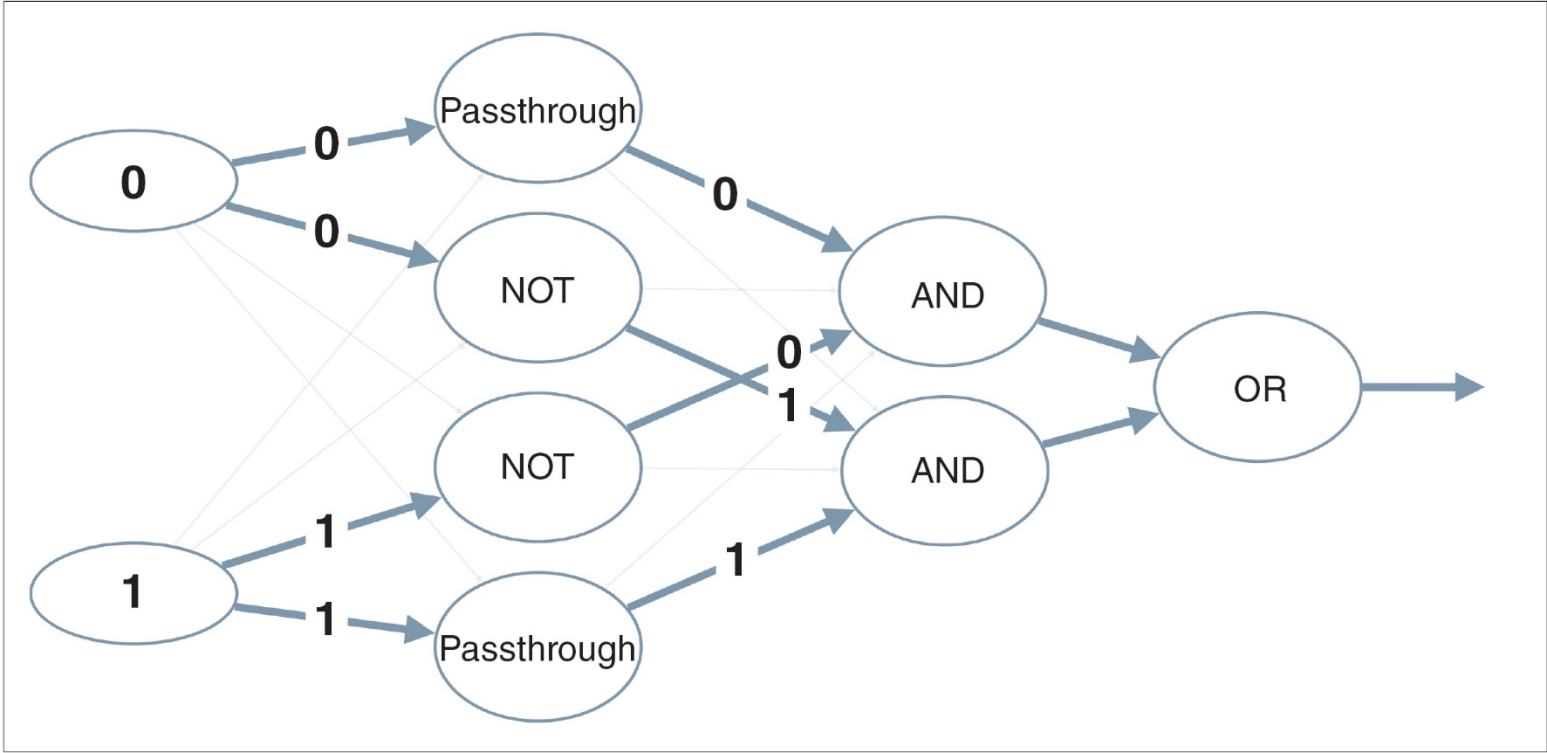
(0,1)    (1,1)

?

(0,0)    (1,0)

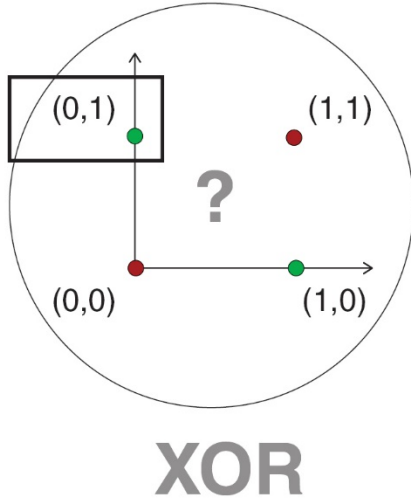**XOR**

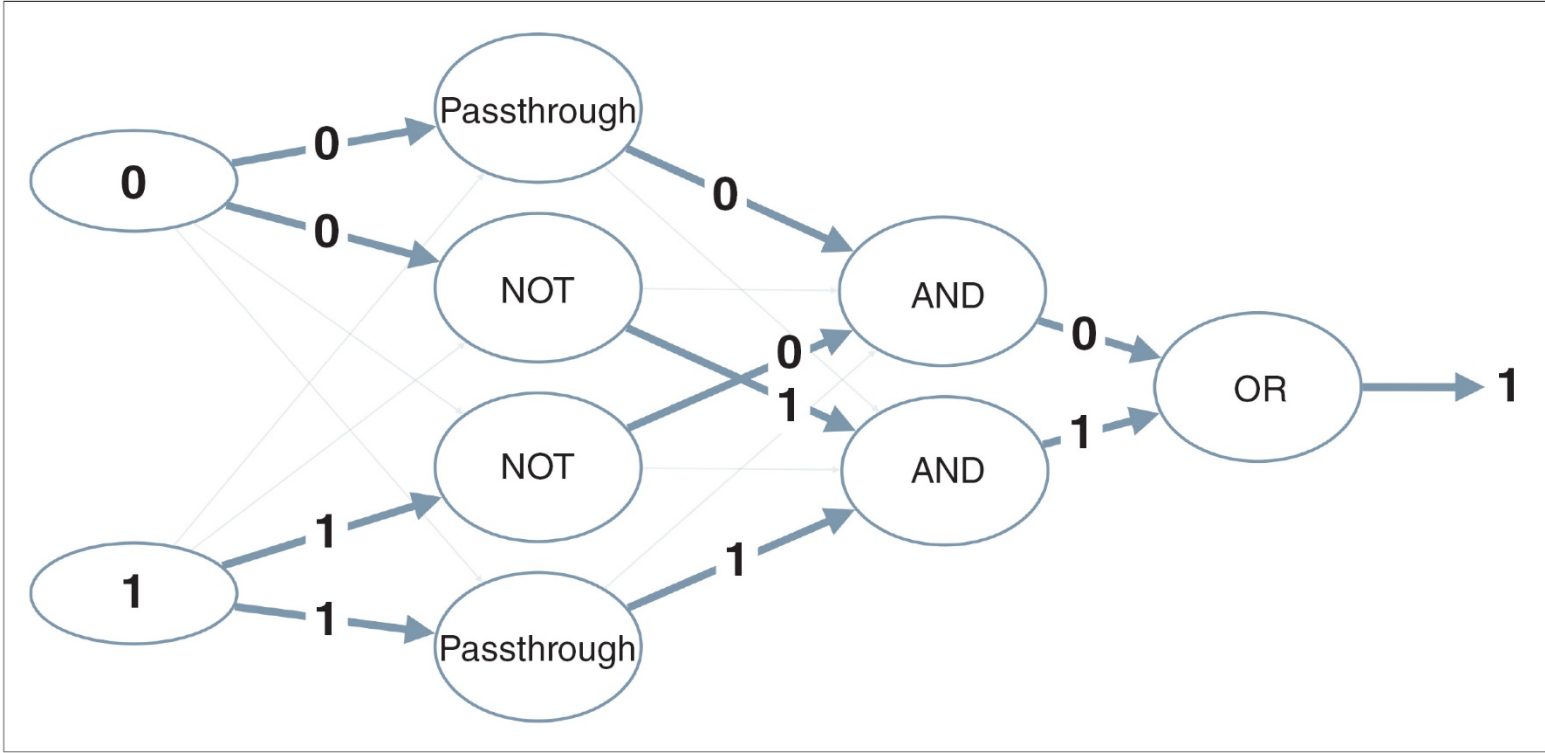# XOR Perceptron

# XOR Perceptron
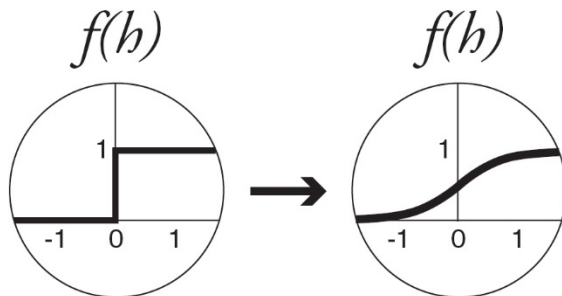
# XOR Perceptron

# XOR Perceptron

What is the intuition behind neural networks?

# How do neural networks learn?

How to train neural networks?

# Backpropagation & Gradient Descent in Neural Networks



*David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams*

## Learning Representations by Back-propagating Errors
**Nature - 1986**

# Iris Dataset

**Inputs**



petal length & width
sepal length & width

**Outputs**



Iris Versicolor      Iris Setosa      Iris Virginica

# How Neural Networks Learn

✔ **Data:** iris dataset

**label**

**features**

petal length

petal width

sepal length

sepal width

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network

**label**

**features**

petal length

petal width

sepal length

sepal width

**sigmoid**

$$S(x) = \frac{1}{1 + e^{-x}}$$

**softmax**

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}} \quad \text{for } j = 1, \ldots, K.$$

sigmoid activation

softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy

↑ **likelihood**

=

↑ **log(likelihood)**

=

↓ **-log(likelihood)**

**label**

cross entropy

**features**

petal length

petal width

sepal length

sepal width

sigmoid activation

softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

**1.** parameter initialization

**label**

**cross entropy**

↓

**gradient descent**

**features**

petal length

petal width

sepal length

sepal width

sigmoid activation

softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

**1.** parameter initialization
**2.** data input



**label**

cross entropy

gradient descent

**features**

petal length  3.2
petal width  2.1
sepal length  3.6
sepal width  1.7

sigmoid activation
softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation

**label**

| 1 |
|---|
| 0 |
| 0 |

→ cross entropy ←
↓
gradient descent

**current**

| 1 |
|---|
| 0 |
| 0 |

0.15
0.55
0.30

-log(0.15)
0.82

**features**

petal length  3.2
petal width  2.1
sepal length  3.6
sepal width  1.7

0.27
0.67
0.13
0.90
0.49

0.15
0.55
0.30

○ sigmoid activation
○ softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation



**label**

cross entropy

↓

gradient descent

**features**

petal length 3.2
petal width 2.1
sepal length 3.6
sepal width 1.7

0.27
0.67
0.13
0.90
0.49

0.15
0.55
0.30

sigmoid activation

softmax activation

**current**    **best**

0.15    1.00
0.55    0.00
0.30    0.00

-log(0.15)    -log(1)
0.82    0.00

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation



**label**

cross entropy
↓
gradient descent

**features**

petal length
petal width
sepal length
sepal width

sigmoid activation
softmax activation

| | current | best | worst |
|---|---|---|---|
| | 0.15 | 1.00 | 0.00 |
| | 0.55 | 0.00 | 0.60 |
| | 0.30 | 0.00 | 0.40 |
| | -log(0.15) | -log(1) | -log(0) |
| | 0.82 | 0.00 | -(-∞) = ∞ |

