



What is Deep Learning ?

Artificial Intelligence

✓ Any technique that enables computers to mimic human behavior

Machine Learning

 \checkmark 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 data?
 ✓ Data-driven approach
- Features in Face Images

 ✓ Low level features: Lines, edges
 ✓ Mid level features: Eyes, Nose, ears
 ✓ High level features: Facial structure

Why Now?

- Neural networks study
 - ✓ 1952: Stochastic gradient descent
 - ✓ 1958: Perceptron
 - ✓ 1986: Backpropagation
 - ✓ 1995: Deep Convolutional NN
- Neural networks date back decades, so why the resurgence?

✓ Big Data

> Larger datasets, easier collection & storage

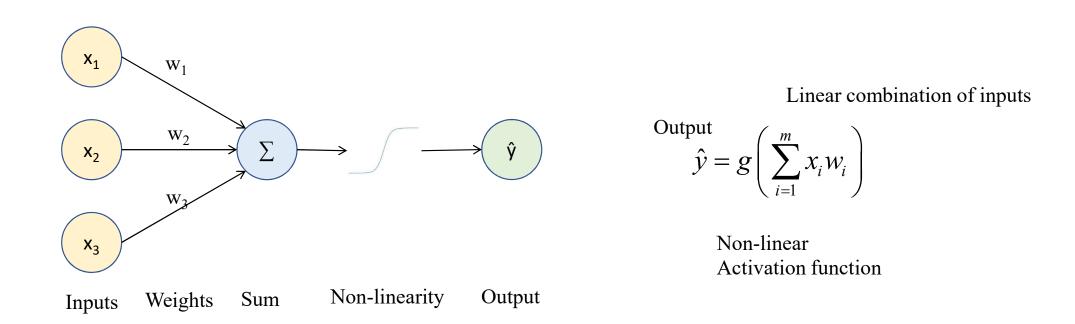
✓ Hardware

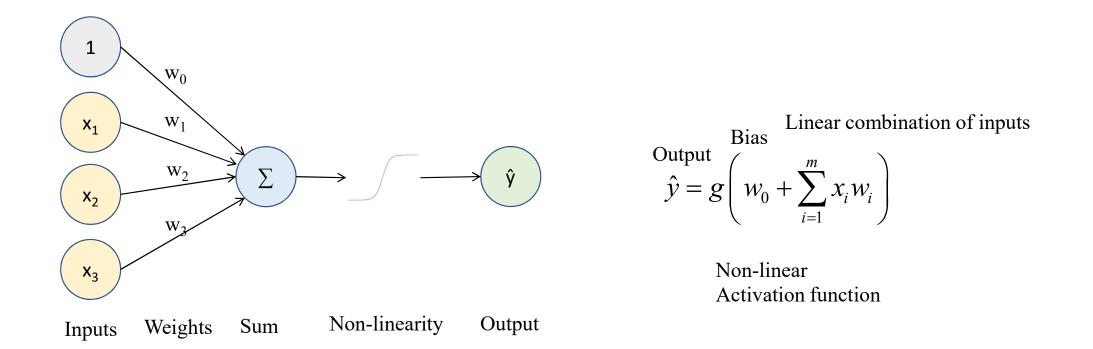
> Graphics processing units (GPUs), massively parallelizable

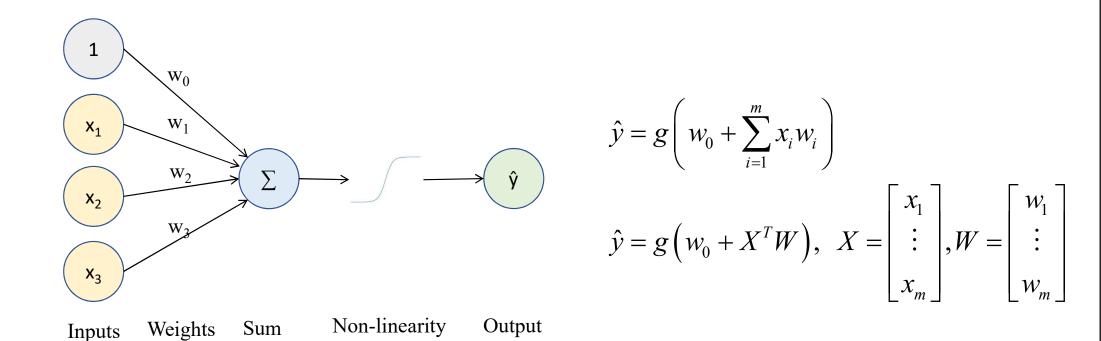
✓ Software

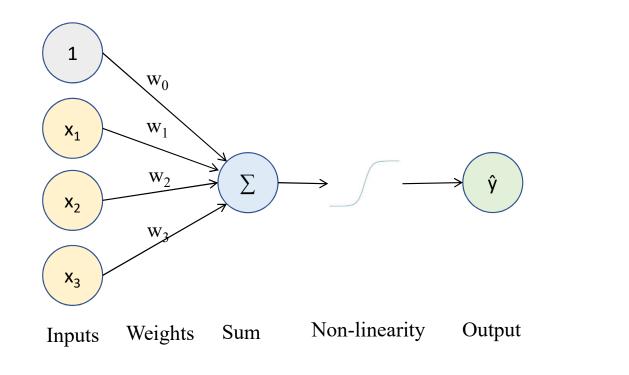
> Improved techniques, new models, toolboxes

