



FSAN/ELEG815: Statistical Learning

Gonzalo R. Arce

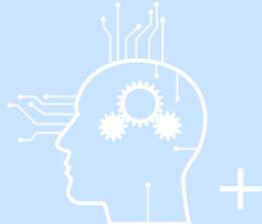
Department of Electrical and Computer Engineering
University of Delaware

X: Neural Networks

What is Deep Learning?

Artificial Intelligence

Any technique that enables computers to mimic human behavior



+

Machine Learning

Ability to learn without explicitly being programmed



Deep Learning

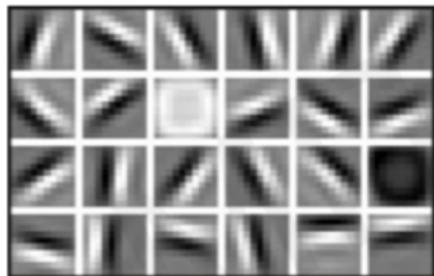
Extract patterns from data using neural networks



Why Deep Learning?

Hand engineered features are time consuming, brittle and not scalable in practice

Can we learn the **underlying features** directly from data?



Line and Edges

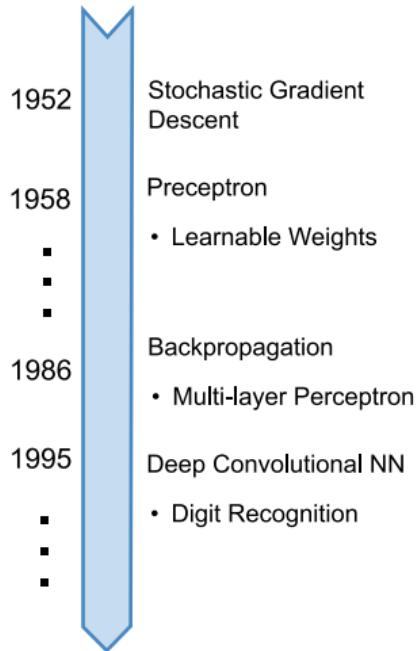


Eyes, Nose and Ears



Facial Structure

Why Now?



Neural Networks date back decades, so why the resurgence?

1. Big Data

- ▶ Large Datasets
- ▶ Easier Collection and Storage



2. Hardware

- ▶ Graphics Processing Units (GPUs)
- ▶ Massively Parallelizable

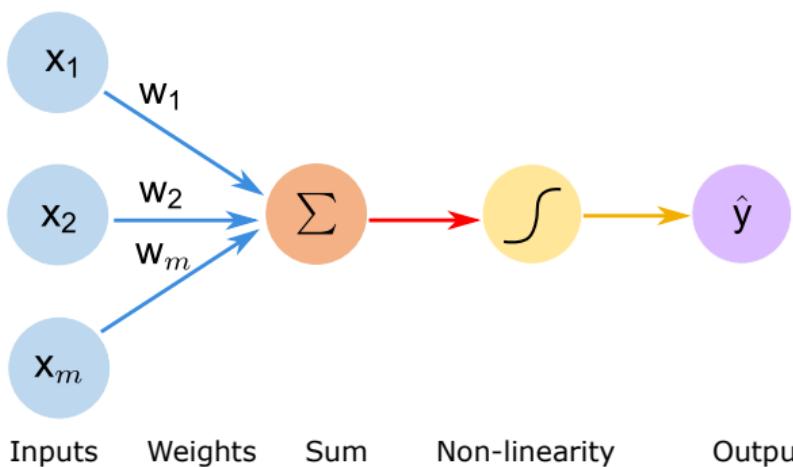


3. Software

- ▶ Improved Techniques
- ▶ New Models
- ▶ Toolboxes



The Perceptron: Forward Propagation



Output

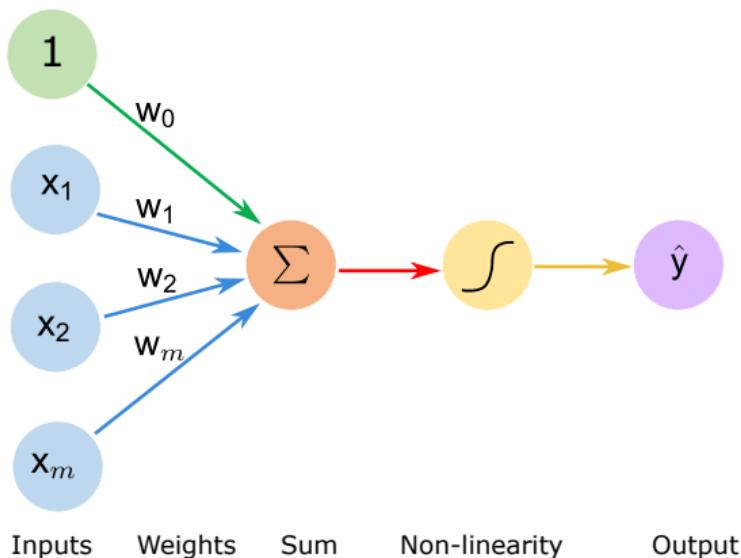
Linear combination of inputs

$\hat{y} = g \left(\sum_{i=1}^m w_i x_i \right)$

Non-linear activation function

The diagram shows the mathematical formula for the perceptron's output. The output is labeled \hat{y} . Above it, the linear combination of inputs is shown as $\sum_{i=1}^m w_i x_i$. Below the linear combination, the non-linear activation function g is indicated by a yellow arrow pointing upwards. A purple arrow points downwards from the output to the activation function, and a red arrow points downwards from the linear combination to the activation function.

The Perceptron: Forward Propagation



Output

Bias

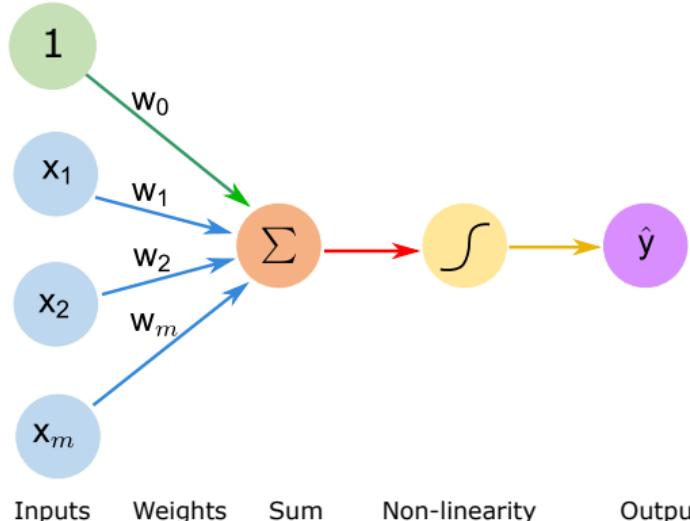
Linear combination of inputs

$\hat{y} = g(w_0 + \sum_{i=1}^m w_i x_i)$

Non-linear activation function

The diagram shows the mathematical formula for the perceptron's output. The output \hat{y} is calculated as the result of applying a non-linear activation function g to the linear combination of the bias w_0 and the weighted sum of the inputs x_i and their corresponding weights w_i .

The Perceptron: Forward Propagation

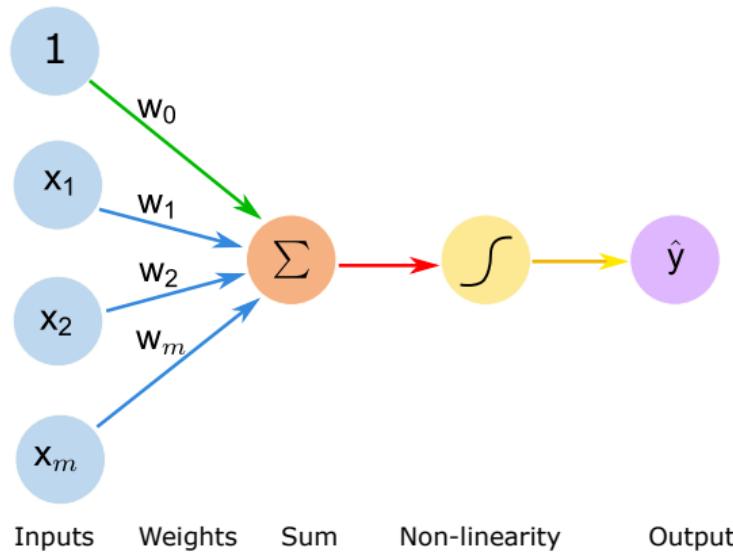


$$\hat{y} = g(w_0 + \sum_{i=1}^m w_i x_i)$$

$$\hat{y} = g(w_0 + \mathbf{X}^\top \mathbf{W})$$

where: $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$

The Perceptron: Forward Propagation

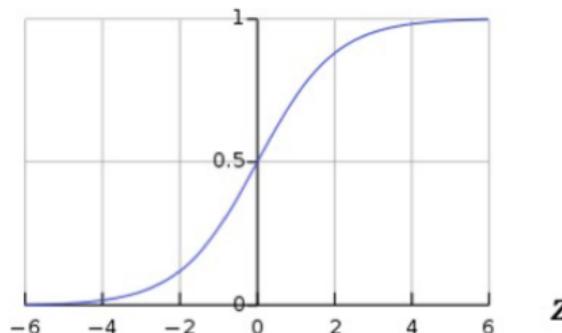


Activation Functions

$$\hat{y} = g(w_0 + \mathbf{x}^\top \mathbf{w})$$

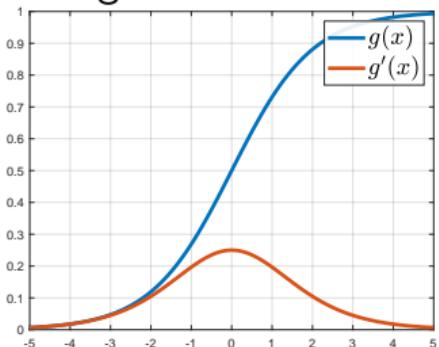
- Example: sigmoid function

$$g(z) = \sigma(z) = \frac{1}{1+e^{-z}}$$



Common Activation Functions

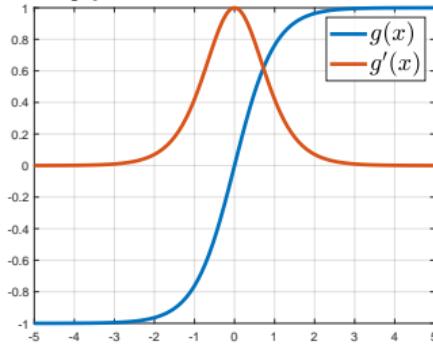
Sigmoid Function



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

$$g'(x) = g(x)(1 - g(x))$$

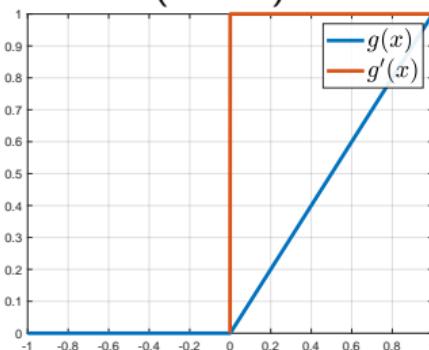
Hyperbolic Function



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

$$g'(x) = 1 - g(x)^2$$

Rectified Linear Unit
(ReLU)

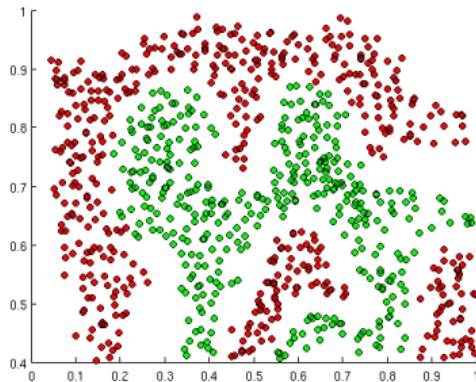


$$g(x) = \max(0, x)$$

$$g'(x) = \begin{cases} 1 & x > 0 \\ 0 & \text{otherwise} \end{cases}$$

The Importance of Activation Functions

The purpose of the activation function is to **introduce non-linearities** into the network

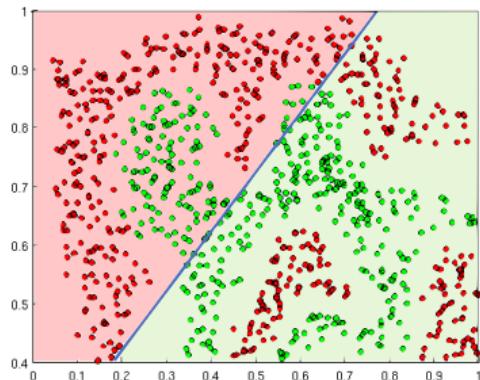


What if we want to build a neural network to distinguish green vs red points?