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

$$\Delta w \propto - \text{gradient}$$

$$= - \frac{\partial L}{\partial w}$$

$$= - \eta \frac{\partial L}{\partial w}$$

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

$$\Delta w \propto - \text{gradient}$$

$$= - \frac{\partial L}{\partial w}$$

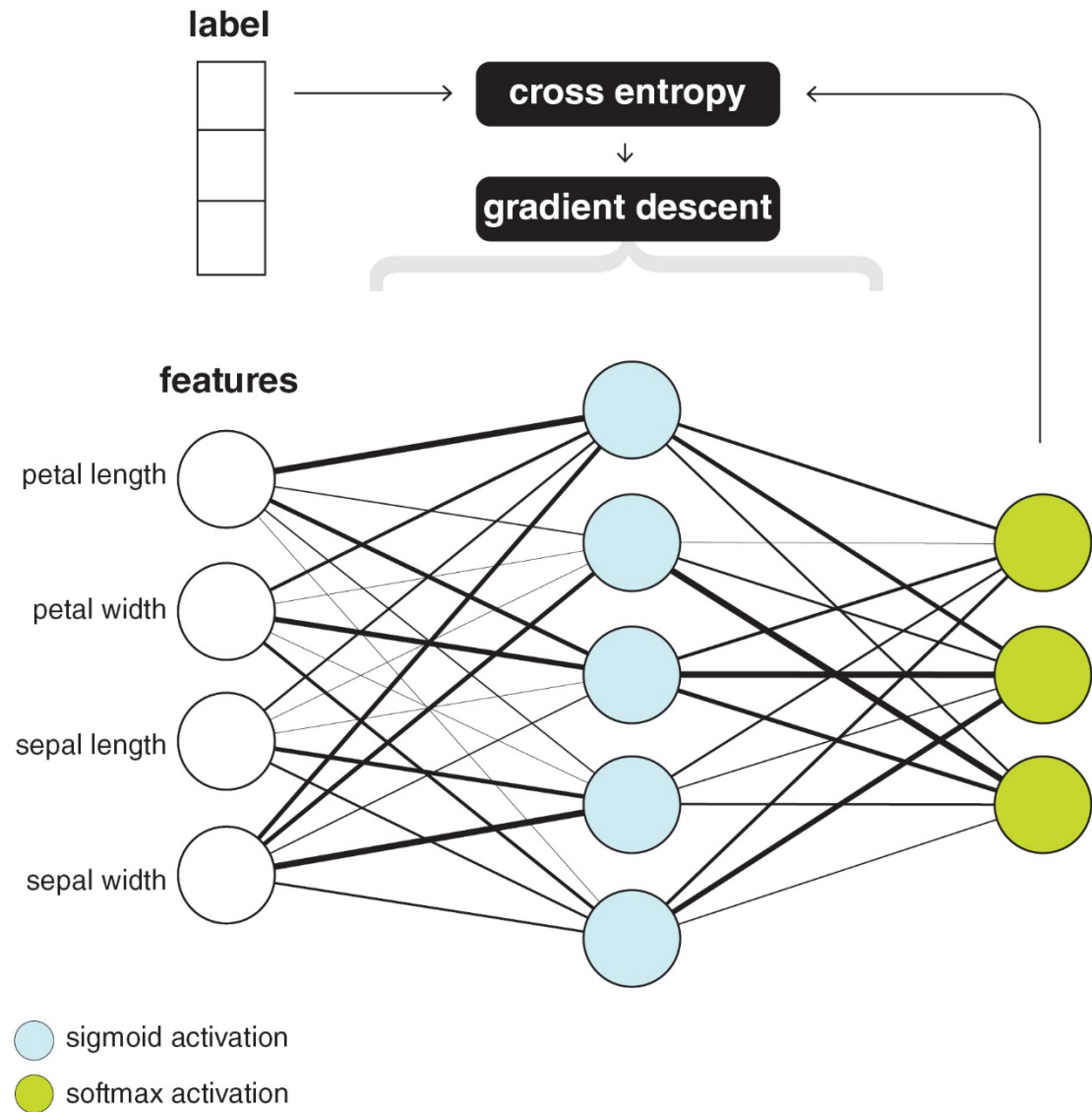$$= - \eta \frac{\partial L}{\partial w}$$

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
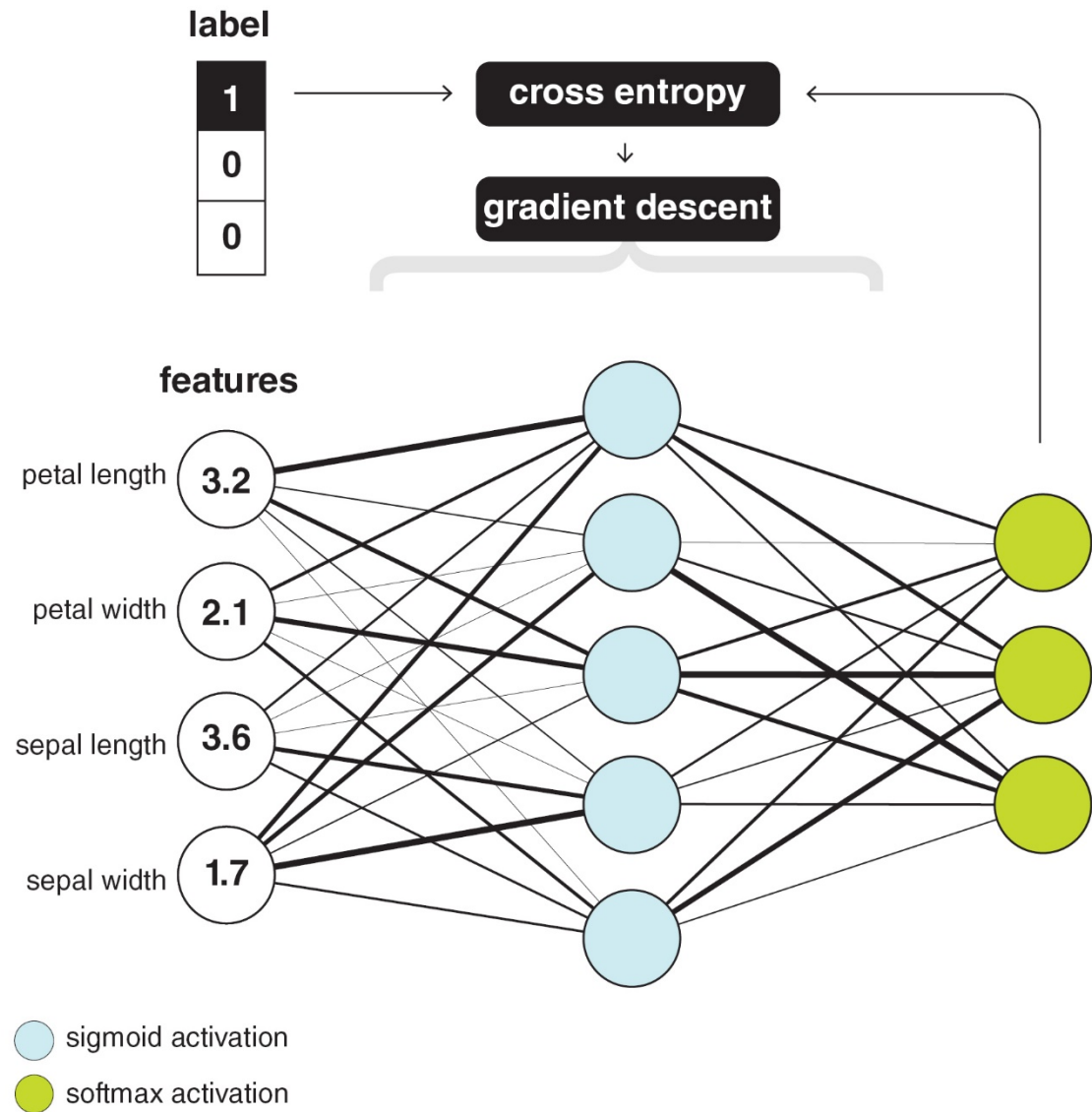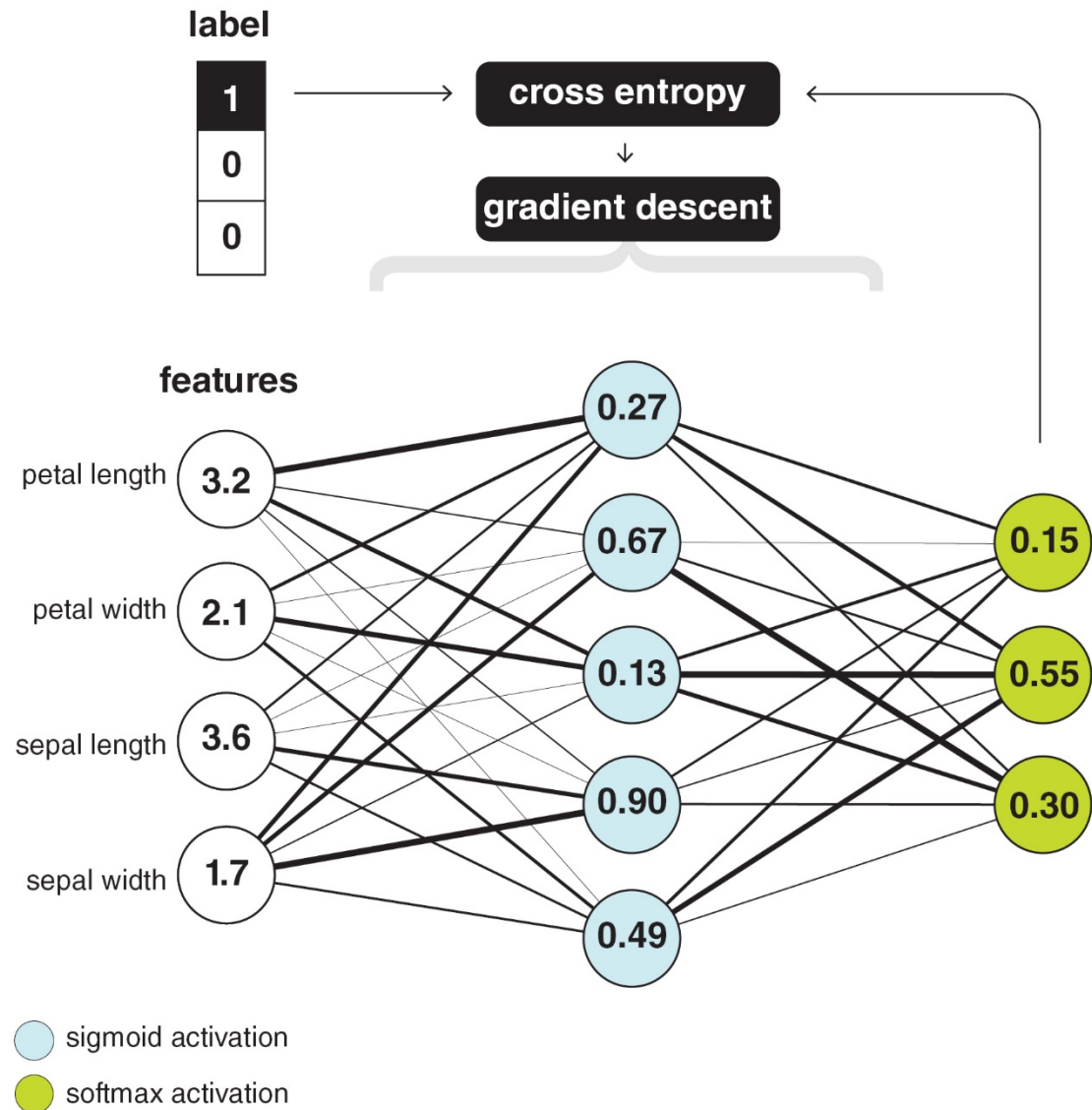✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

f( f( input ) )

**label**

| 1 |
|---|
| 0 |
| 0 |

cross entropy
↓
gradient descent

**features**

petal length **3.2**
petal width **2.1**
sepal length **3.6**
sepal width **1.7**

0.27
0.67
0.13
0.90
0.49

0.15
0.55
0.30

sigmoid activation
softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
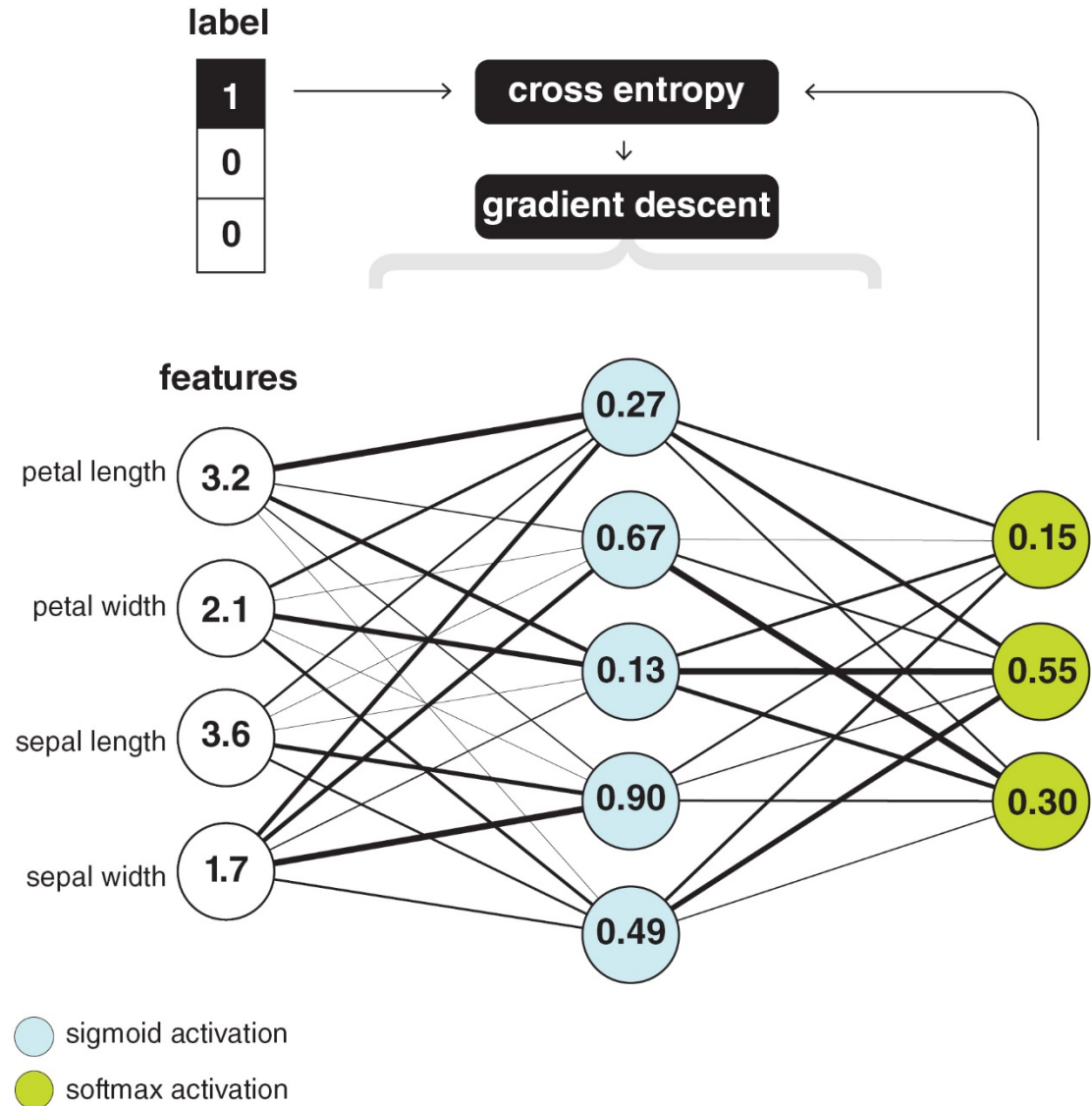✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
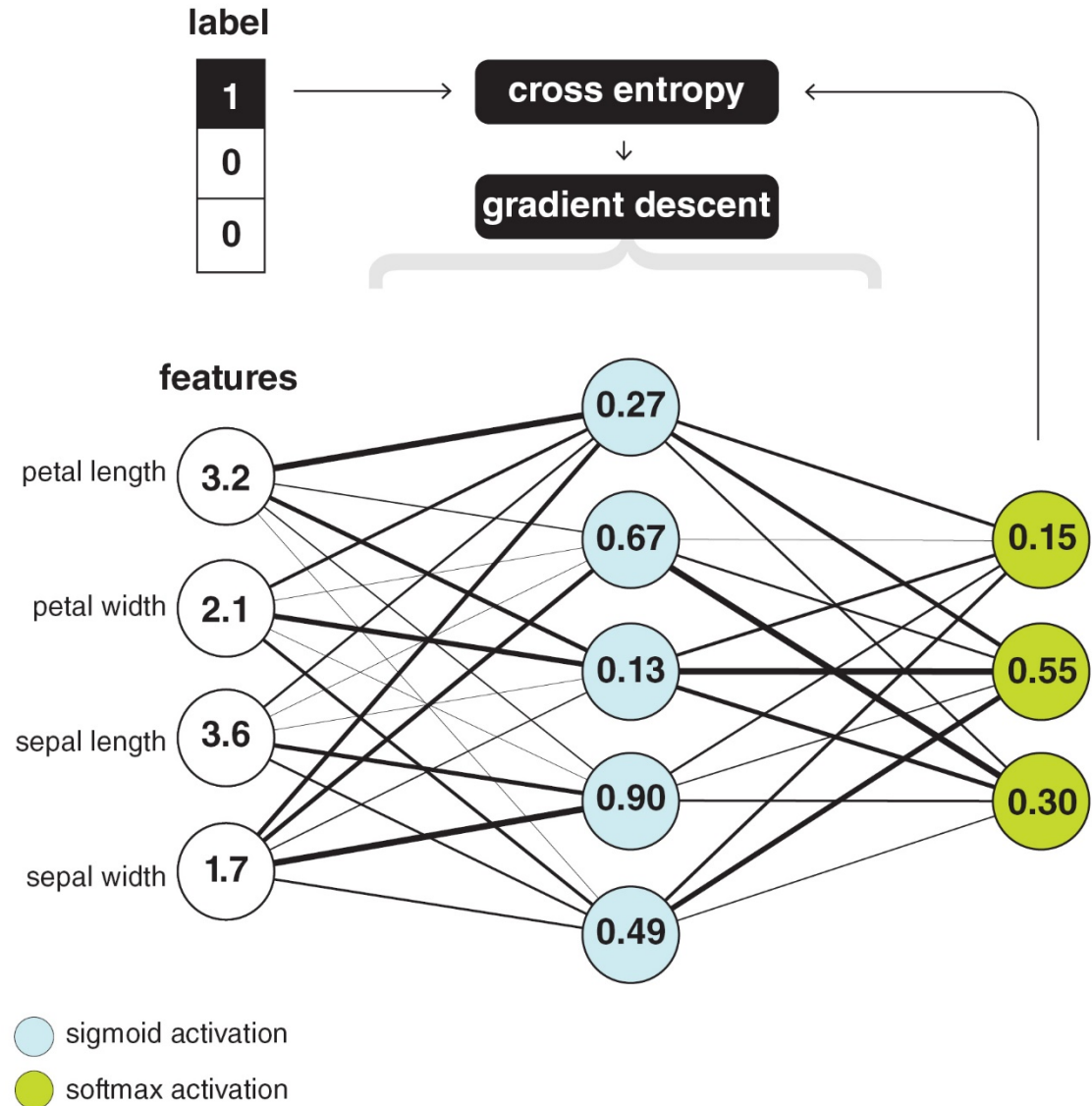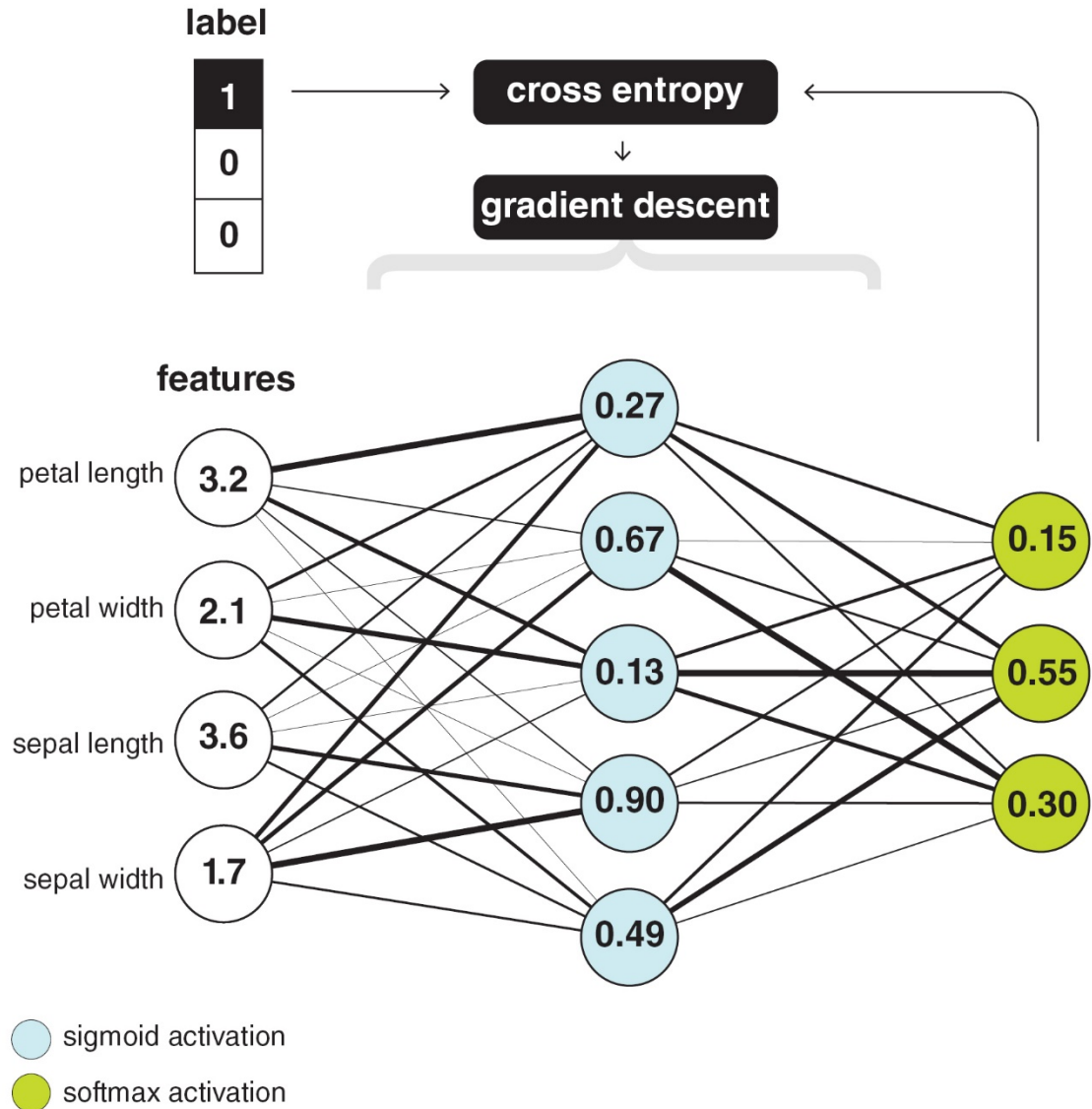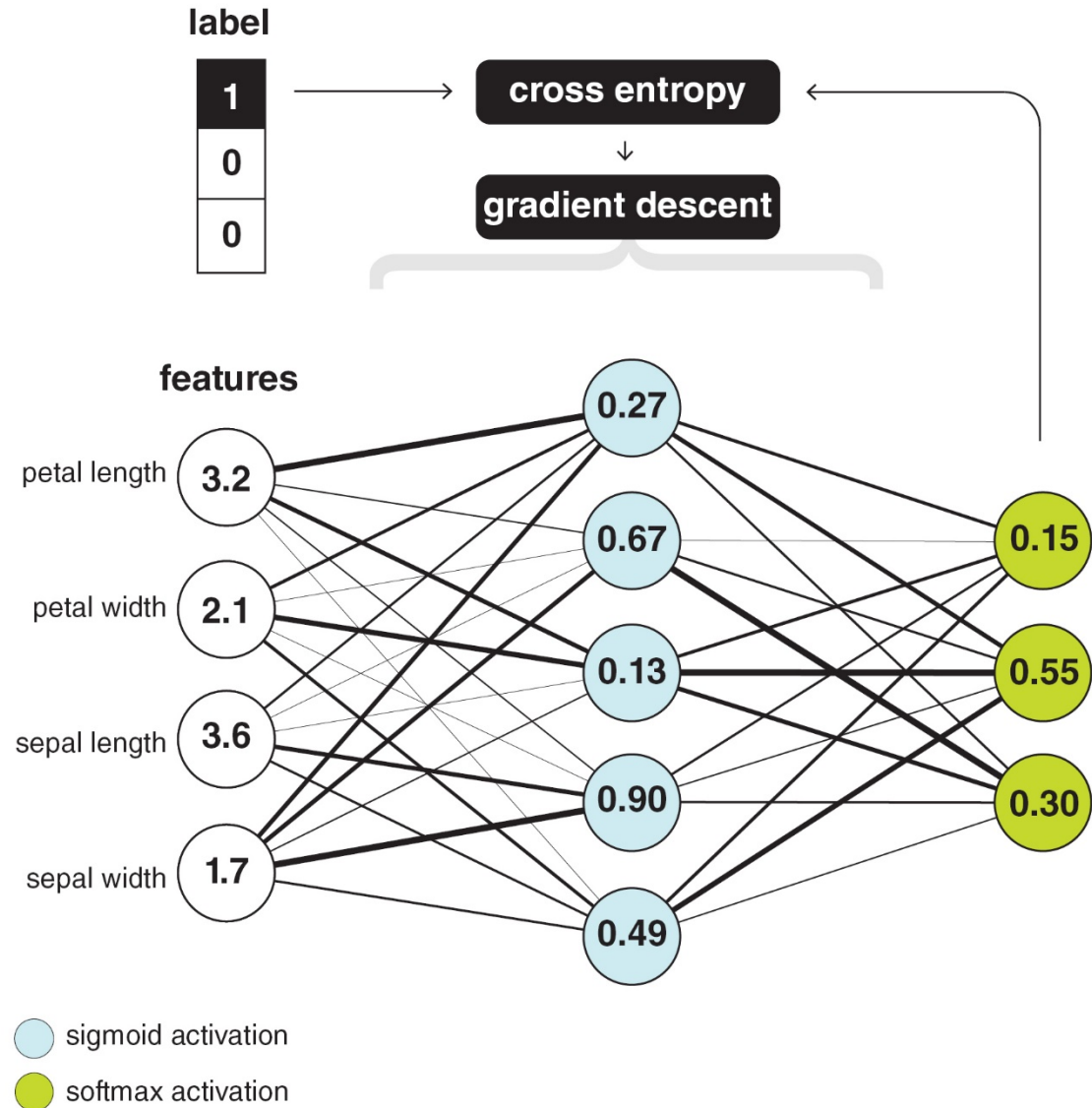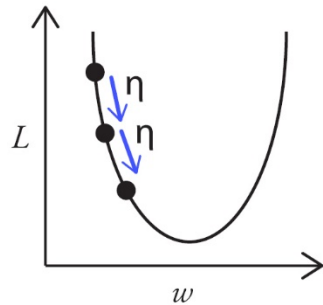5. backpropagation + updates