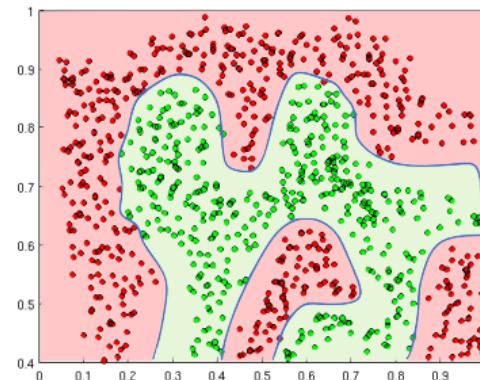
The Importance of Activation Functions

The purpose of the activation function is to **introduce non-linearities** into

the network

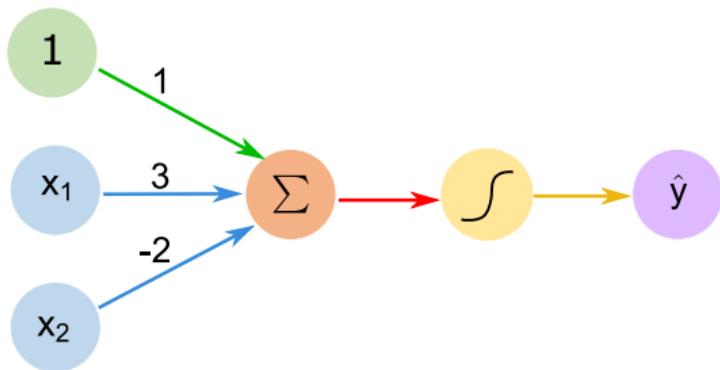


Linear activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions

The Perceptron: Example



We have $w_0 = 1$ and $\mathbf{W} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$
 $\hat{y} = g(w_0 + \mathbf{X}^\top \mathbf{W})$

$$\hat{y} = g\left(w_0 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right)$$

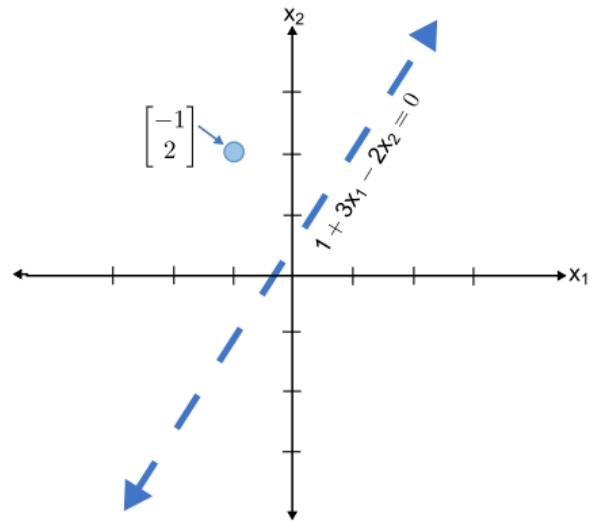
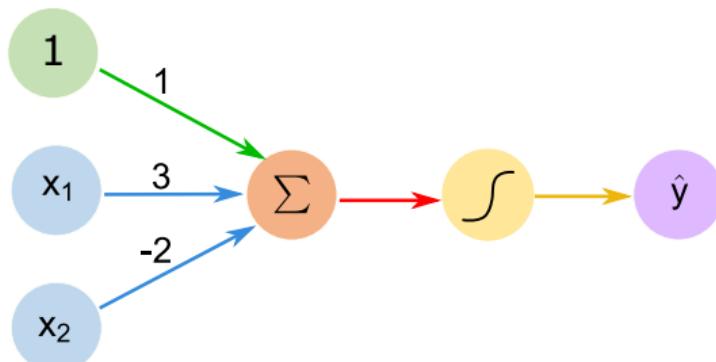
$$\hat{y} = g(\underbrace{1 + 3x_1 - 2x_2}_{})$$

This is just a line in 2D

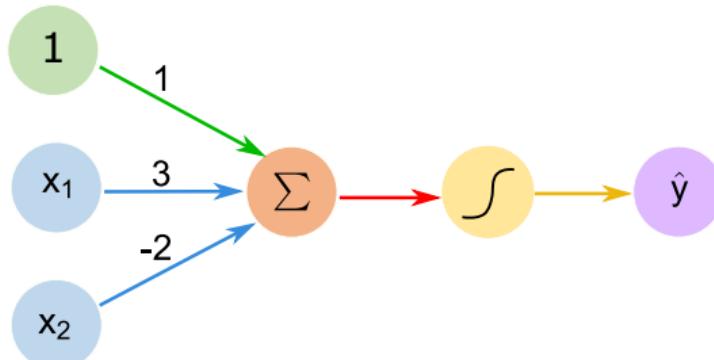
The Perceptron: Example

Assume $\mathbf{X} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$:

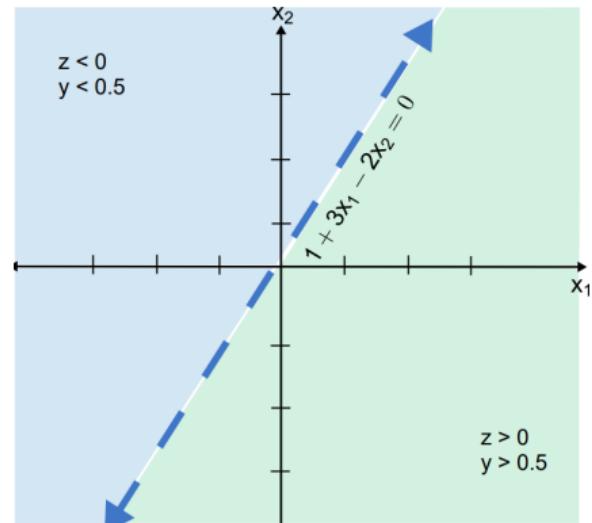
$$\hat{y} = g(1 + (3 \times -1) - (2 \times 2))$$
$$\hat{y} = g(-6) \approx 0.0002$$



The Perceptron: Example

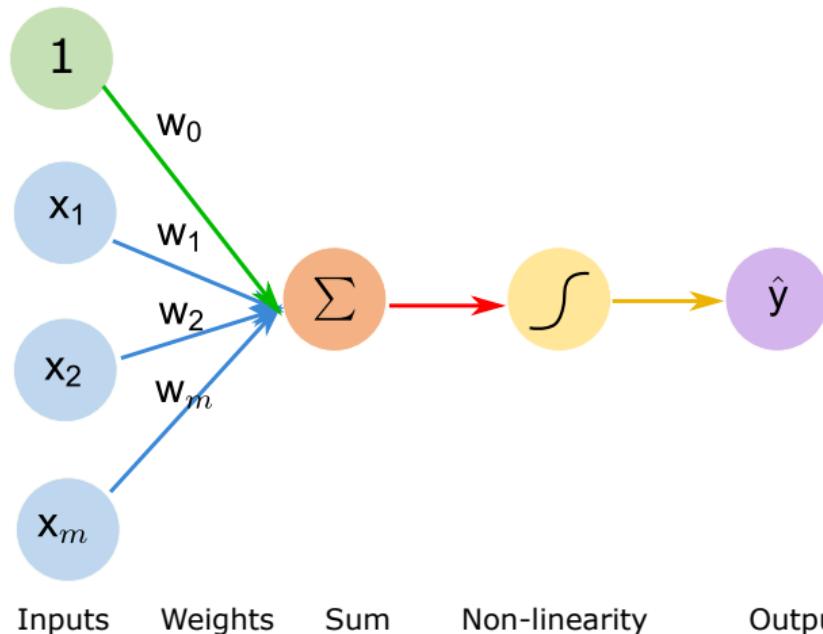


$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

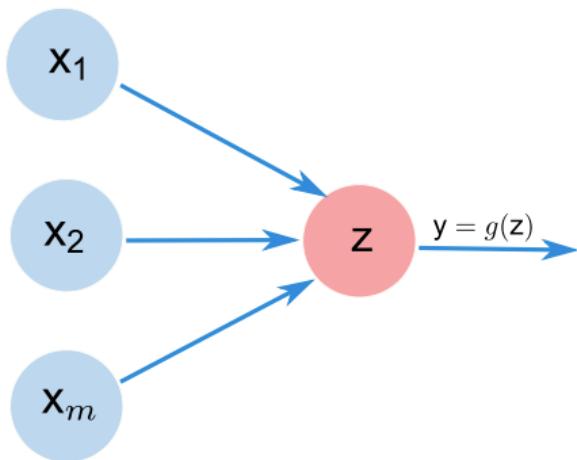


Building Neural Network with Perceptron

$$\hat{y} = g(w_0 + \mathbf{x}^\top \mathbf{w})$$



Building Neural Network with Perceptron



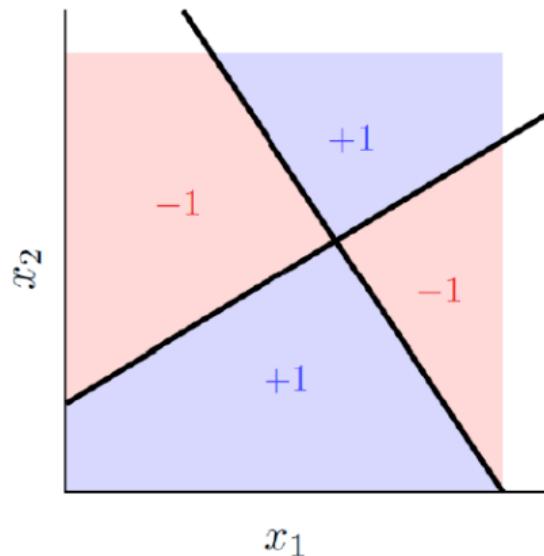
$$z = w_0 + \sum_{j=1}^m x_j w_j$$

Combining Perceptrons

Consider the target function in the figure which is a Boolean XOR function.

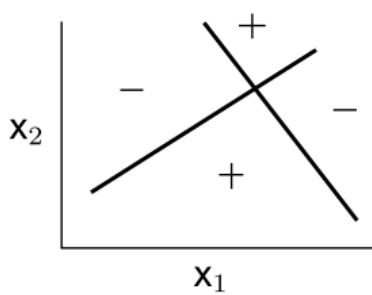
The perceptron cannot implement this classification.

Decompose the f into two perceptrons, corresponding to the lines in the figure.

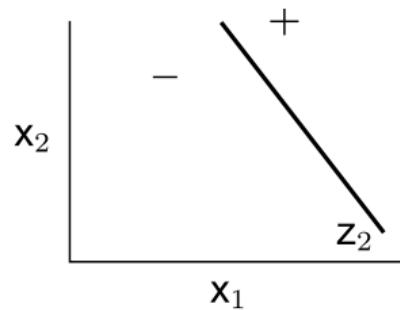
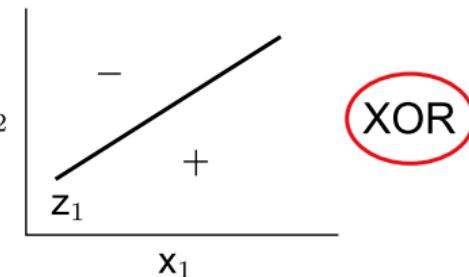


Combining Perceptrons

$$f = \text{XOR}(z_1, z_2)$$



=

 f

$$g_1(\mathbf{x}) = \text{sign}(\mathbf{w}_1^\top \mathbf{x})$$

$$g_2(\mathbf{x}) = \text{sign}(\mathbf{w}_2^\top \mathbf{x})$$

Rewrite f using OR and AND operations:

$$f = g_1 \overline{g_2} + \overline{g_1} g_2$$

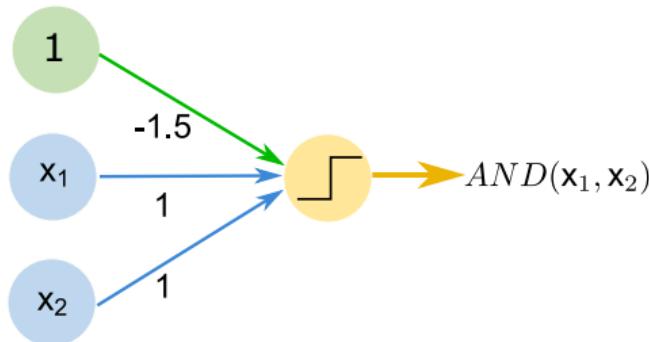
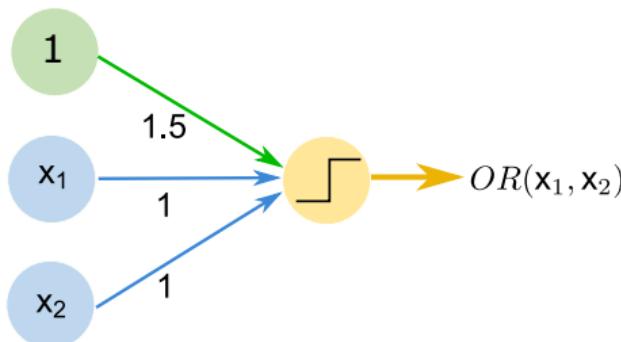
where AND is represented by multiplication, OR by addition and overbar for negation.