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

---

**1.** parameter initialization
**2.** data input
**3.** forward propagation
**4.** loss calculation
**5.** backpropagation + updates

---

$$\Delta w_n = -\eta \, \frac{\partial L}{\partial w_n}$$

$$= -\eta \, \frac{\partial L}{\partial so} \frac{\partial so}{\partial n} \frac{\partial n}{\partial w_n}$$

**label**

**cross entropy**

↓

**gradient descent**

**features**

petal length **3.2**

petal width **2.1**

sepal length **3.6**

sepal width **1.7**

0.27

0.67

0.13

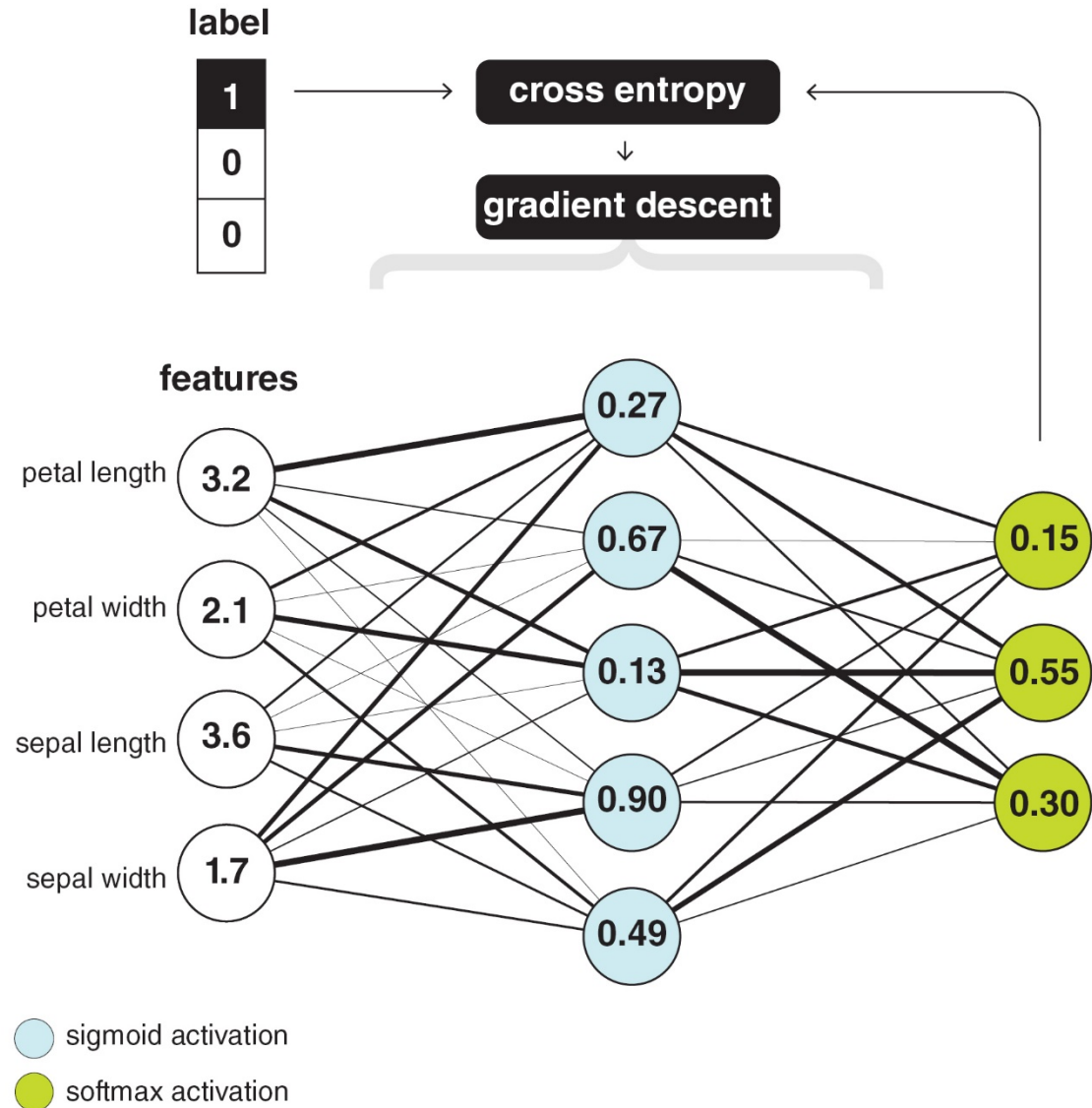0.90

0.49

$L$

$w_n$

$n$

0.15

0.55

0.30

sigmoid activation

softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
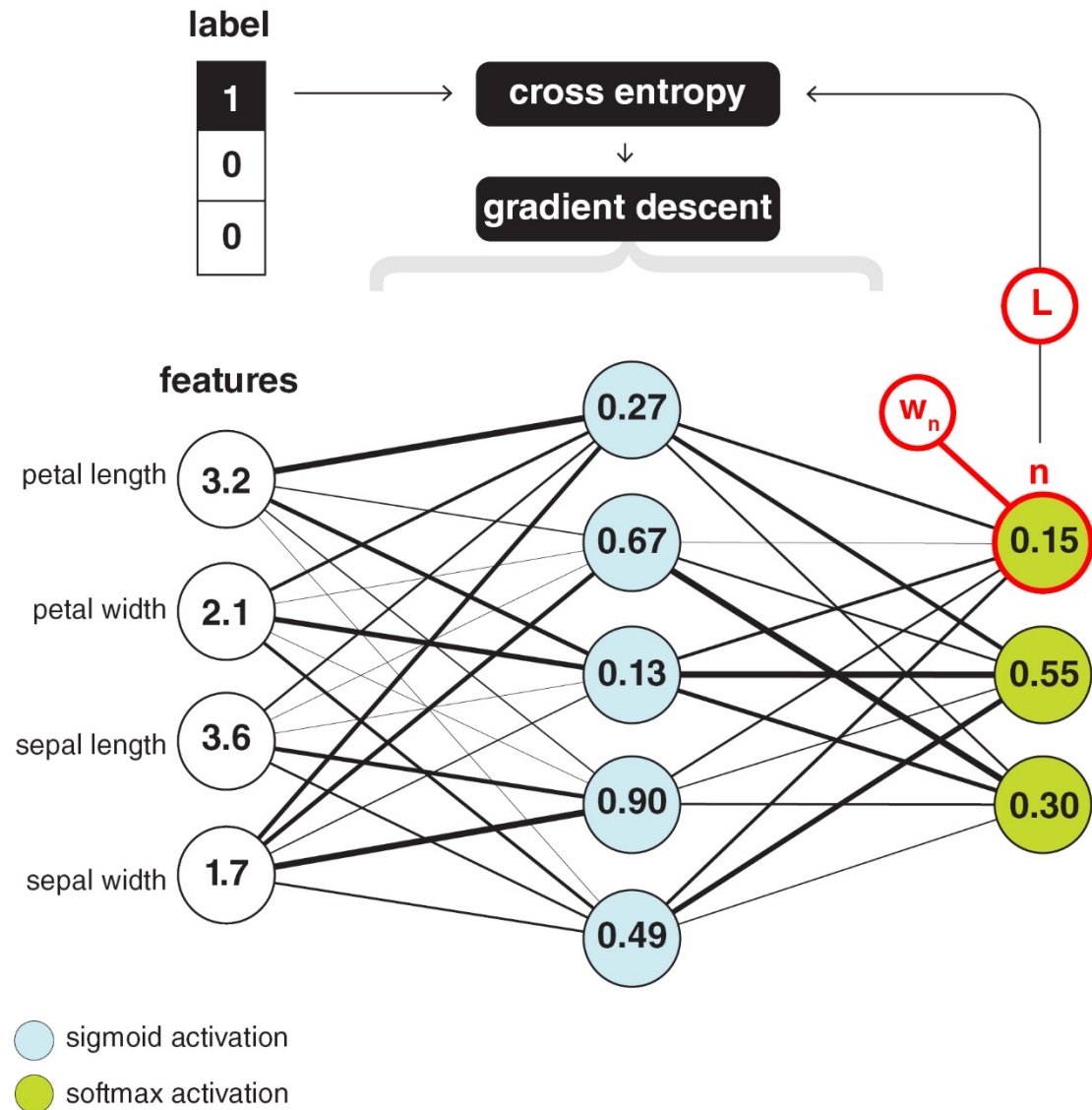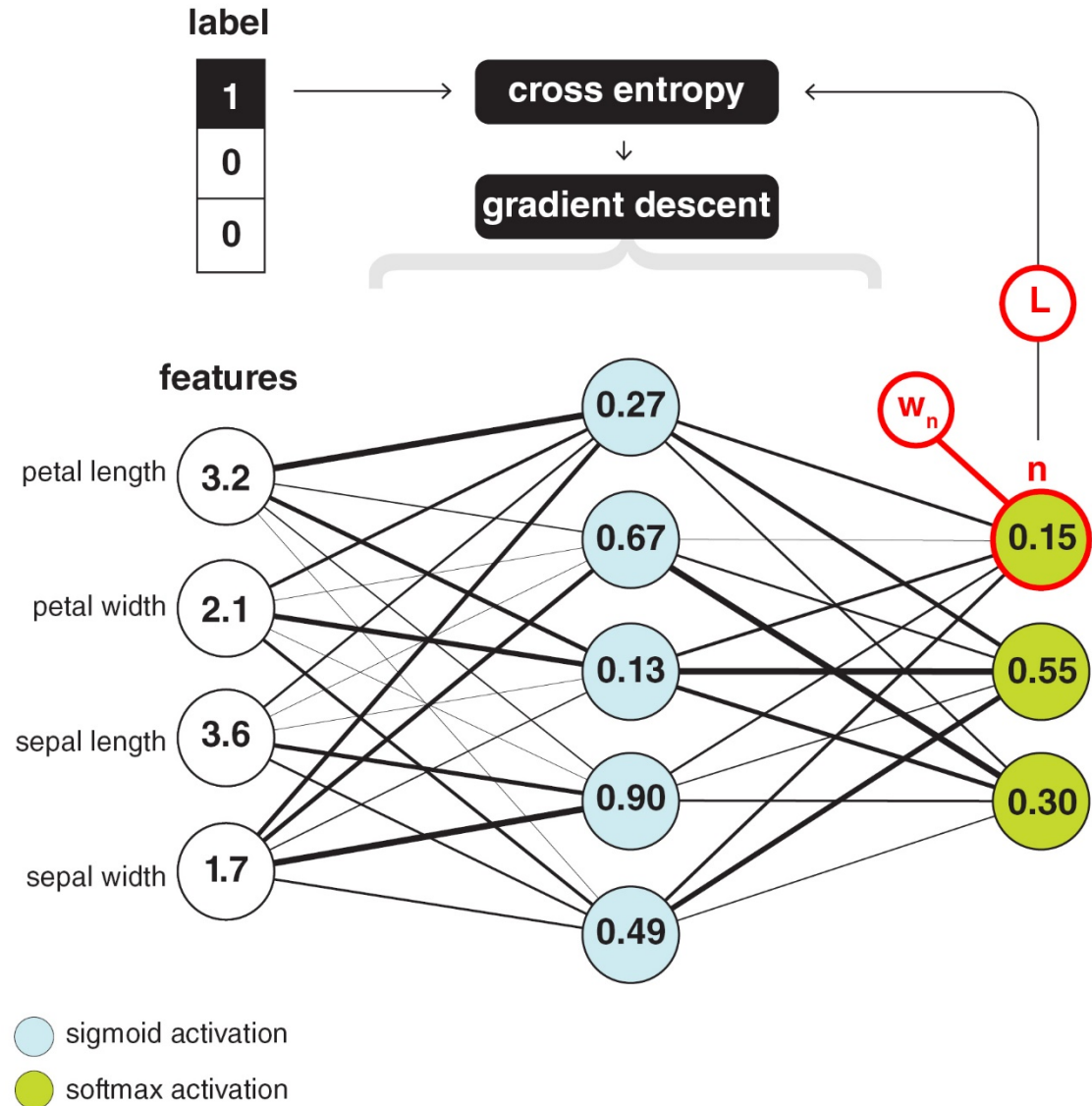4. loss calculation
5. backpropagation + updates

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

---

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
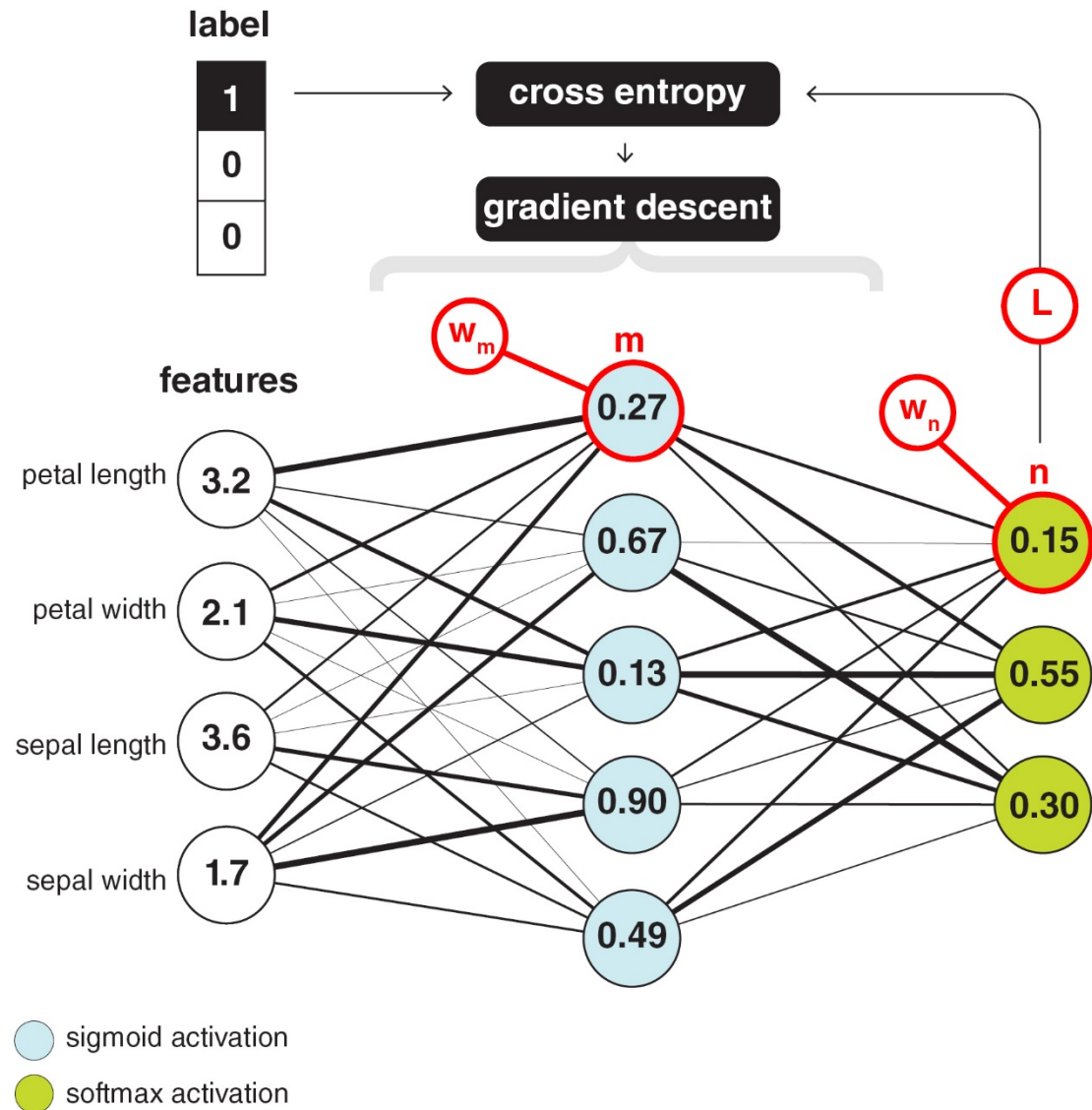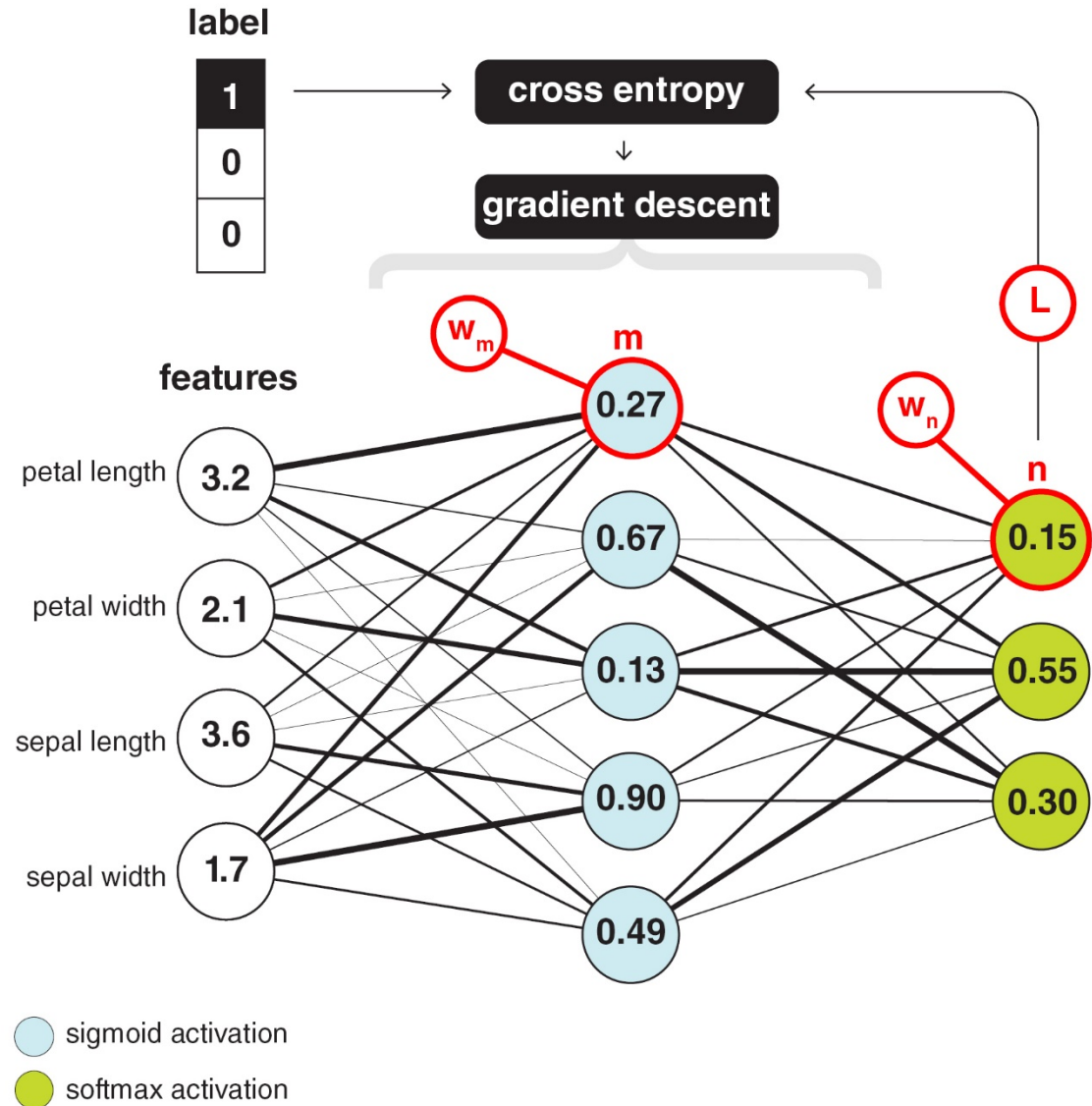5. backpropagation + updates

---

$$\Delta w_m = -\eta \, \frac{\partial L}{\partial w_m}$$

$$= -\eta \, \frac{\partial L}{\partial so} \frac{\partial so}{\partial n} \frac{\partial n}{\partial si} \frac{\partial si}{\partial m} \frac{\partial m}{\partial w_m}$$



sigmoid activation

softmax activation

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
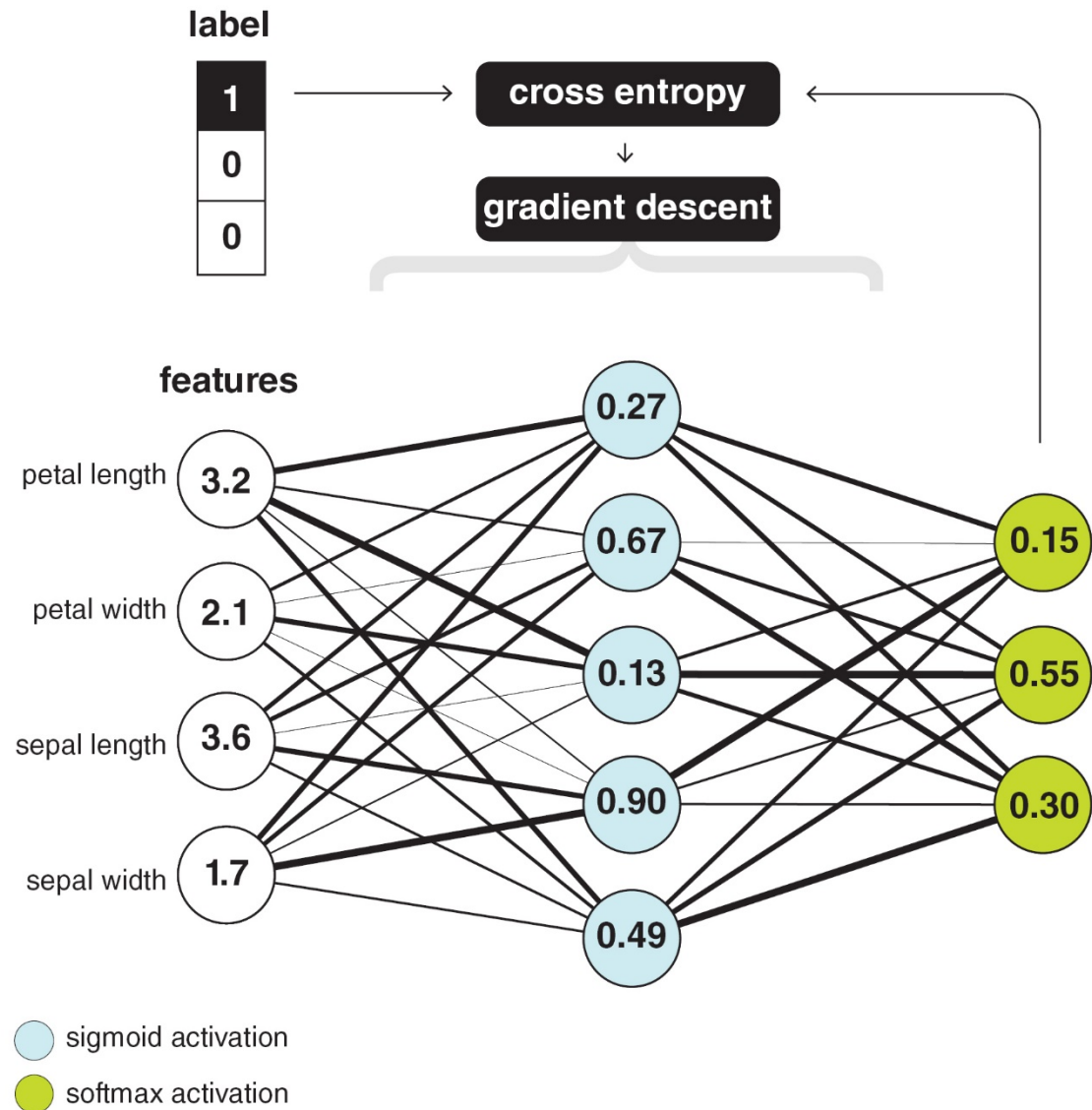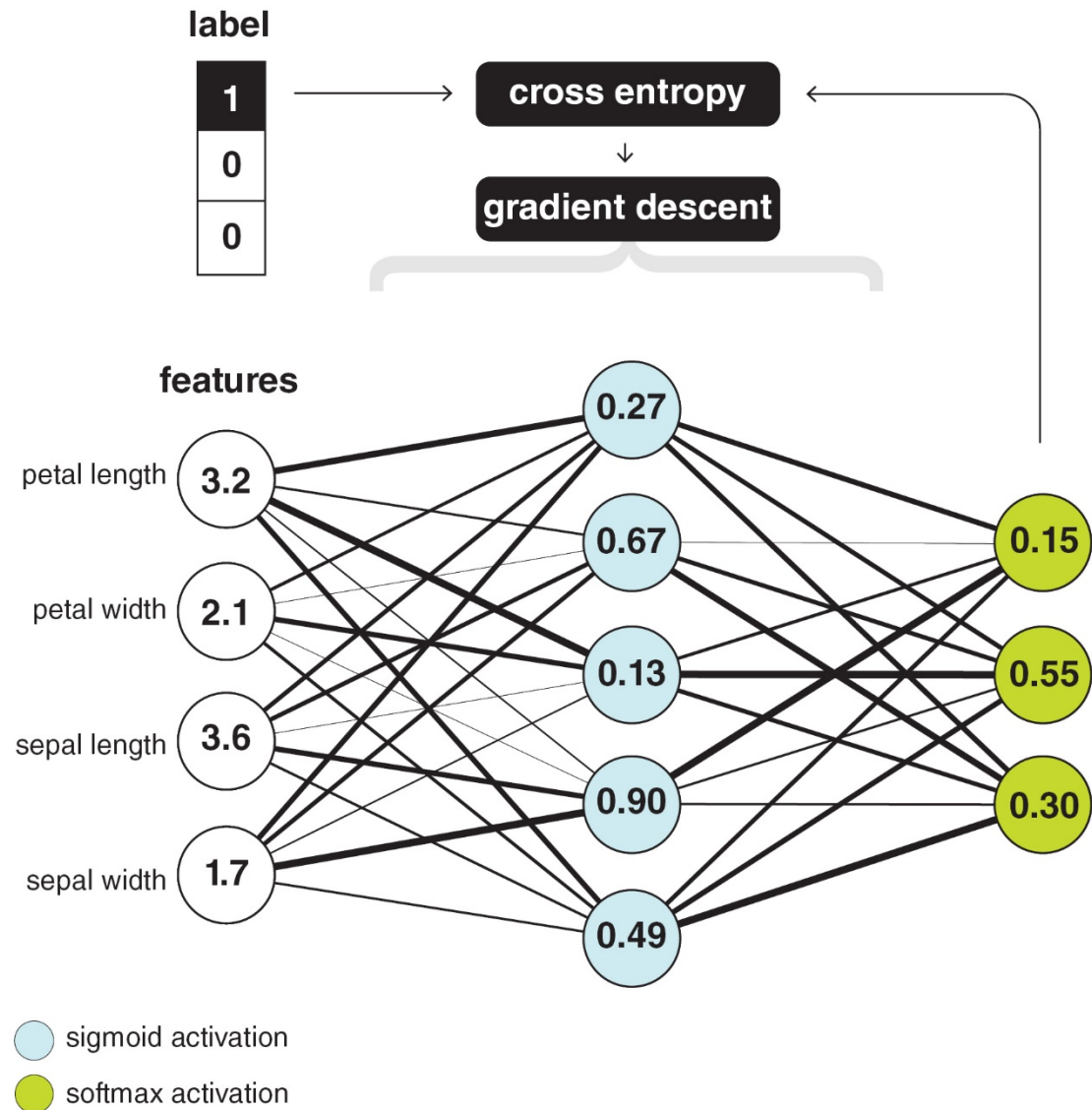4. loss calculation
5. backpropagation + updates