Combining Perceptrons

OR and AND can be implemented by the perceptron:

$$\text{OR}(x_1, x_2) = \text{sign}(x_1 + x_2 + 1.5)$$

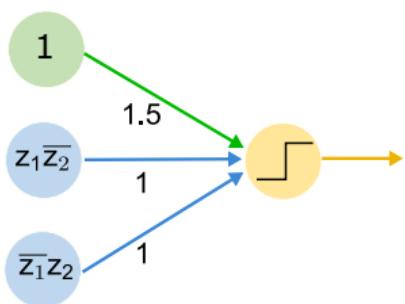
$$\text{AND}(x_1, x_2) = \text{sign}(x_1 + x_2 - 1.5)$$



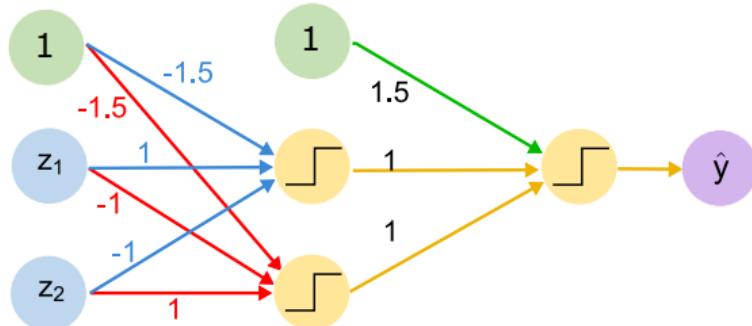
Everything coming to a node is summed and then transformed by $\text{sign}(\cdot)$ to get the final output

Creating Layers

$$f = z_1 \bar{z}_2 + \bar{z}_1 z_1$$



(a)

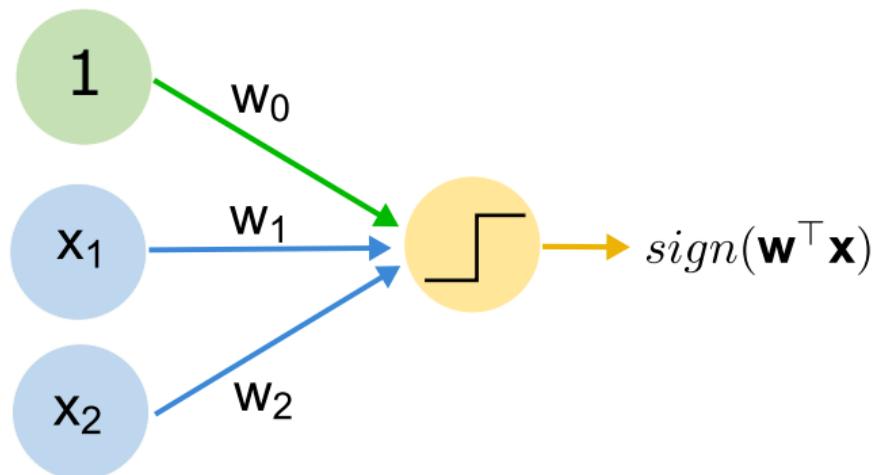


(b)

- (a) OR of the inputs $z_1 \bar{z}_2$ and $\bar{z}_1 z_1$ (b) The blue and red weights simulate the required ANDS ($z_1 \bar{z}_2$ and $\bar{z}_1 z_1$). Negative in the weights handle negations.

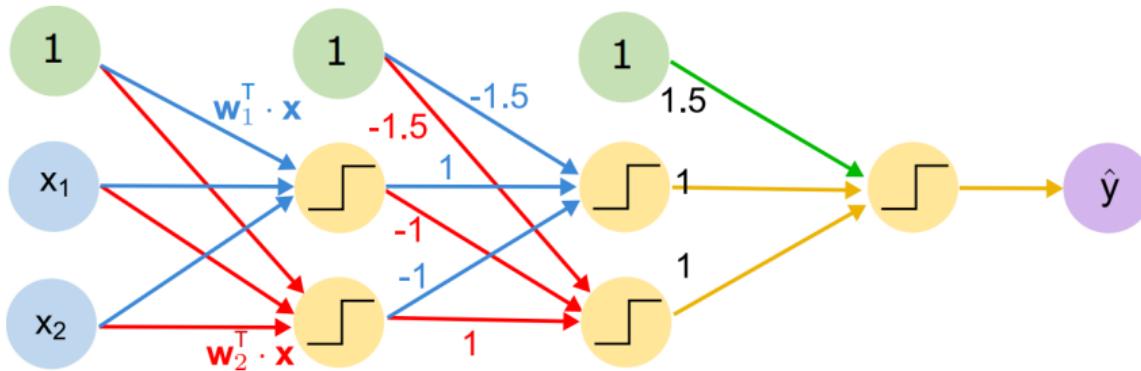
The Multilayer Perceptron

$z_1(\mathbf{x}) = \mathbf{w}_1^\top \mathbf{x}$ and $z_2(\mathbf{x}) = \mathbf{w}_2^\top \mathbf{x}$ are perceptrons:



The Multilayer Perceptron

$z_1(\mathbf{x}) = \mathbf{w}_1^\top \mathbf{x}$ and $z_2(\mathbf{x}) = \mathbf{w}_2^\top \mathbf{x}$ are perceptrons:



3 layers

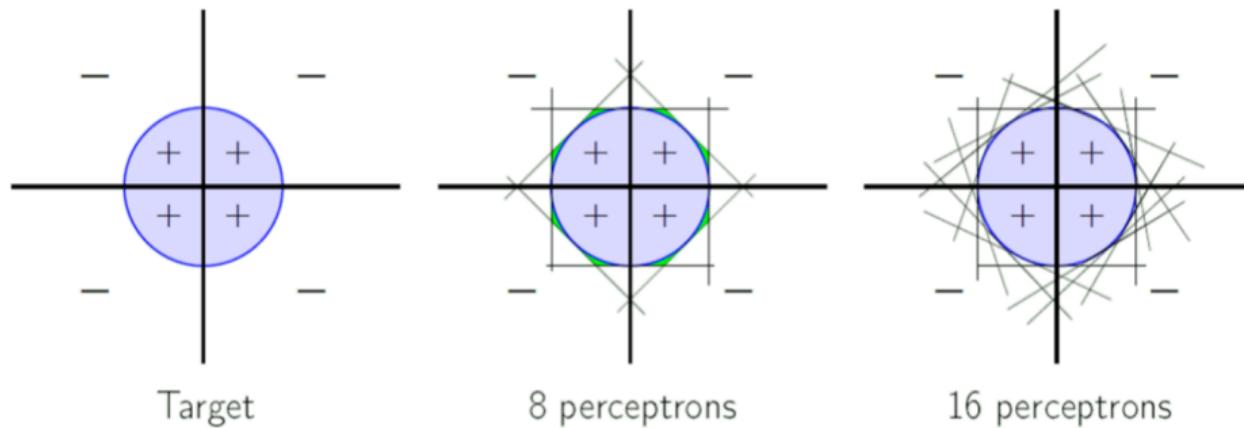
compare to the perceptron
(one)

"Feedforward"

No backward pointing arrows and no
jumps to other layers

A Powerful Model

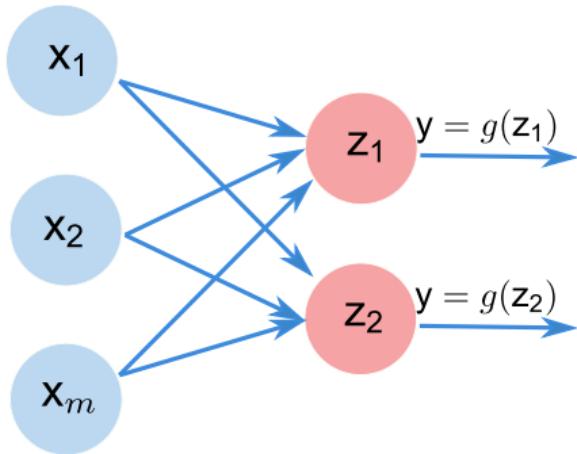
Let's consider now a dis target function:



We can model more complex functions by adding more nodes (hidden units) in the hidden layers (more perceptrons in the decomposition of f)

Multi-output Perceptron

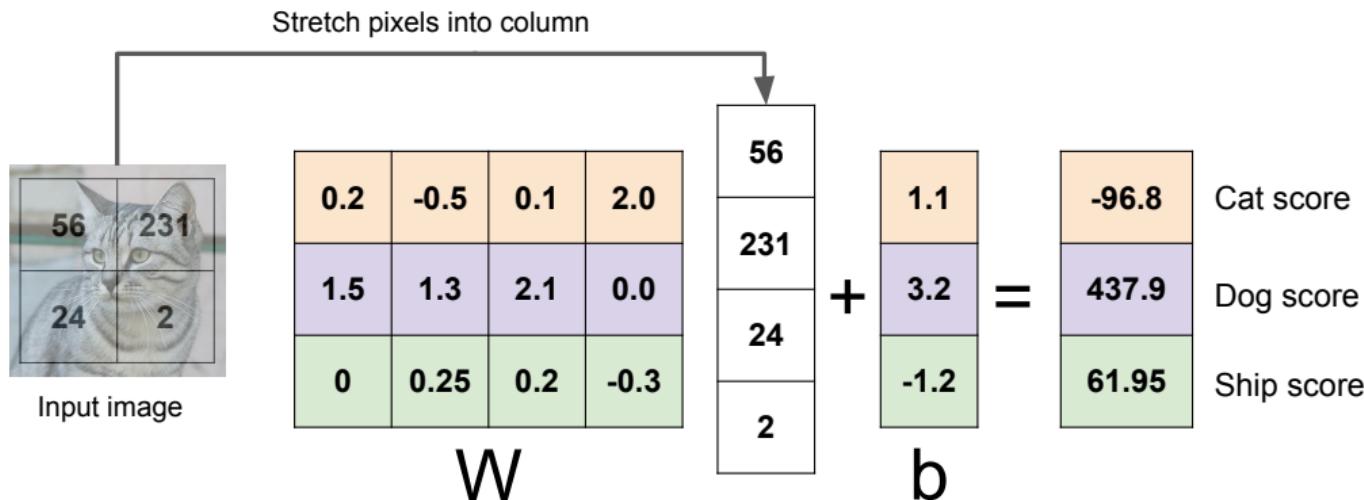
Because all inputs are densely connected to all outputs, these layers are called **Dense** layers:



$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

Multi-output Perceptron

Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)

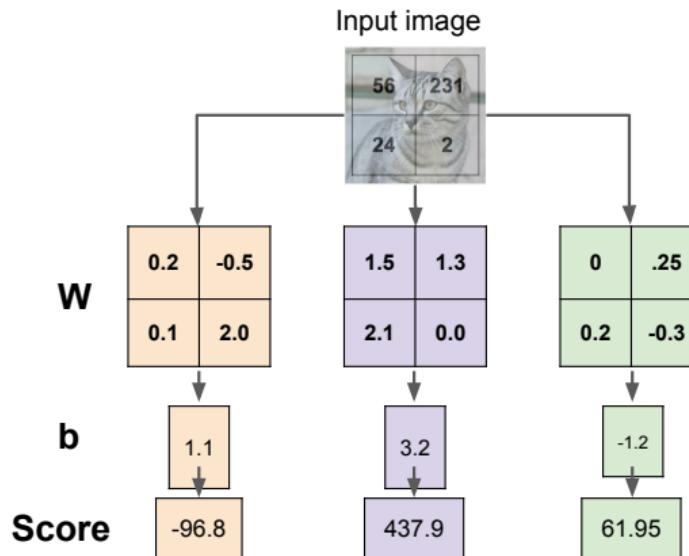
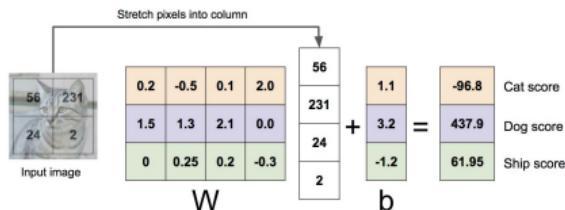


Multi-output Perceptron

Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)

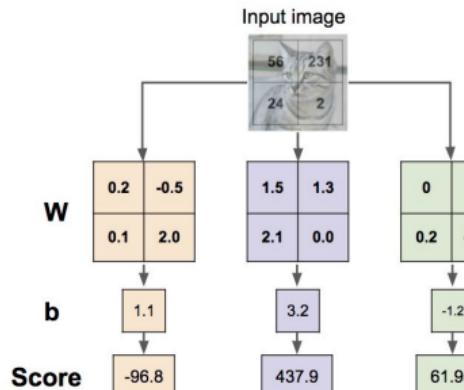
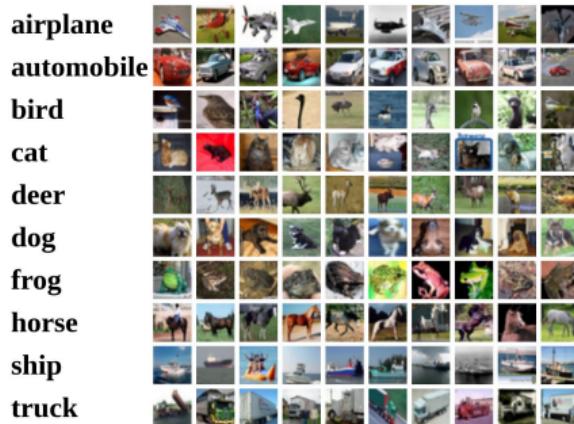
Algebraic Viewpoint

$$f(x, W) = Wx$$



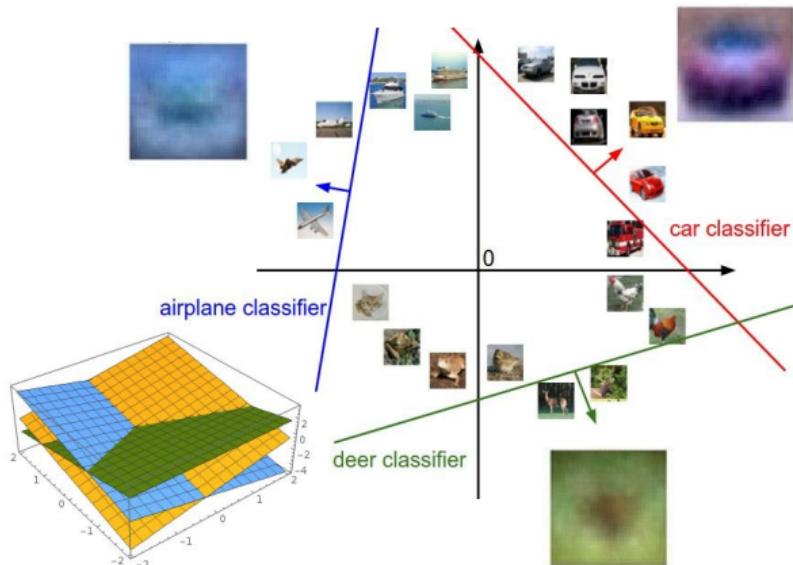
Multi-output Perceptron

Interpreting a Linear Classifier: Visual Viewpoint



Multi-output Perceptron

Interpreting a Linear Classifier: Geometric Viewpoint



$$f(x, W) = Wx + b$$



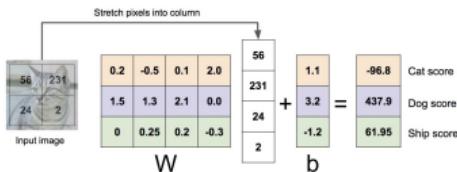
Array of **32x32x3** numbers
(3072 numbers total)

Multi-output Perceptron

Linear Classifier: Three Viewpoints

Algebraic Viewpoint

$$f(x, W) = Wx$$



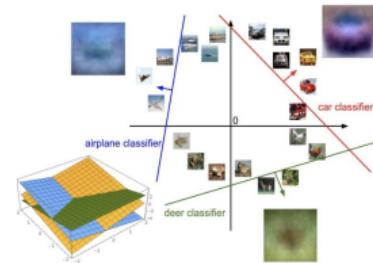
Visual Viewpoint

One template per class

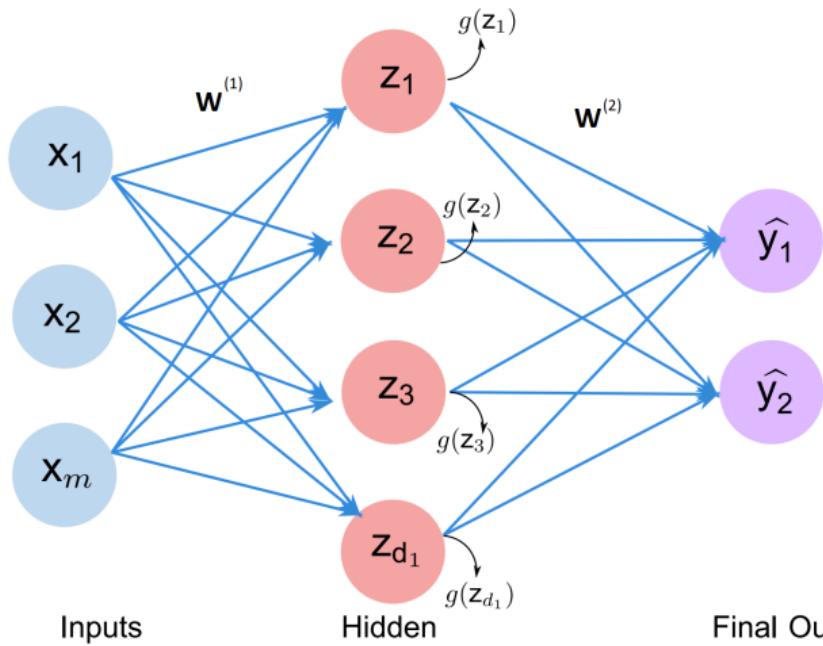


Geometric Viewpoint

Hyperplanes cutting up space



Multi-output Multi-layer perceptron

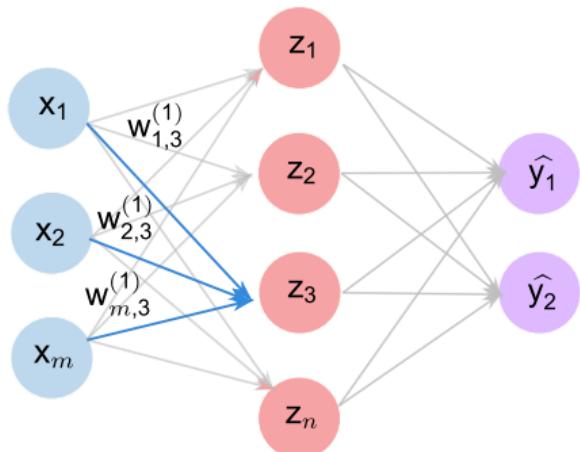


$$\mathbf{z}_i^{(1)} = \mathbf{w}_{0,i}^{(1)} + \sum_{j=1}^m \mathbf{x}_j^{(0)} \mathbf{w}_{j,i}^{(1)}$$

$$\hat{\mathbf{y}}_i = g(\mathbf{w}_{0,i}^{(2)} + \sum_{j=1}^m \mathbf{x}_j^{(1)} \mathbf{w}_{j,i}^{(2)})$$

$$\mathbf{x}_j^{(1)} = g(\mathbf{z}_j^{(1)})$$

Multi-output Multi-layer perceptron



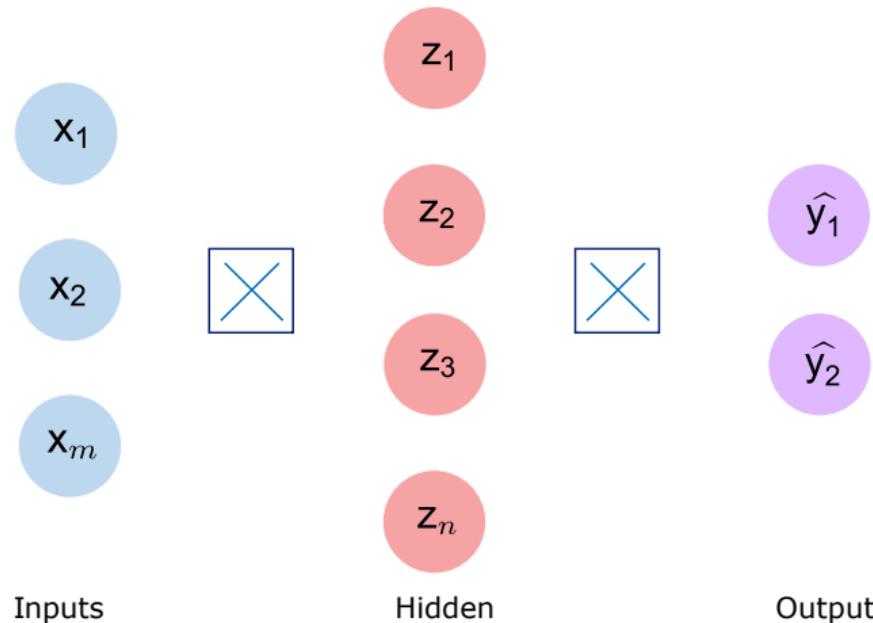
$$z_3 = w_{0,3}^{(1)} + \sum_{j=1}^m x_j w_{j,3}^{(1)}$$

$$z_3 = w_{0,3}^{(1)} + x_1 w_{1,3}^{(1)} + x_2 w_{2,3}^{(1)} + x_3 w_{3,3}^{(1)}$$

Generally

$$\mathbf{z} = \mathbf{W}^{(1)\top} \mathbf{x}^{(0)} + \mathbf{w}_0^{(1)}$$

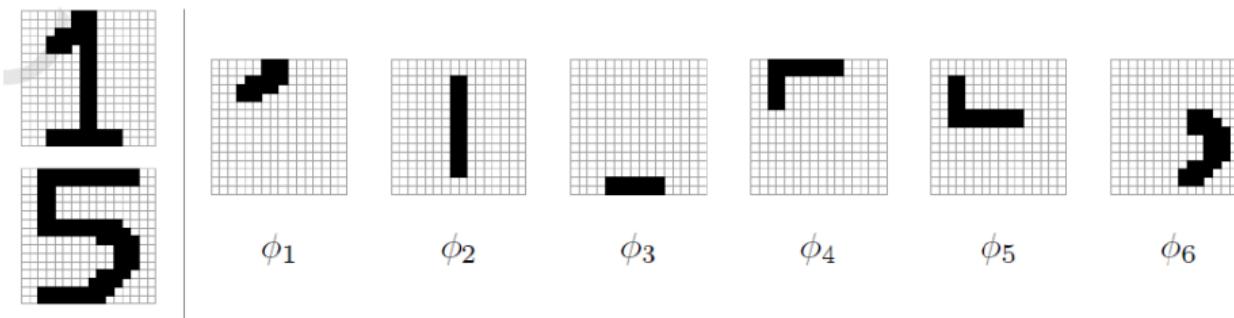
Single Layer Neural Network



Network with many layers - Example

Classify "1" vs "5". Decompose the two digits into basic components:

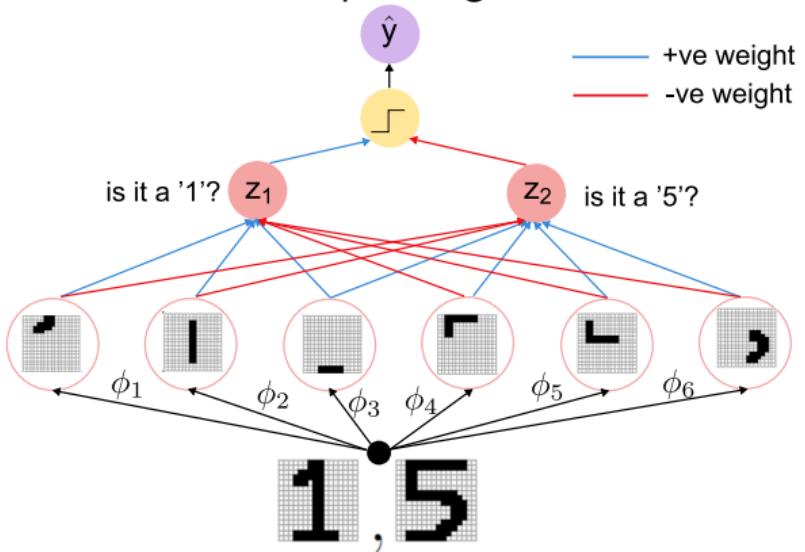
- ▶ Every "1" should contain ϕ_1 , ϕ_2 and ϕ_3
- ▶ Every "5" should contain ϕ_3 , ϕ_4 , ϕ_5 and ϕ_6 , perhaps a little of ϕ_1



These shapes are *features* of the input. We want ϕ_1 to be large (close to 1) if the corresponding feature is in the input image and small (close to -1) if not.

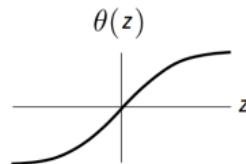
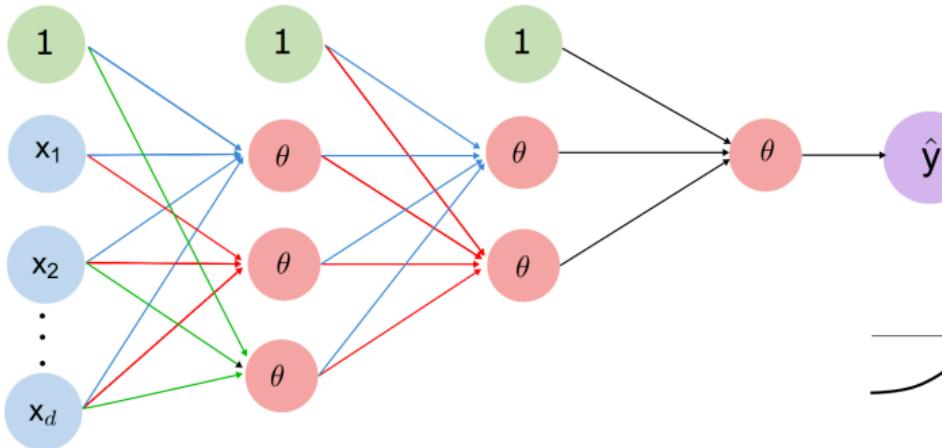
Network with many layers - Example

ϕ_1 is feature function which computes the presence (+1) and absence (-1) of the corresponding feature.



If we feed in "1", ϕ_1 , ϕ_2 and ϕ_3 compute +1 and ϕ_4 , ϕ_5 and ϕ_6 compute -1. Combining with the signs of the weights, z_1 will be positive and z_2 will be negative.

Deep Neural Network

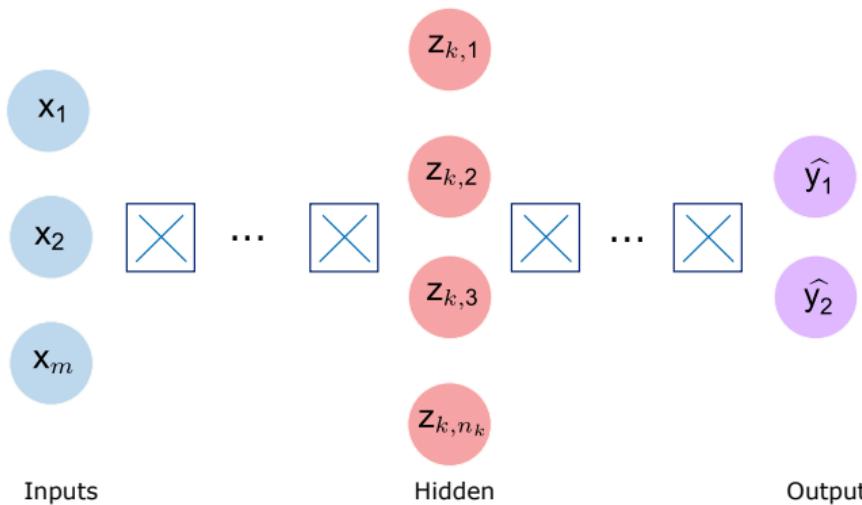


input \mathbf{x} hidden layers $1 \leq l \leq L$ output layer $l = L$

When f is not strictly decomposable into perceptrons, but the decision boundary is smooth (θ), then a multilayer perceptron can be close to implementing f .