# How Neural Networks Learn

✔ **Data:** iris dataset
✔ **Model:** 3-layer neural network
✔ **Loss:** cross entropy
✔ **Optimizer:** gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates
6. repeat 2,3,4, & 5

# Gradient Descent Flavors

**vanilla gradient descent** - entire dataset
**stochastic gradient descent** - random batch of samples (IID)
**online gradient descent** - (need not be IID)

# Gradient Descent Flavors

**vanilla gradient descent** - entire dataset

**stochastic gradient descent** - random batch of samples (IID)

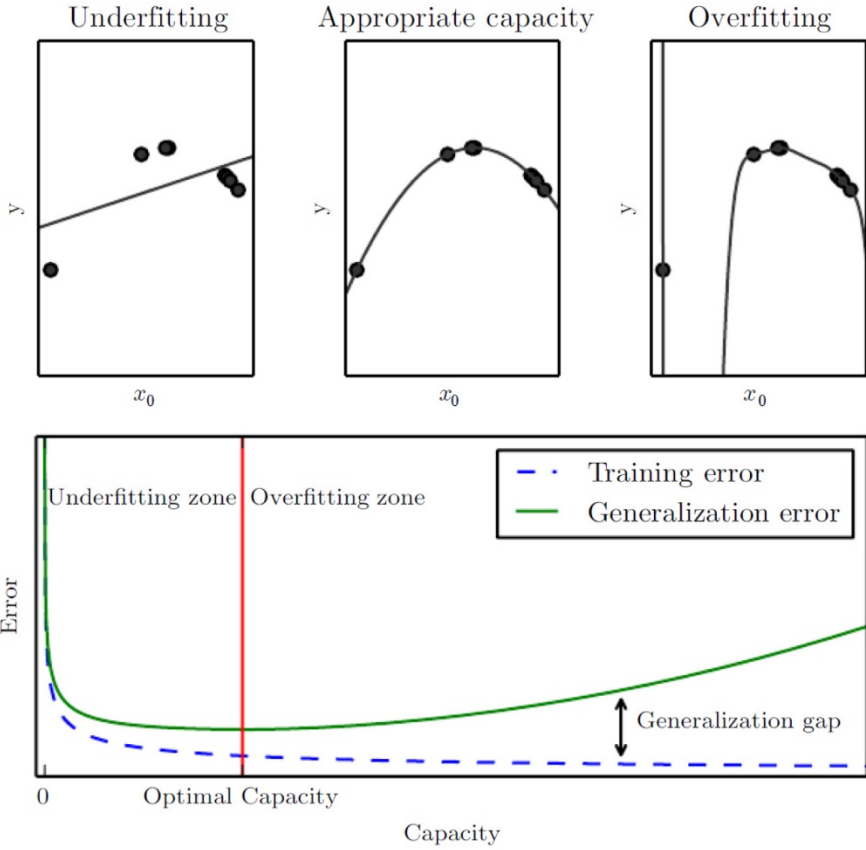**online gradient descent** - (need not be IID)

**learning rate**   **batch size**   **# of epochs**

What is the intuition behind neural networks?

How do neural networks learn?

How to train neural networks?

# The Perfect Fit

# The Perfect Fit



parameters **vs** hyperparameters

*Ian Goodfellow, Yoshua Bengio & Aaron Courville*

Deep Learning
**MIT Press - 2016**

# Hyperparameters

**model-specific** **vs** **optimizer-specifc**

| architecture | learning rate |
| activations | batch size |
| initializations | # of epochs |
| loss functions | |
| optimizers | |
| regularizers | |

# Architecture



# of layers
# of units/layer

*Christian Szegedy, Wei Liu, Yangqing Jia, et al.*

Going Deeper with Convolutions (GoogleNet/Inception)
**CVPR - 2015**

# Activations



**step**

**✗** non-differentiable

# Activations

### step

✘ non-differentiable

### sigmoid

✔ smooth + step-like
✔ good activations close to 0
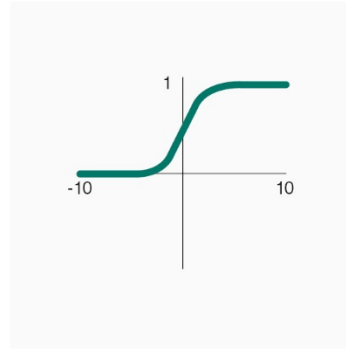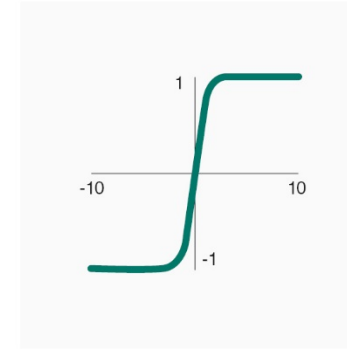✔ activations are bound 0~1
✘ vanishing gradients

# Activations

**step**

✘ non-differentiable

**sigmoid**

✔ smooth + step-like
✔ good activations close to 0
✔ activations are bound 0~1
✘ vanishing gradients

**tanh**

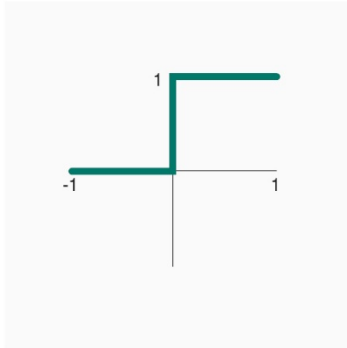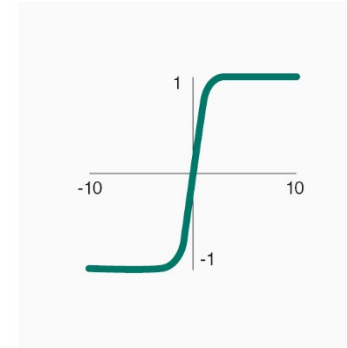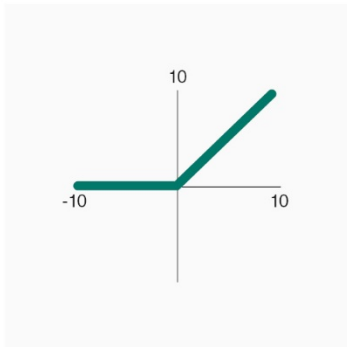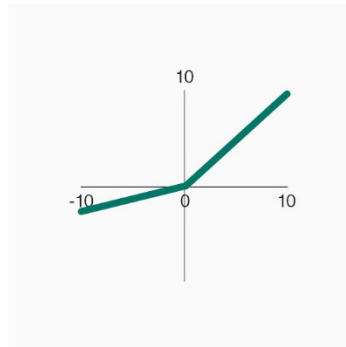✔ scaled sigmoid
✔ stronger activations
✘ vanishing gradients

*Yann LeCun, Leon Bottou, Genevieve B. Orr & Klaus -Robert Müller*

Efficient BackProp
**Neural Networks: Tricks of the Trade - 1998**

# Activations

**step**

✘ non-differentiable

**sigmoid**

✔ smooth + step-like
✔ good activations close to 0
✔ activations are bound 0~1
✘ vanishing gradients

**tanh**

✔ scaled sigmoid
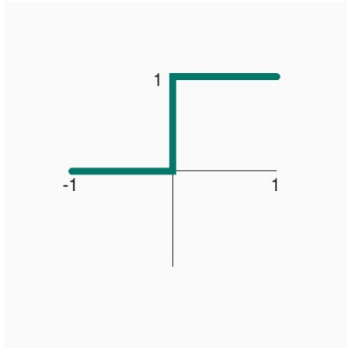✔ stronger activations
✘ vanishing gradients

**ReLU**

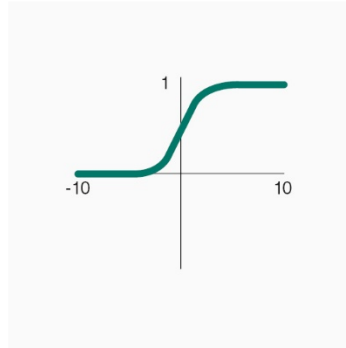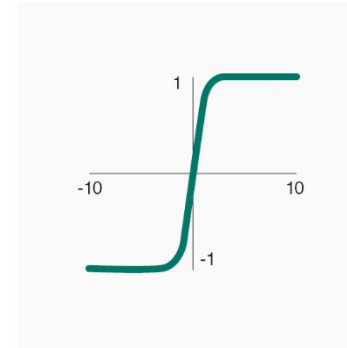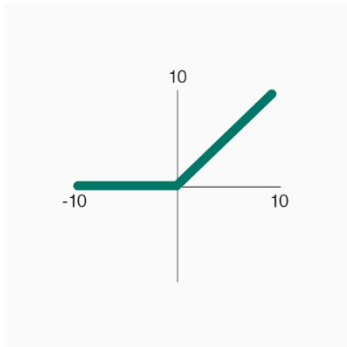✔ sparse activations
✔ efficient
✘ dead nodes
✘ not bound

*Alex Krizhevsky, Ilya Sutskever & Geoffrey E. Hinton*

ImageNet Classification with Deep Convolutional Neural Networks
**Advances in Neural Information Processing - NIPS 2012**

# Activations

**step**
- ✘ non-differentiable

**sigmoid**
- ✔ smooth + step-like
- ✔ good activations close to 0
- ✔ activations are bound 0~1
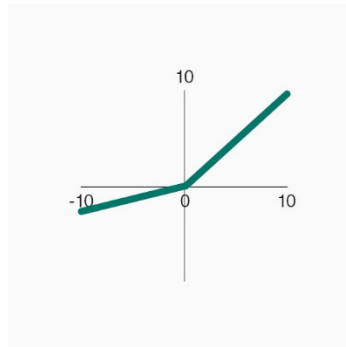- ✘ vanishing gradients

**tanh**
- ✔ scaled sigmoid
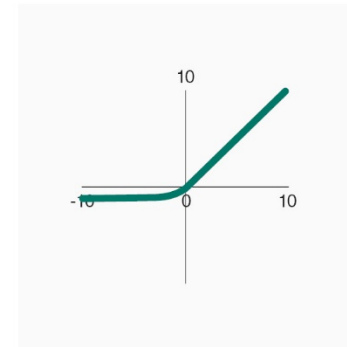- ✔ stronger activations
- ✘ vanishing gradients

**ReLU**
- ✔ sparse activations
- ✔ efficient
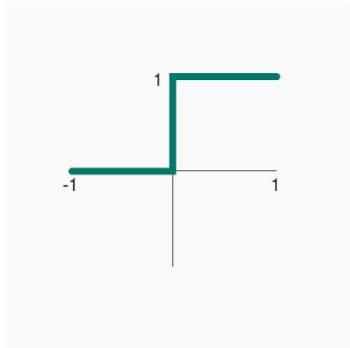- ✘ dead nodes
- ✘ not bound
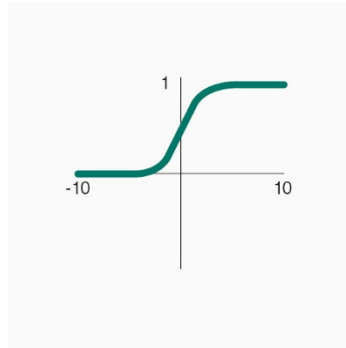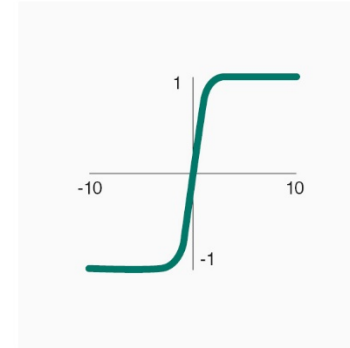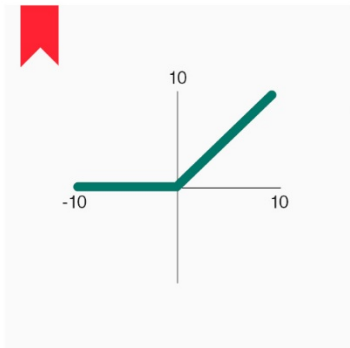
**Leaky ReLU**
- ✔ no dead nodes

# Activations

**step**



✘ non-differentiable

**sigmoid**



✔ smooth + step-like
✔ good activations close to 0
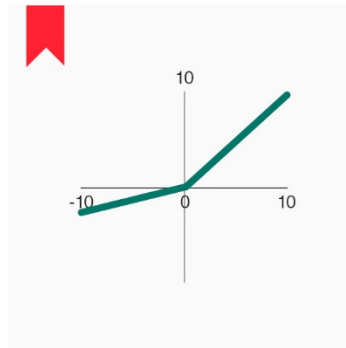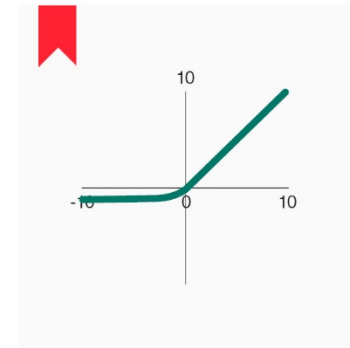✔ activations are bound 0~1
✘ vanishing gradients

**tanh**



✔ scaled sigmoid
✔ stronger activations
✘ vanishing gradients

**ReLU**



✔ sparse activations
✔ efficient
✘ dead nodes
✘ not bound

**Leaky ReLU**



✔ no dead nodes

**ELU**



✔ robust to noise
✘ expensive

# Activations

### step

✘ non-differentiable

### sigmoid

✔ smooth + step-like
✔ good activations close to 0
✔ activations are bound 0~1
✘ vanishing gradients

### tanh

✔ scaled sigmoid
✔ stronger activations
✘ vanishing gradients

### ReLU

✔ sparse activations
✔ efficient
✘ dead nodes
✘ not bound

### Leaky ReLU

✔ no dead nodes

### ELU

✔ robust to noise
✘ expensive

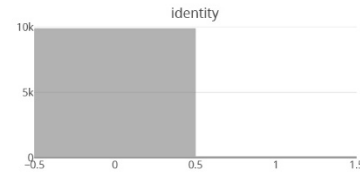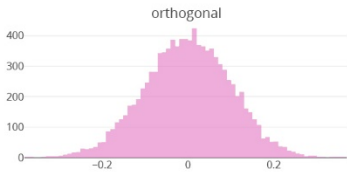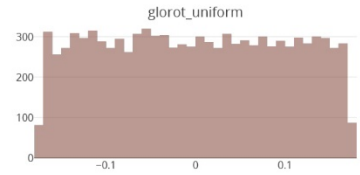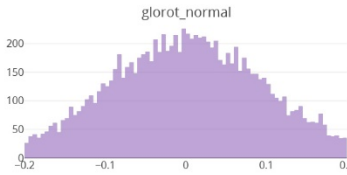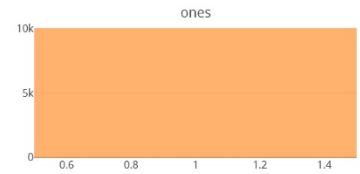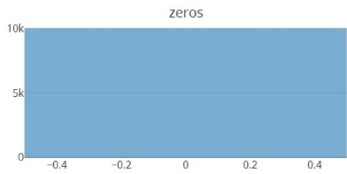# Initializations

**0 -** stuck at a saddle point

**constants -** difficult to break the symmetry

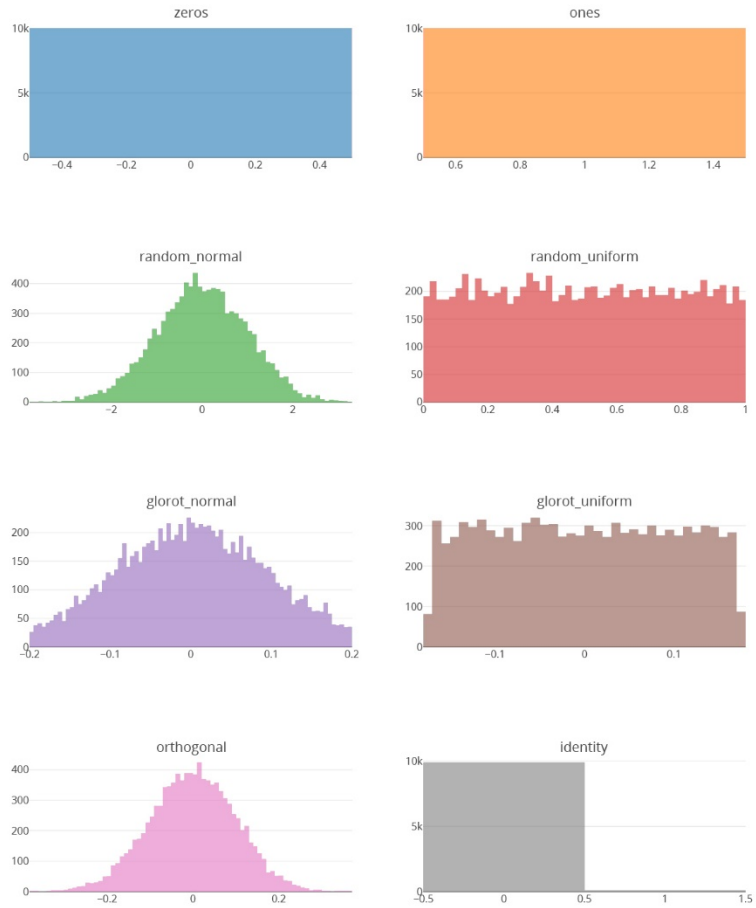**large random values -** small gradients, slow convergence

# Initializations



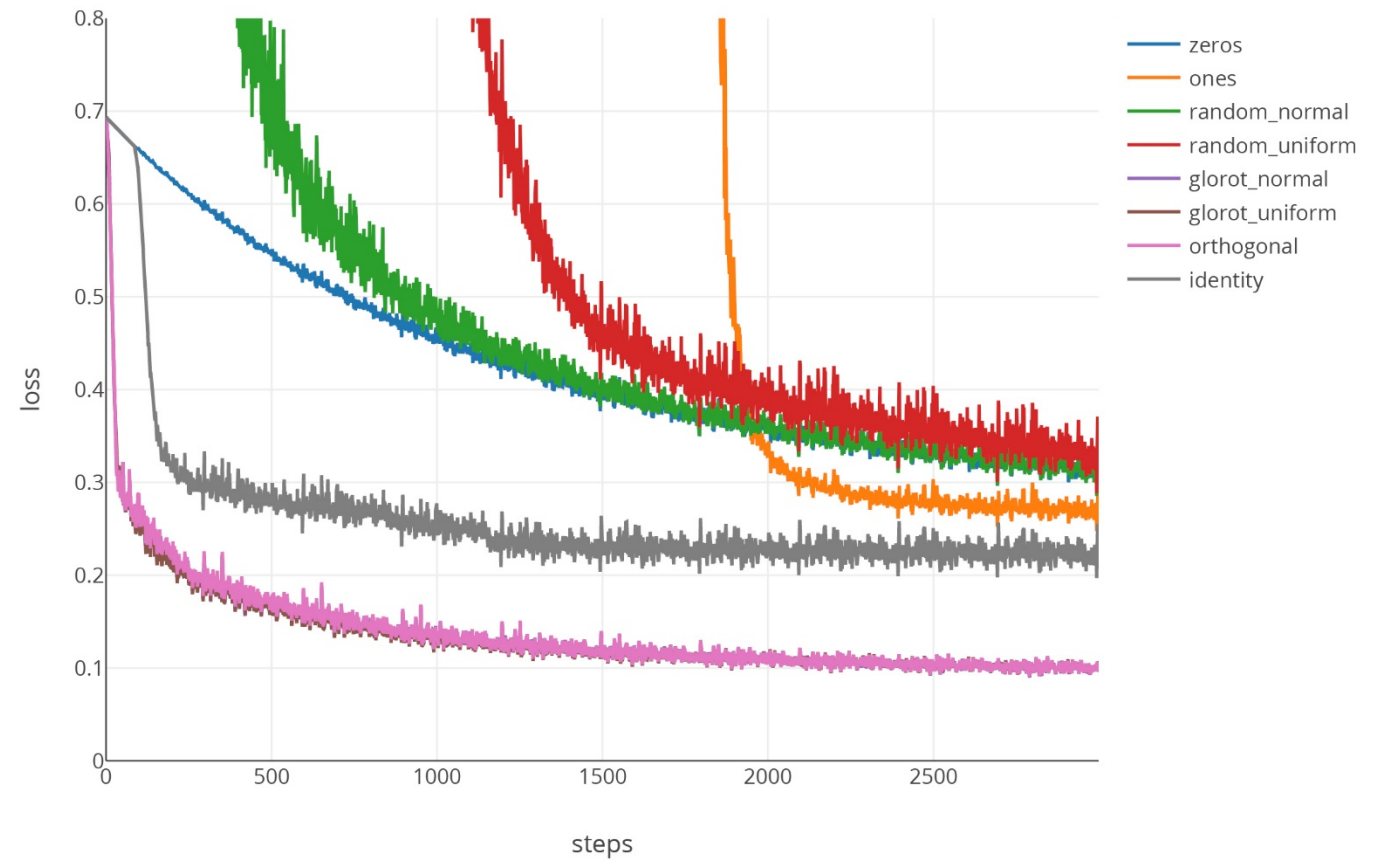Tensorflow initializer distribution - 10k samples

# Initializations



Tensorflow initializer distribution - 10k samples

Step loss with different weight initialization

# Initializations

| Name | $\alpha$ | $\beta$ | $\gamma$ | Reference |
|------|----------|---------|----------|-----------|
| Constant | $\alpha = 0$ | $\beta = 0$ | $\gamma \geq 0$ | used by [ZF14] |
| → Xavier/Glorot uniform | $\alpha = \sqrt{\frac{6}{n_{in}+n_{out}}}$ | $\beta = 0$ | $\gamma = 0$ | [GB10] |
| → Xavier/Glorot normal | $\alpha = 0$ | $\beta = \left(\frac{2}{(n_{in}+n_{out})}\right)^2$ | $\gamma = 0$ | [GB10] |
| → He | $\alpha = 0$ | $\beta = \frac{2}{n_{in}}$ | $\gamma = 0$ | [HZRS15b] |
| Orthogonal | — | — | $\gamma = 0$ | [SMG13] |
| LSUV | — | — | $\gamma = 0$ | [MM15] |

Table B.2.: Weight initialization schemes of the form $w \sim \alpha \cdot \mathcal{U}[-1, 1] + \beta \cdot \mathcal{N}(0, 1) + \gamma$.
$n_{in}, n_{out}$ are the number of units in the previous layer and the next layer. Typically, biases are initialized with constant 0 and weights by one of the other schemes to prevent unit-coadaptation. However, dropout makes it possible to use constant initialization for all parameters.
LSUV and Orthogonal initialization cannot be described with this simple pattern.