Deep Neural Network

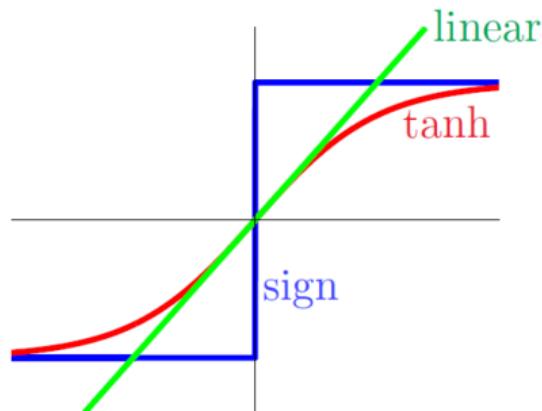
$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^m x_j^{(k-1)} w_{j,i}^{(k)}$$



Combining Perceptrons

Consider a simple perceptron $z(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$:

- ▶ $\theta(z) = \text{sign}(z)$: Learning the weights is hard combinatorial problem (not smooth)
- ▶ $\theta(z) = \tanh(z)$: differentiable approximation to $\text{sign}(\cdot)$ that allows analytic methods for learning

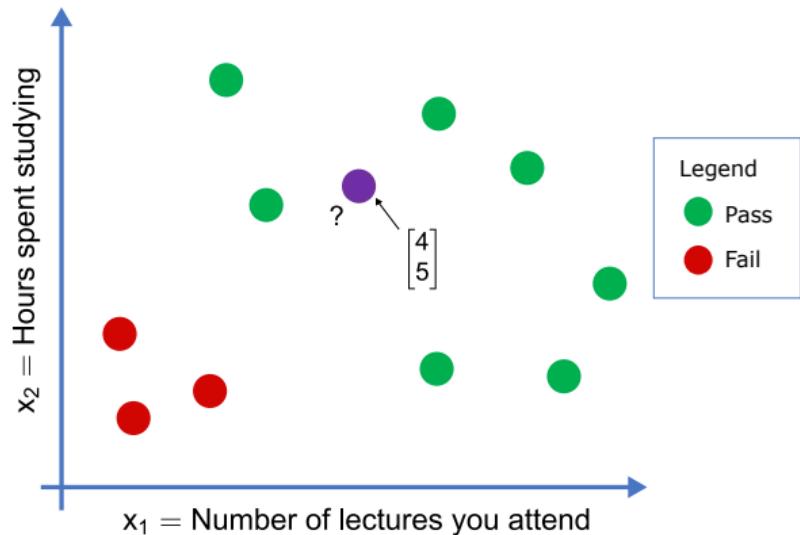


$$\theta(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

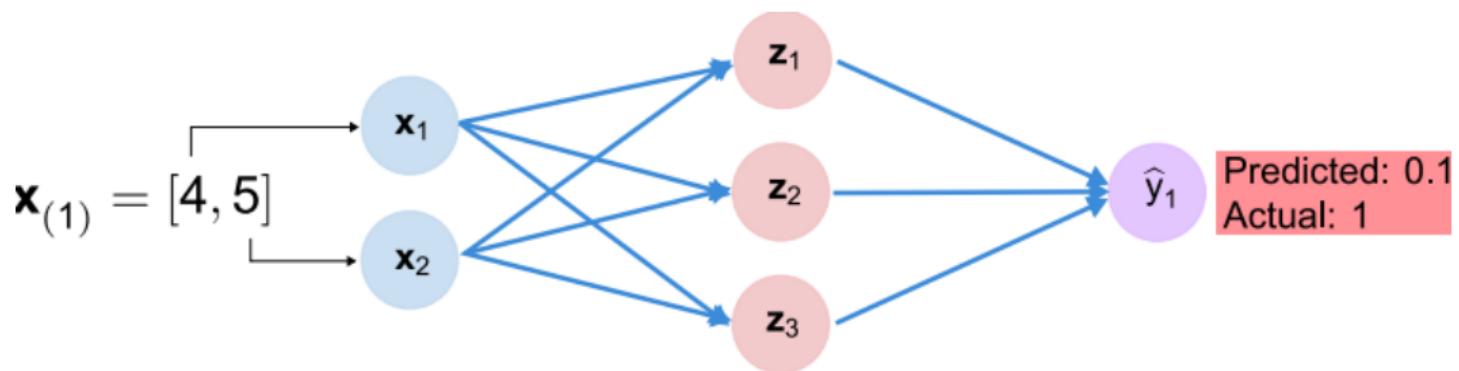
Example problem

Will I pass this class? Let's start with a simple two feature model:

- ▶ x_1 Number of lectures you attend
- ▶ x_2 Hours spent studying

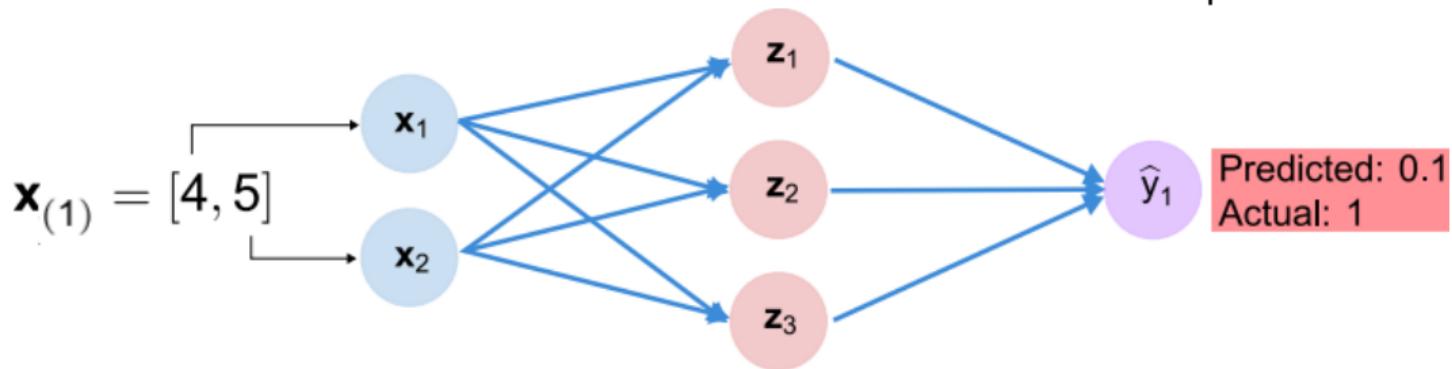


Example problem: Will I pass this class?



Quantifying Loss

The **loss** of a network measures the cost incurred from incorrect predictions



$$\mathcal{L}(f(\mathbf{x}_{(i)}; \mathbf{W}), \mathbf{y}_{(i)})$$

Predicted Actual

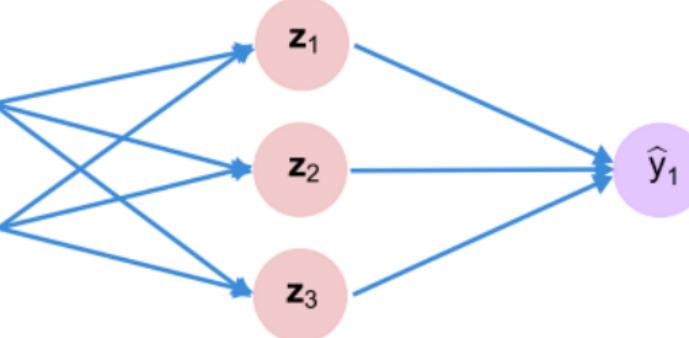
Empirical Loss

The **empirical loss** measures the total loss over the entire dataset

$$\mathbf{X} = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$

\mathbf{x}_1

\mathbf{x}_2



$$\begin{array}{c} f(\mathbf{x}) \\ \left[\begin{array}{c} 0.1 \\ 0.8 \\ 0.6 \\ \vdots \end{array} \right] \end{array} \quad \begin{array}{c} \mathbf{y} \\ \left[\begin{array}{c} \times \\ \times \\ \checkmark \\ \vdots \end{array} \right] \end{array}$$

Also Known as:

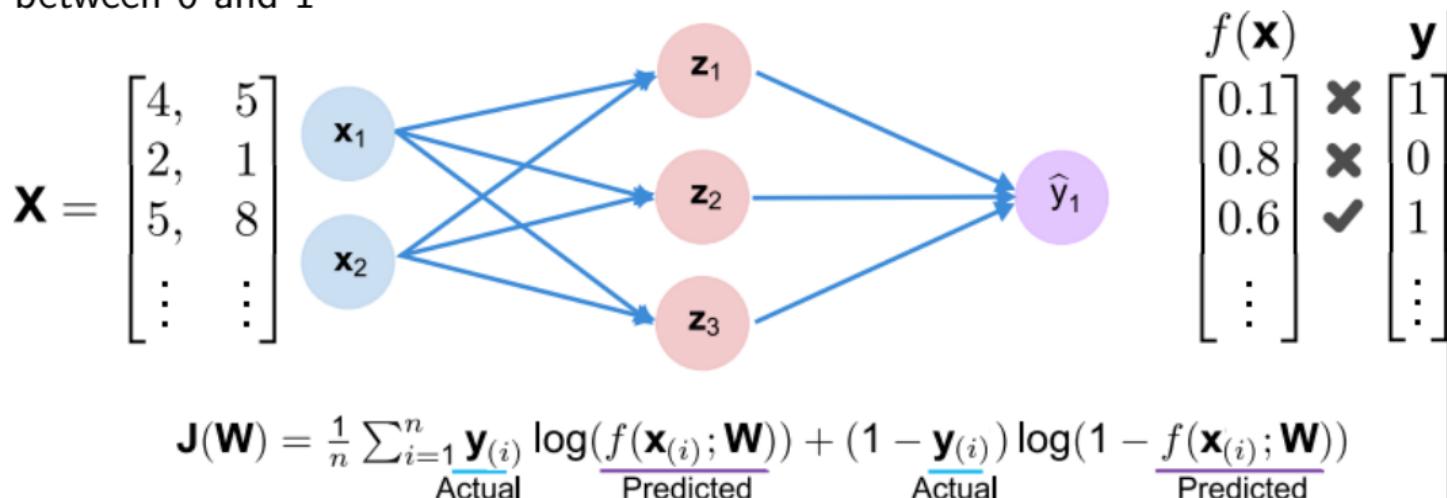
- Objective function
- Cost function
- Empirical Risk

$$\mathbf{J}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(\mathbf{x}_{(i)}; \mathbf{W}), \mathbf{y}_{(i)})$$

Predicted Actual

Binary Cross Entropy Loss

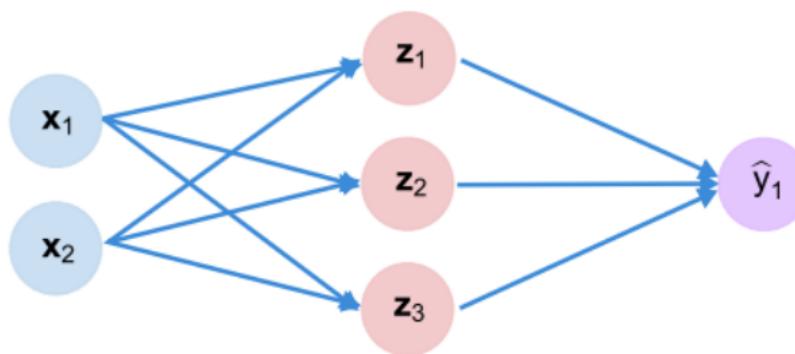
Cross entropy loss can be used with models that output a probability between 0 and 1



Mean Square Error Loss

Mean square error loss can be used with models that output a continuous real numbers

$$\mathbf{X} = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$

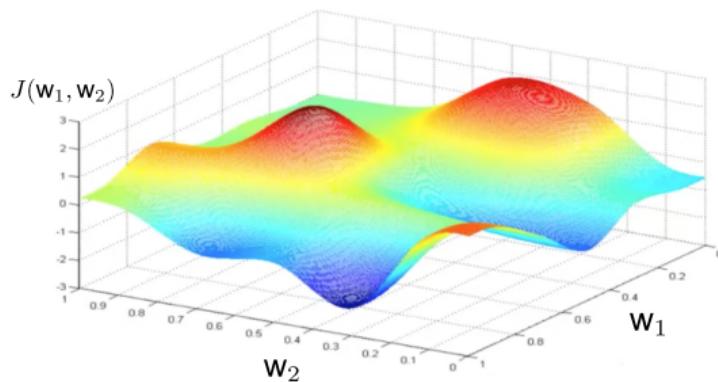


$f(\mathbf{x})$	\mathbf{y}
[30]	[90]
[80]	[20]
[85]	[95]
\vdots	\vdots

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \frac{(\mathbf{y}_{(i)} - f(\mathbf{x}_{(i)}; \mathbf{W}))^2}{\text{Actual} \quad \text{Predicted}}$$

Training Neural Networks: Loss Optimization

Find the network weights that **achieve the lowest loss**



$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(\mathbf{x}_i; \mathbf{W}), y_{(i)})$$

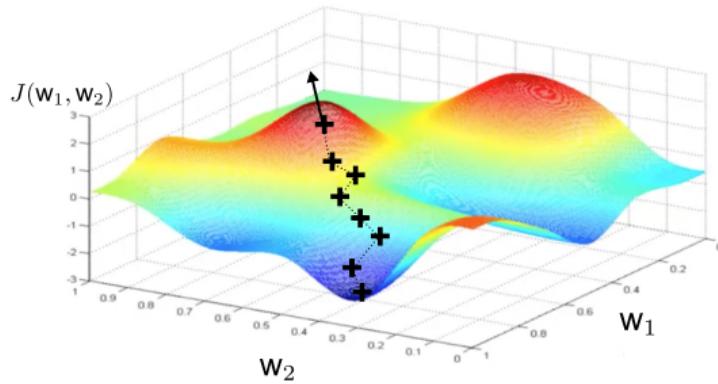
$$\mathbf{W}^* = \arg \min_{\mathbf{W}} J(\mathbf{W})$$

Remember:

$$\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$$

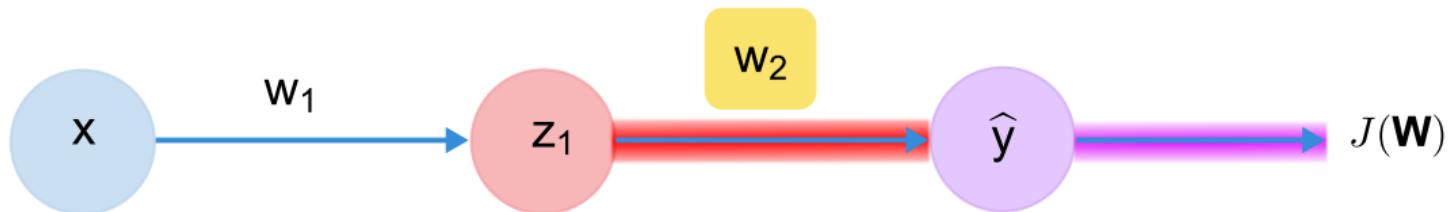
The loss function is a function of the network weights

Training Neural Networks: Loss Optimization



1. Randomly pick an initial (w_1, w_2)
2. Compute the gradient:
$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$
3. Take a small step in the opposite direction of the gradient
4. Repeat until convergence

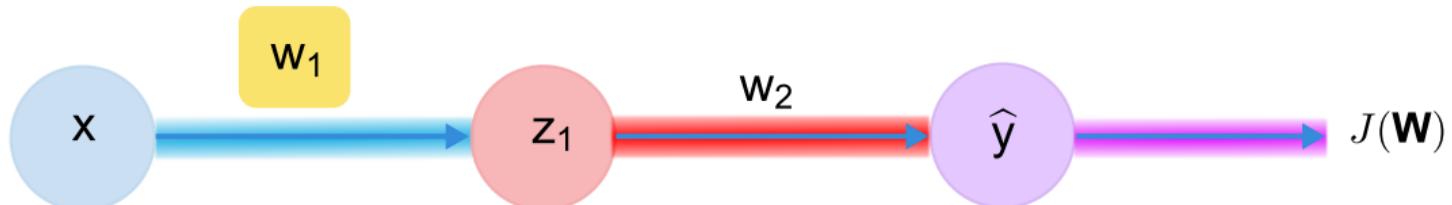
Computing Gradients: Backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_2} = \text{Use the chain rule!}$$

$$\frac{\partial J(\mathbf{W})}{\partial w_2} = \underline{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}} \underline{\frac{\partial \hat{y}}{\partial w_2}}$$

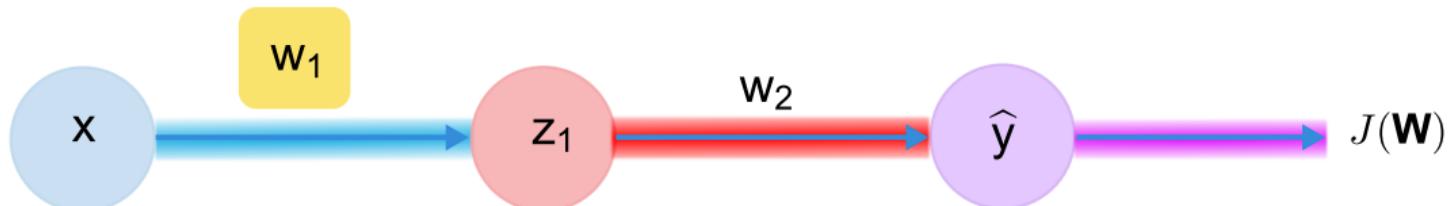
Computing Gradients: Backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underline{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}} \underline{\frac{\partial \hat{y}}{\partial z_1}} \underline{\frac{\partial z_1}{\partial w_1}}$$

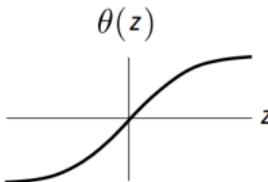
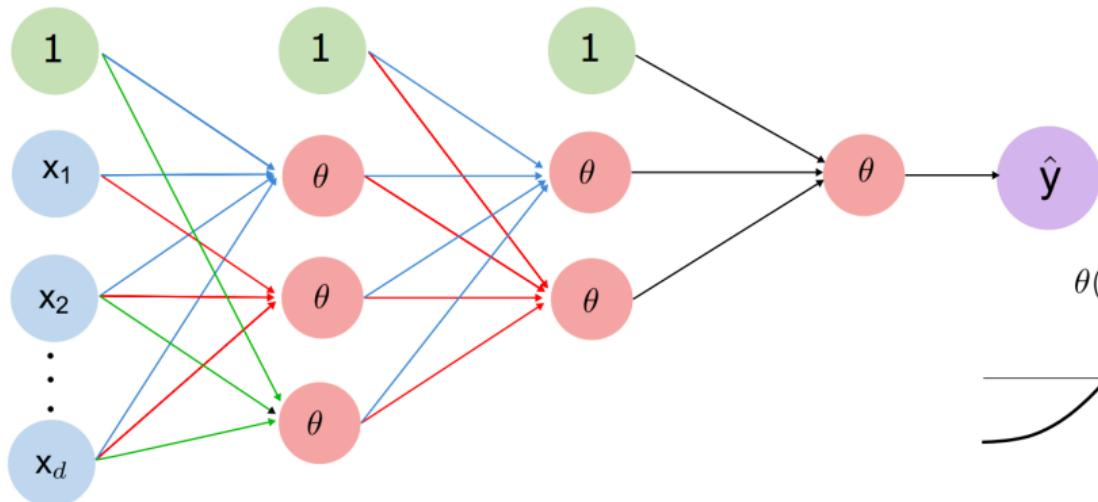
Computing Gradients: Backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underline{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}} \underline{\frac{\partial \hat{y}}{\partial z_1}} \underline{\frac{\partial z_1}{\partial w_1}}$$

How the network operates

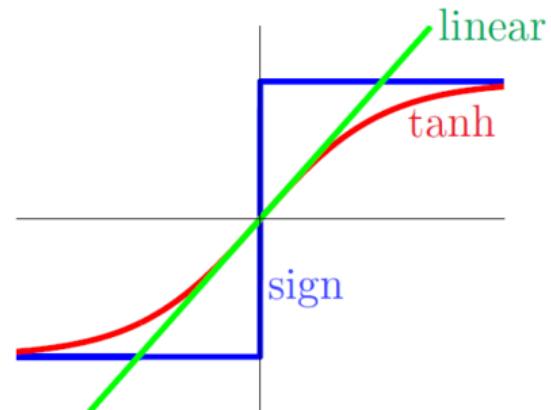


input \mathbf{x} hidden layers $1 \leq l \leq L$ output layer $l = L$

How the network operates

$$w_{ij}^{(l)} = \begin{cases} 1 \leq l \leq L & \text{Layers} \\ 0 \leq i \leq d^{(l-1)} & \text{Layers} \\ 1 \leq j \leq d^{(l)} & \text{Layers} \end{cases}$$

$$x_j^{(l)} = \theta(z_j^{(l)}) = \theta(\sum_{i=0}^{d^{(l-1)}} w_{ij}^l x_i^{(l-1)})$$



Apply \mathbf{x} to $x_1^{(0)} \dots x_{d^{(0)}}^{(0)} \rightarrow x_1^{(L)} = h(\mathbf{x})$

$$\theta(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Apply SGD

All the weights $\mathbf{W} = \{w_{ij}^{(l)}\}$ determine $f(\mathbf{x})$

Error on example (\mathbf{x}_n, y_n) is

$$\mathbf{J}(f(\mathbf{x}_n), y_n) = \mathbf{J}(\mathbf{W})$$

To implement SGD, we need the gradient

$$\nabla \mathbf{J}(\mathbf{w}) = \frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}}$$

for all i, j, l

Computing $\nabla J(\mathbf{W})$

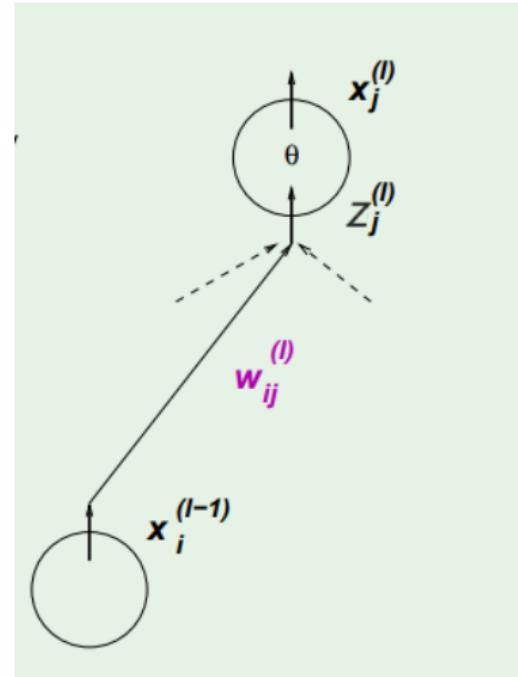
We can evaluate $\frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}}$ one by one: analytically or numerically

A trick for efficient computation:

$$\frac{\partial J(\mathbf{W})}{\partial w_{ij}^{(l)}} = \frac{\partial J(\mathbf{W})}{\partial z_j^{(l)}} \times \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}}$$

We have $\frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)}$

W only need: $\frac{\partial J(\mathbf{W})}{\partial z_j^{(l)}} = \delta_i^{(l-1)}$



δ for the final layer

For the final layer $l = L$ and $j = 1$

$$\delta_i^{(L-1)} = \frac{\partial J(\mathbf{W})}{\partial z_j^{(L)}}$$

$$J(\mathbf{W}) = (x_1^{(L)} - y_n)^2$$

$$\delta_i^{(L-1)} = \frac{\partial J(\mathbf{W})}{\partial z_j^{(L)}}$$

$$J(\mathbf{W}) = (\theta(z_1^{(L)}) - y_n)^2$$

Since $\theta'(z) = 1 - \theta^2(z)$ for the tanh

$$\frac{\partial J(\mathbf{W})}{\partial z_1^{(L)}} = 2(\theta(z_1^{(L)}) - y_n)(1 - \theta^2(z^{(L)}))$$

$$\frac{\partial J(\mathbf{W})}{\partial z_1} = 2(x_1^{(L)} - y_n)(1 - (x_1^{2(L)}))$$

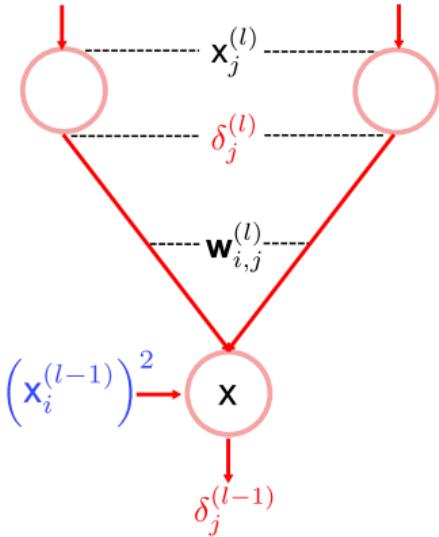
Backpropagation of δ

$$\delta_i^{(l-1)} = \frac{\partial e(\mathbf{w})}{\partial z_j^{(l)}}$$

$$\delta_i^{(l-1)} = \sum_{j=1}^{d^{(l)}} \frac{\partial e(\mathbf{w})}{\partial z_j^{(l)}} \times \frac{\partial z_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial z_i^{(l-1)}}$$

$$\delta_i^{(l-1)} = \sum_{j=1}^{d^{(l)}} \delta_j^{(l)} \times w_{ij}^{(l)} \times \theta'(z_i^{(l-1)})$$

$$\delta_i^{(l-1)} = (1 - (x_i^{(l-1)})^2) \sum_{j=1}^{d^{(l)}} w_{ij}^{(l)} \delta_j^{(l)}$$



Backpropagation Algorithm

- ▶ Initialize all weights $w_{ij}^{(l)}$ at small values chosen **at random**
- ▶ Pick one sample $n \in \{1, 2, \dots, N\}$ uniformly at random
- ▶ **Forward Part:** Compute all $x_j^{(l)}$
- ▶ **Backward Part:** Compute all $\delta_j^{(l)}$
- ▶ Update the weights:

$$\begin{aligned} w_{ij}^{(l)} &\leftarrow w_{ij}^{(l)} - \eta \frac{\partial J(\mathbf{w})}{\partial w_{ij}^{(l)}} \\ w_{ij}^{(l)} &\leftarrow w_{ij}^{(l)} - \eta x_i^{(l-1)} \delta_j^{(l)} \end{aligned}$$

- ▶ Iterate until it is time to stop.