Analysis and Optimization of Convolutional Neural Network Architectures
https://arxiv.org/pdf/1707.09725.pdf

# Loss Functions

**regression -** mean squared error

**multiclass classification -** categorical cross entropy

**pixel classification -** dice/ Wasserstein dice coefficient

# Optimizers

**stochastic gradient descent + momentum**

# Optimizers

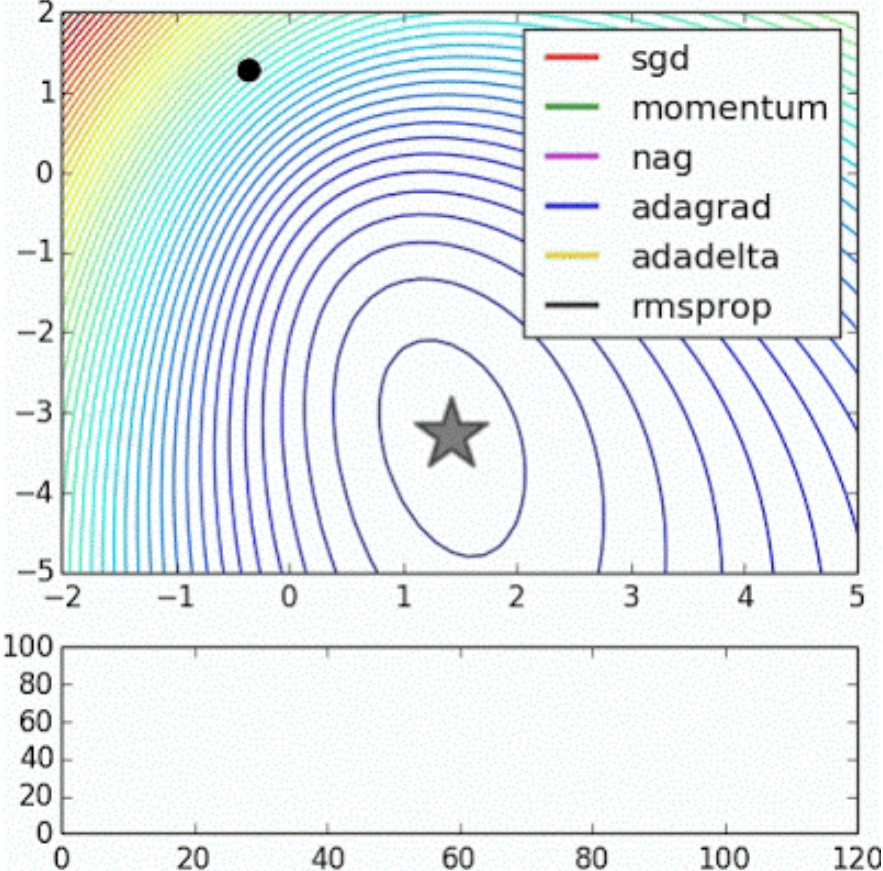**stochastic gradient descent + momentum**

**adaptive gradient (AdaGrad)**

# Optimizers

**stochastic gradient descent + momentum**

**adaptive gradient (AdaGrad)**

**root mean square propagation (RMSProp)**

Coursera: Neural Networks for Machine Learning - Lecture 6

http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

# Optimizers

# Regularizers

## L1, L2 regularization

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2}|W|$$

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2}W^2$$

# Regularizers

## L1, L2 regularization

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2}|W|$$

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2}W^2$$

# Regularizers

## dropout



(a) Standard Neural Net      (b) After applying dropout.

Figure 1: Dropout Neural Net Model. **Left**: A standard neural net with 2 hidden layers. **Right**: An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

# Regularizers

## batch normalization

# Regularizers

## batch normalization

# Optimizer-specifc Hyperparameters

**learning rate**

0.1, 0.01, 0.001, 0.0001,...

# Optimizer-specifc Hyperparameters

**learning rate**

0.1, 0.01, 0.001, 0.0001,...

# Optimizer-specifc Hyperparameters

**learning rate**

0.1, 0.01, 0.001, 0.0001,...

# Optimizer-specifc Hyperparameters

## learning rate
0.1, 0.01, 0.001, 0.0001,...

# Optimizer-specifc Hyperparameters

**learning rate**

0.1, 0.01, 0.001, 0.0001,...



**batch size**

16, 32, 64, 128,...

# Optimizer-specifc Hyperparameters

**learning rate**

0.1, 0.01, 0.001, 0.0001,...



**batch size**

16, 32, 64, 128,...

**# of epochs**

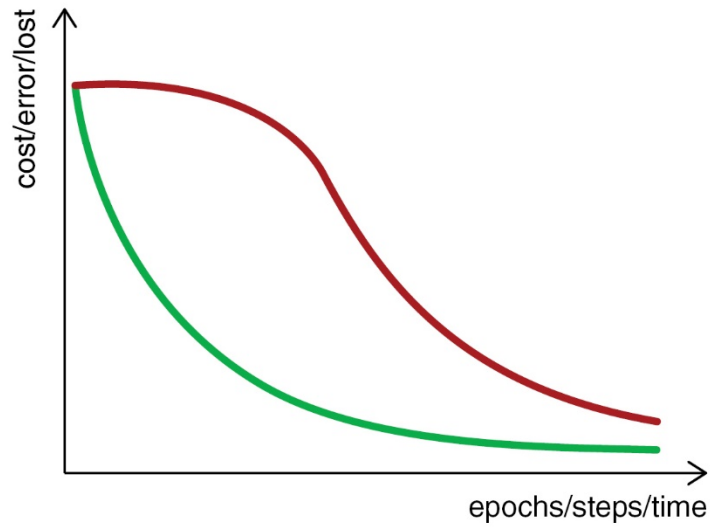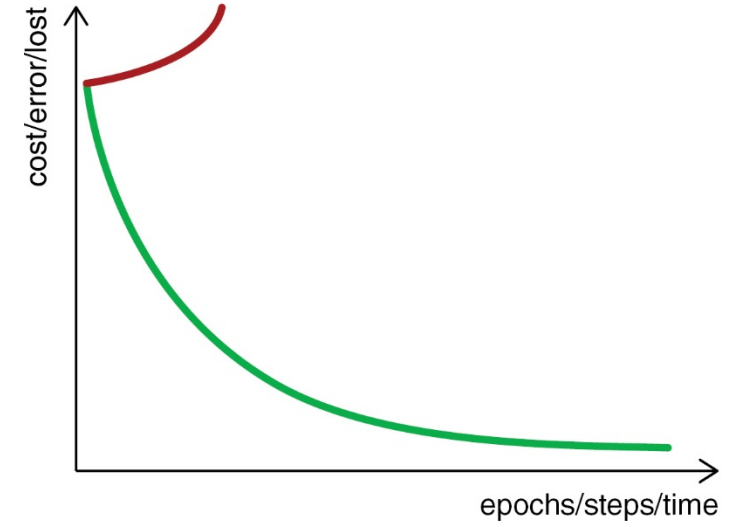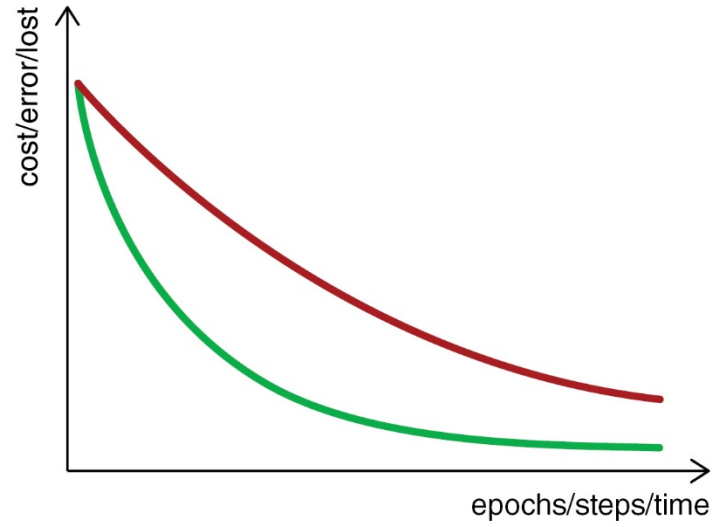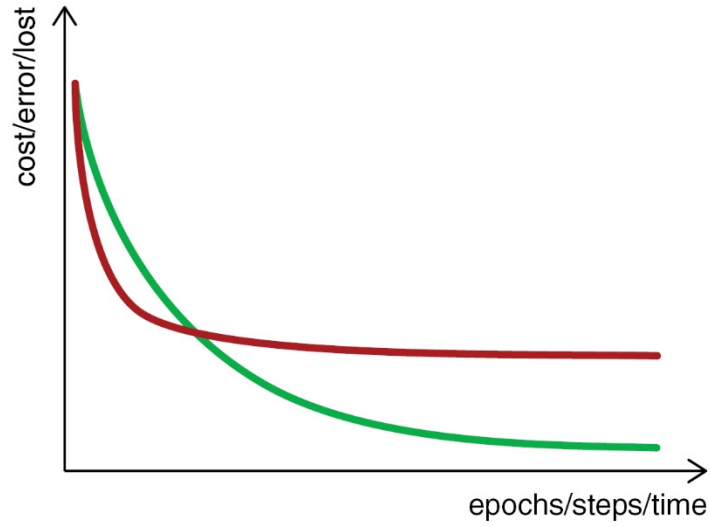early stopping

# Babysitting your Network



https://lossfunctions.tumblr.com/

# Debugging through Learning Curves

# Debugging through Learning Curves
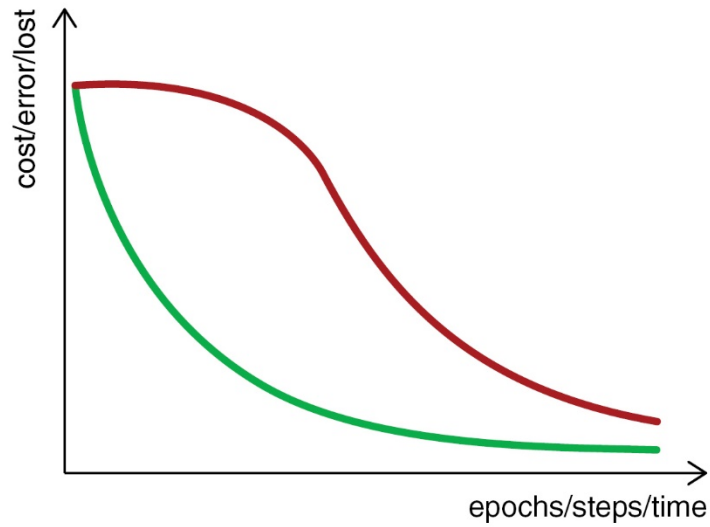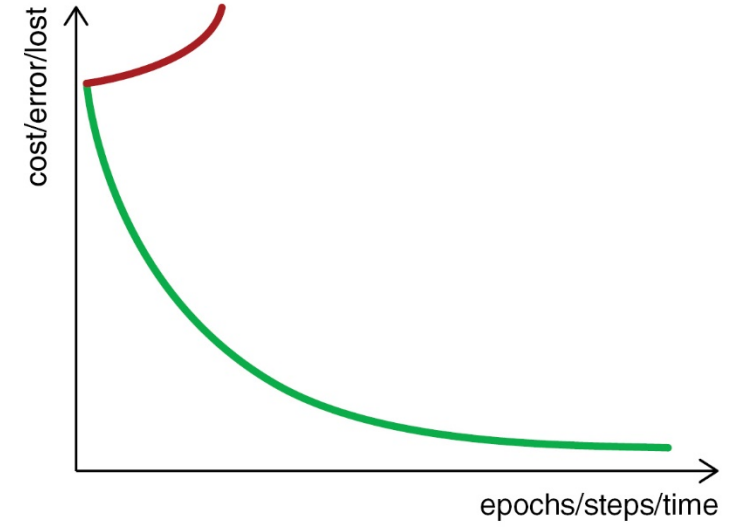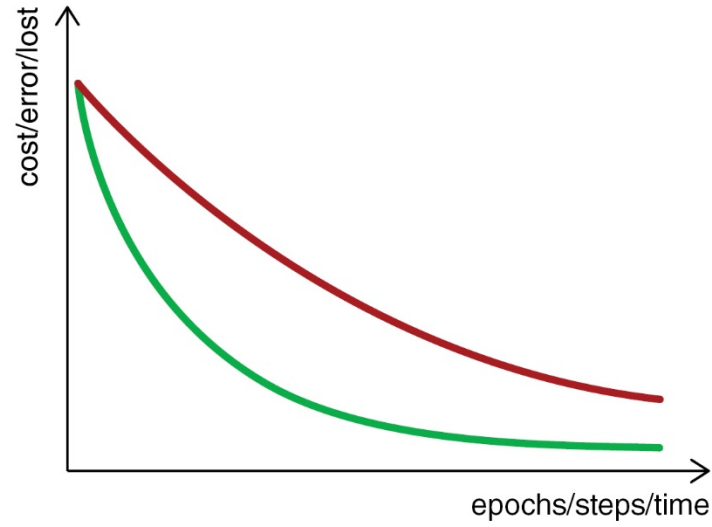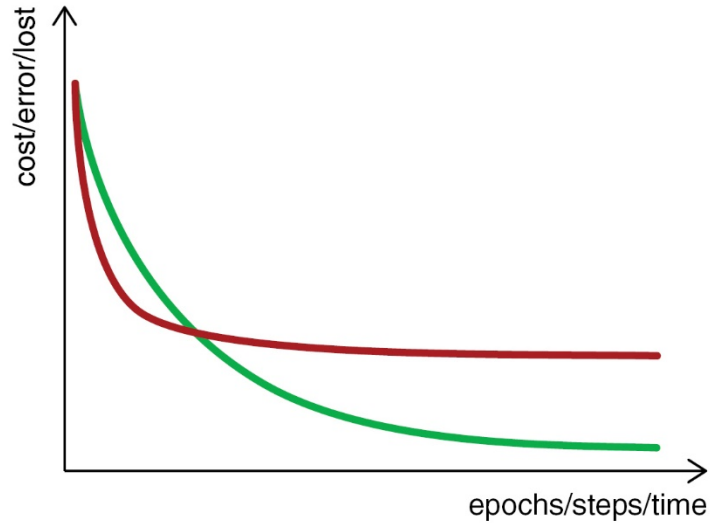
# Debugging through Learning Curves

# Debugging through Learning Curves

# Debugging through Learning Curves

# Debugging through Learning Curves

# Debugging through Learning Curves

*Thank you!*