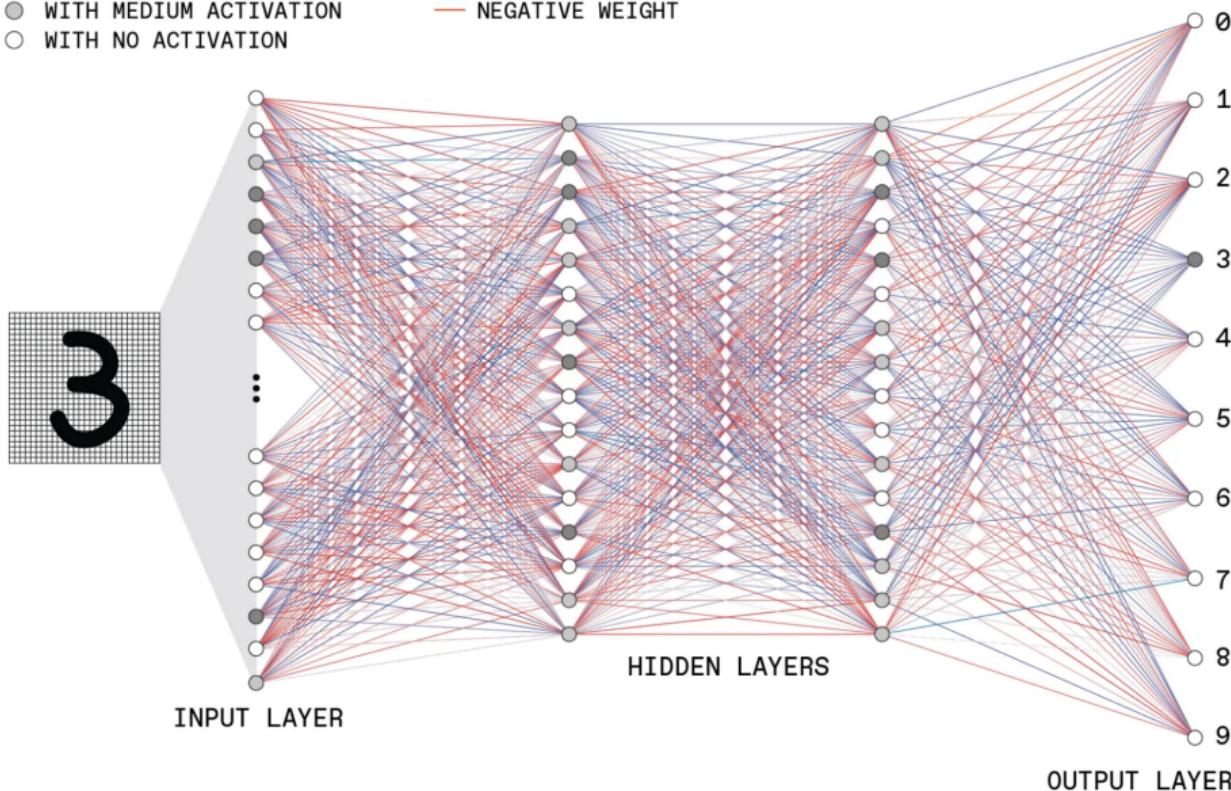
When is the best time to stop?

Considering a marginal error improvement, a maximum error and also a maximum number of iterations is a good practice.

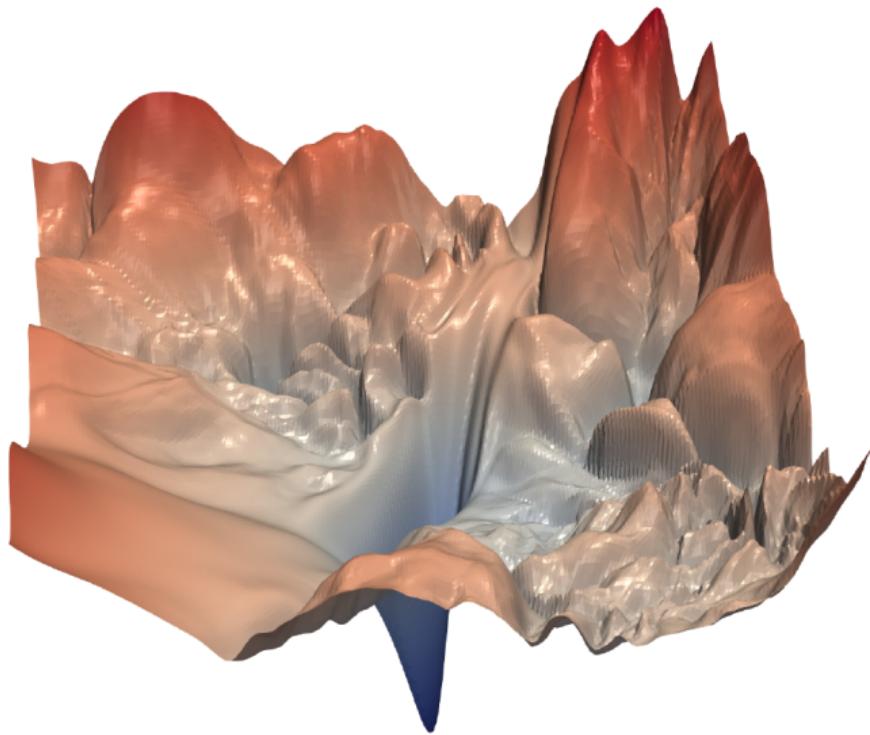
NEURONS:

- WITH HIGH ACTIVATION
- WITH MEDIUM ACTIVATION
- WITH NO ACTIVATION

— POSITIVE WEIGHT
— NEGATIVE WEIGHT



Training Neural Networks



Li, H. et al. Visualizing the loss landscape of neural nets (2017).

Loss Functions Can Be Difficult to Optimize

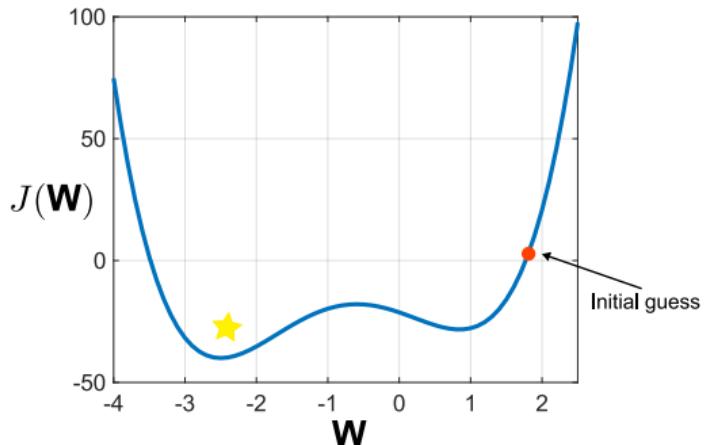
Remember: Optimization through gradient descent:

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$

$\eta \rightarrow$ How can we set the learning rate?

Training Neural Networks: Loss Optimization

- ▶ **Small learning rate** converges slowly and gets stuck in false local minima
- ▶ **Large learning rates** overshoot, become unstable and diverge
- ▶ **Stable learning rates** converge smoothly and avoid local minima



How to select the learning rate

Idea 1:

Try a lots of different learning rates and see what works "just right"

Idea 2:

Do something smarter!

Design an adaptive learning rate that "adapts" to the landspace

Adaptive Learning Rates

- ▶ Learning Rates are not longer fixed
- ▶ Can be made larger or smaller depending on:
 - ▶ how large the gradient is
 - ▶ how fast learning is happening
 - ▶ size of particular weights
 - ▶ etc...

Neural Networks in Practice: Mini-batches

Gradient Descent:

Compute:

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} \rightarrow \text{Can be very computational intensive to compute!}$$

Stochastic Gradient Descent:

Pick a single point i and compute:

$$\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}} \rightarrow \text{Easy to compute but very noisy (stochastic)!}$$

Neural Networks in Practice: Mini-batches

Mini-batches:

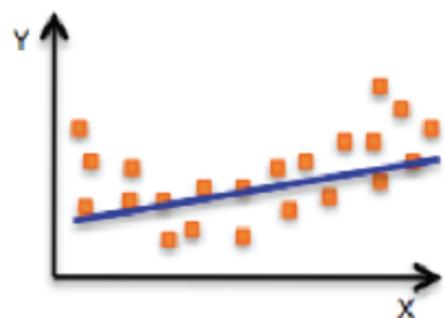
Pick B pints and compute:

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$$

- ▶ More accurate estimation of the gradient
 - ▶ Smoother convergence
 - ▶ Allows for large learning rates
- ▶ Lead to fast training
 - ▶ Can parallelize
 - ▶ Achieve significant speed increases on GPU's

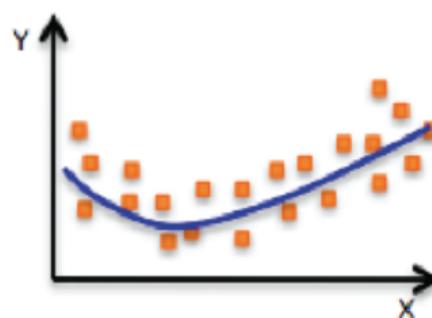
Neural Networks in Practice: Overfitting

The overfitting problem:

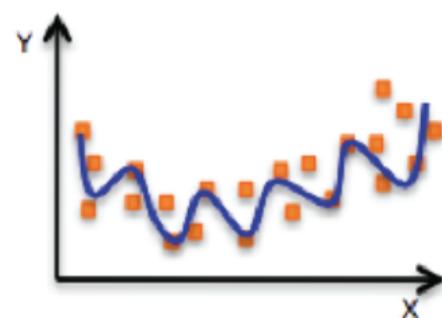


Underfitting

Model does not have capacity
to fully learn the data



Ideal Fit



Overfitting

Too complex, extra parameters,
does not generalize well

Neural Networks in Practice: Regularization

What is it?

Technique that constrains our optimization problem to discourage complex models

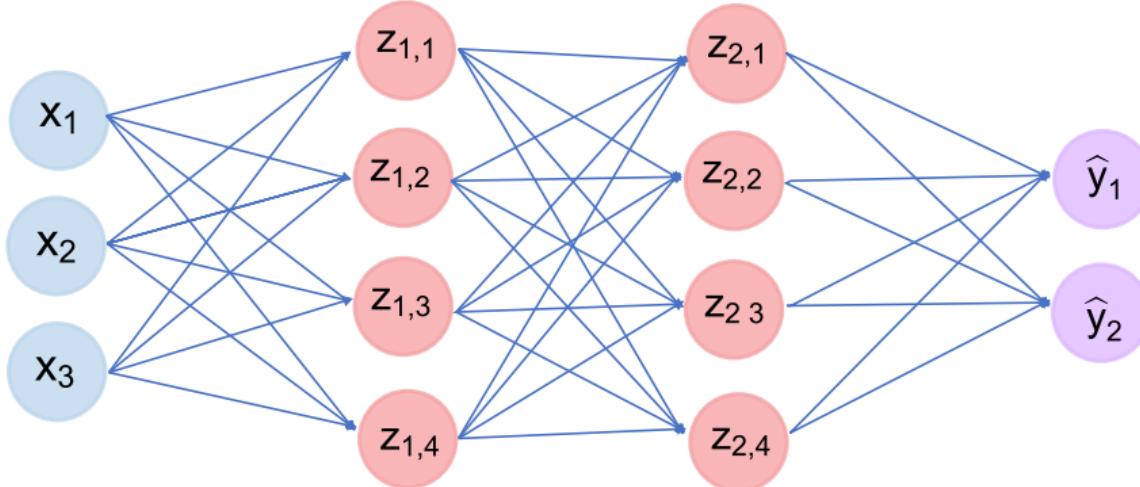
Why we need it?

Improve generalization of our models on unseen data

Neural Networks in Practice: Regularization 1

Dropout:

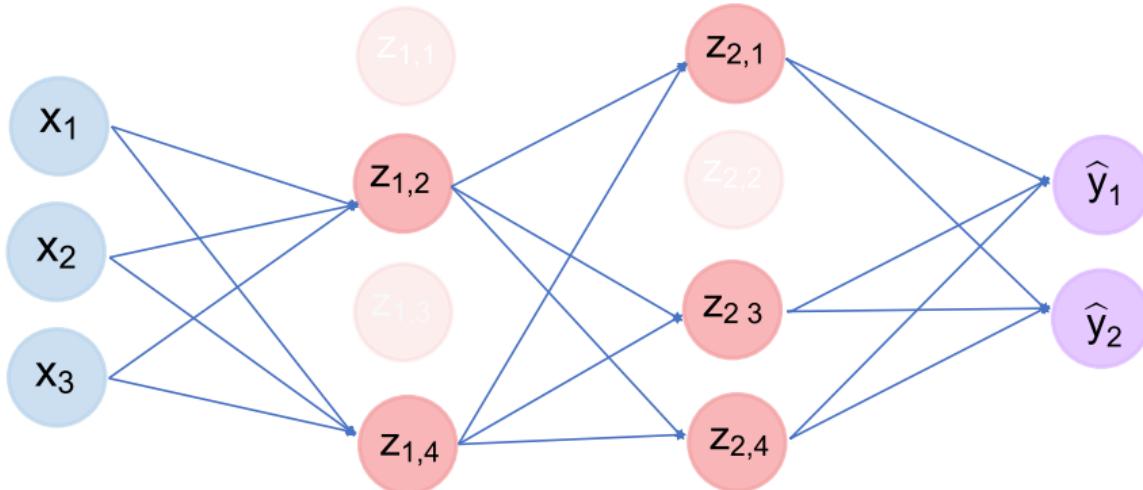
During training, at each iteration, randomly set some activations to 0



Neural Networks in Practice: Regularization 1

Dropout:

During training, at each iteration, randomly set some activations to 0

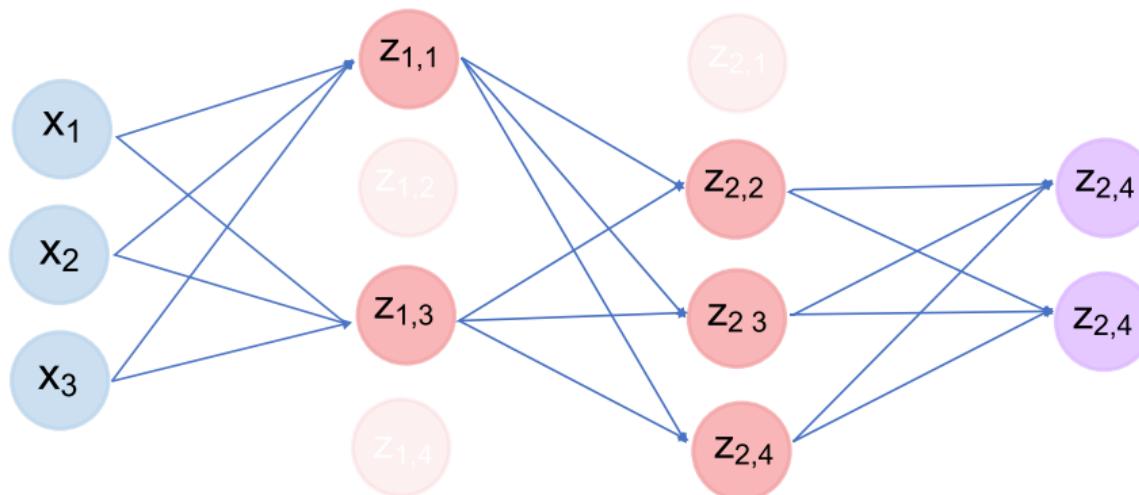


Neural Networks in Practice: Regularization 1

Dropout:

During training, at each iteration, randomly set some activations to 0

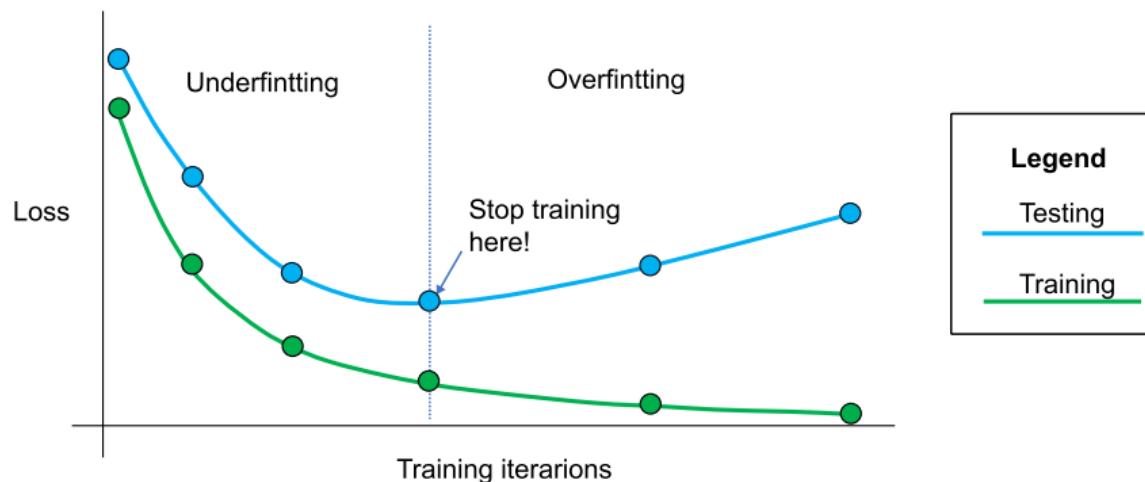
- ▶ Typically "drop" 50% of activations in layer
- ▶ Forces network to not rely on any node



Neural Networks in Practice: Regularization 2

Early Stopping

Stop training before we have a chance to overfit



Working with Real Data Sets: The Iris Data set

\mathbf{x} is a $150 \times 4 \rightarrow$ it has 4 features and we have 150 samples

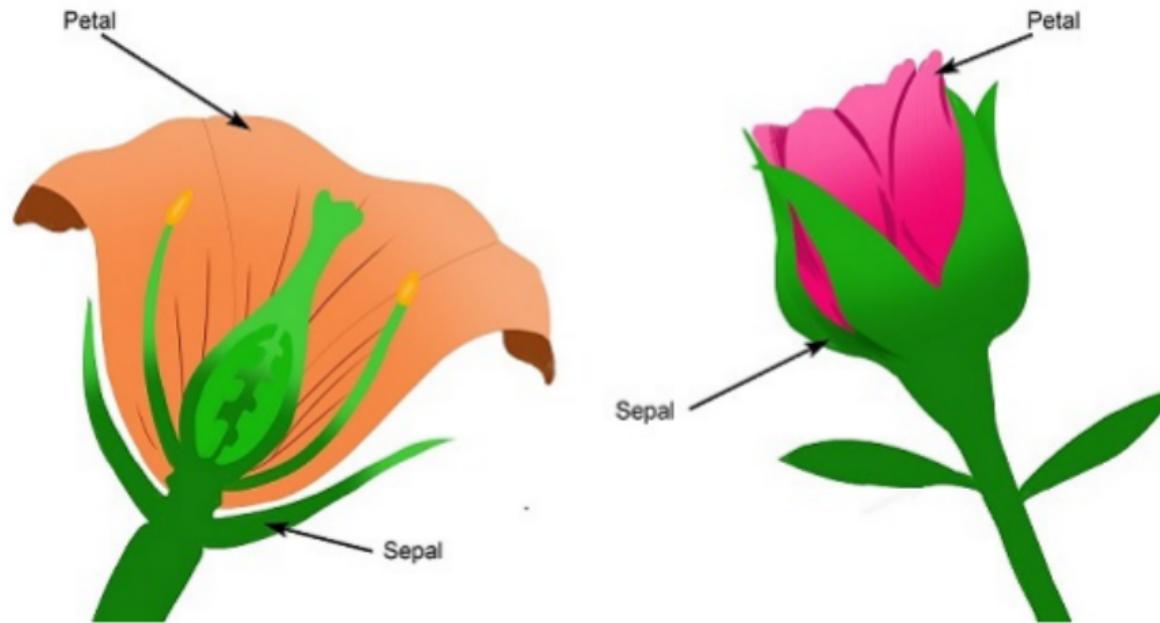
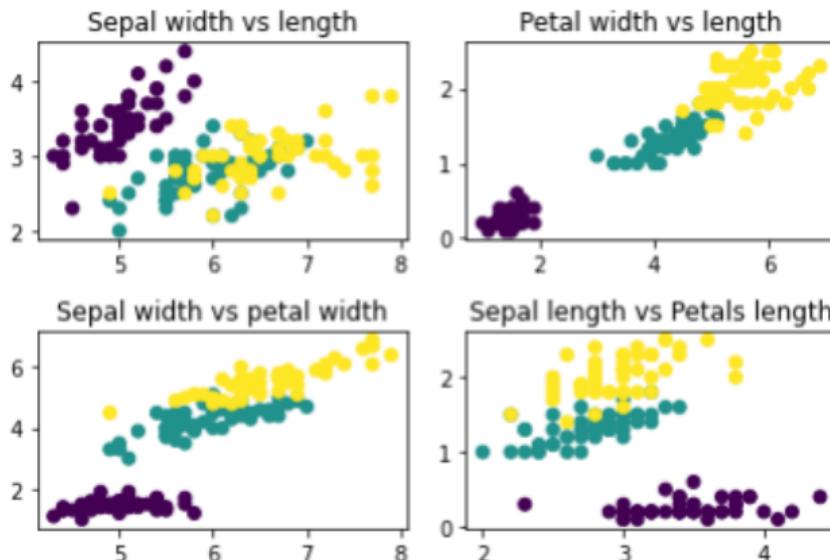


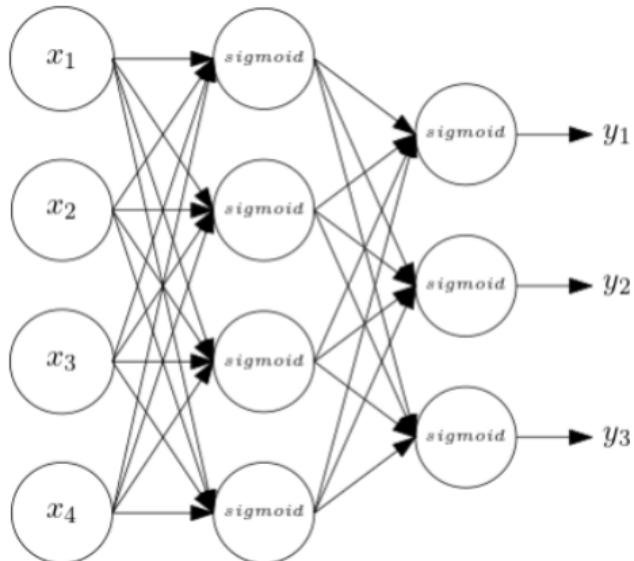
Image Features

1. Width
 - ▶ Petal
 - ▶ Sepal



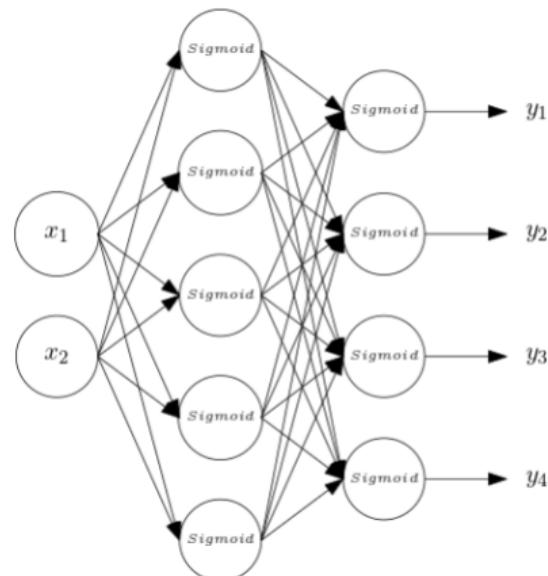
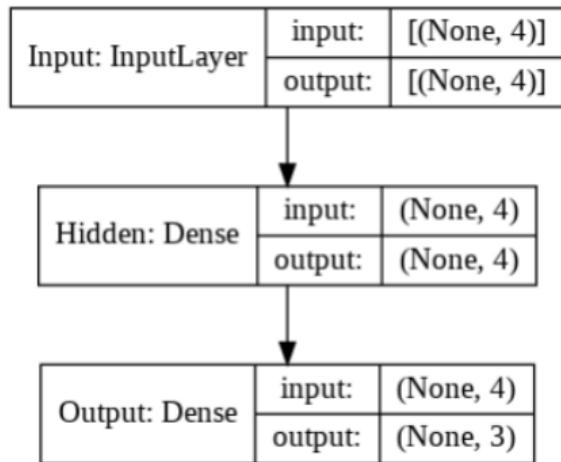
2. Length
 - ▶ Petal
 - ▶ Sepal

Approach



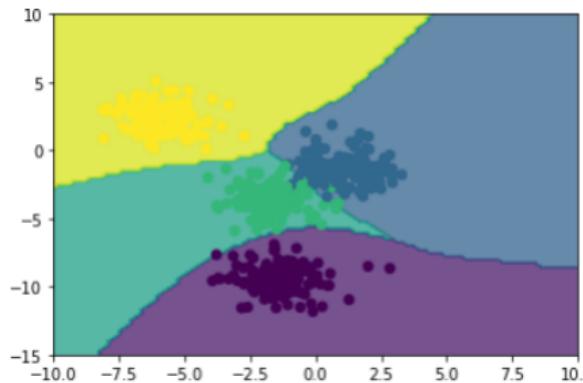
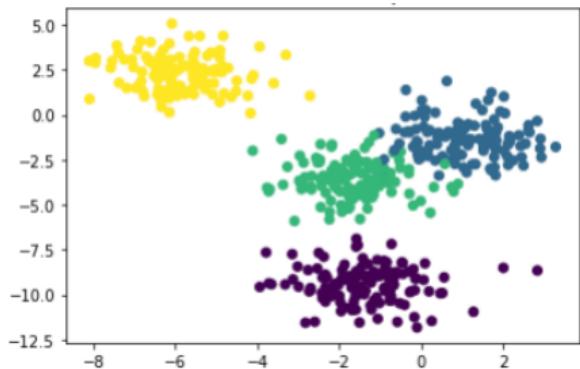
Generate fully connected layers with multiple activation functions to solve a multi feature classification problem.

An efficient coding method: Keras



Solutions

Accurate boundary lines for Classification



Code

```
input = keras.layers.Input(4, name="Input")
hidden = keras.layers.Dense(4, activation='sigmoid', name="Hidden")(input)
output = keras.layers.Dense(3, activation='sigmoid', name="Output")(hidden)

modelIRIS = keras.Model(input, output)

modelIRIS.layers[1].set_weights([w1,b1])
modelIRIS.layers[2].set_weights([w2,b2])

keras.utils.plot_model(modelIRIS, show_shapes=True)

# print(modelIRIS.predict(np.array([0,1,1,2],ndmin=2)))
```

```
X,y = datasets.make_blobs(centers = 4, random_state=2,n_samples = 500)

plt.scatter(X[:,0],X[:,1],c=y)
plt.title("Hand-made Example")
plt.show()
```

```
w1 = np.array([[ -1.5132886 ,  0.78692156, -0.59150934,  1.3457433 ,  0.11323632],
               [ -0.9712865 , -1.5561647 , -1.578427 , -0.7123345 ,  0.75948274]])
b1 = np.array([-2.8332472,  2.0457697, -3.736685 , -0.8151084,  4.7346272])

w2 = np.array([[ 1.3780868 , -4.0084467 ,  0.0883714 ,  2.5190408 ],
               [ 0.37726903,  0.76979834,  1.5470964 , -5.40783   ],
               [ 1.3585504 , -1.799717 ,  2.6277127 , -1.8256785 ],
               [ 2.7409966 ,  1.6214521 , -2.7639852 , -3.9925196 ],
               [-5.5055437 ,  2.8043787 ,  1.1856915 ,  3.2509704 ]])
b2 = np.array([-0.54020935,  0.57120425, -0.5757609 ,  0.8131391 ])
```

Acknowledgement

Alexander Amini and Ava Soleimany, MIT 6.S191: Introduction to Deep Learning, IntroToDeepLearning.com