The Perceptron (The structural building block of deep learning)





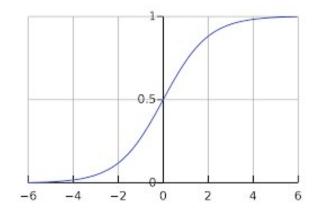




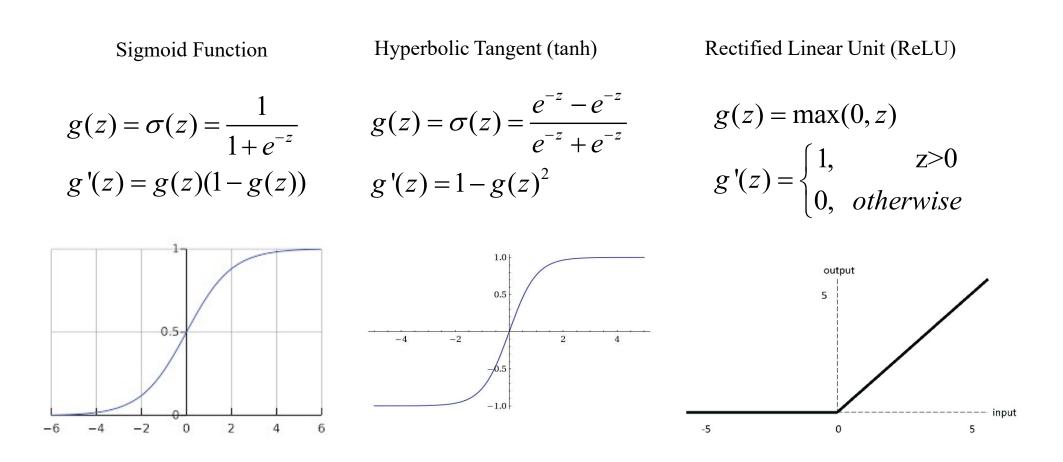
Activation Function

$$\hat{y} = g\left(w_0 + X^T W\right)$$

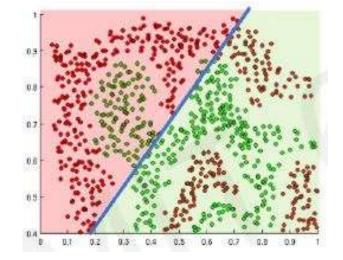
Sigmoid Function $g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$

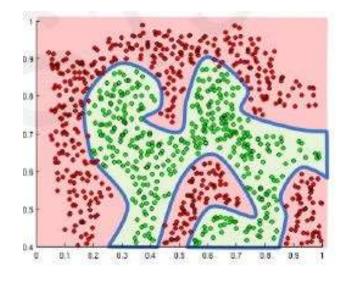


Common Activation Functions



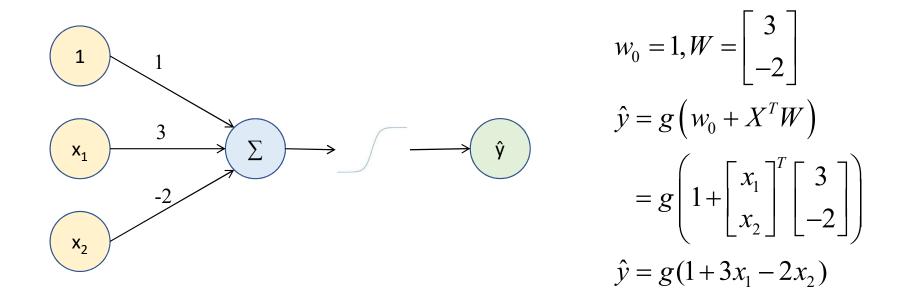
Importance of Activation Functions



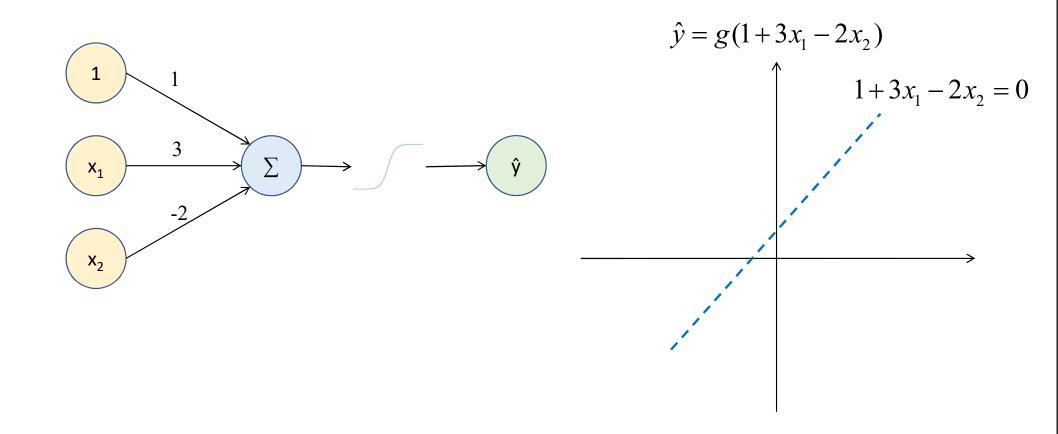


Linear activation functions produce linear decisions no matter the network size Non-linearity allows us to approximate arbitrarily complex functions

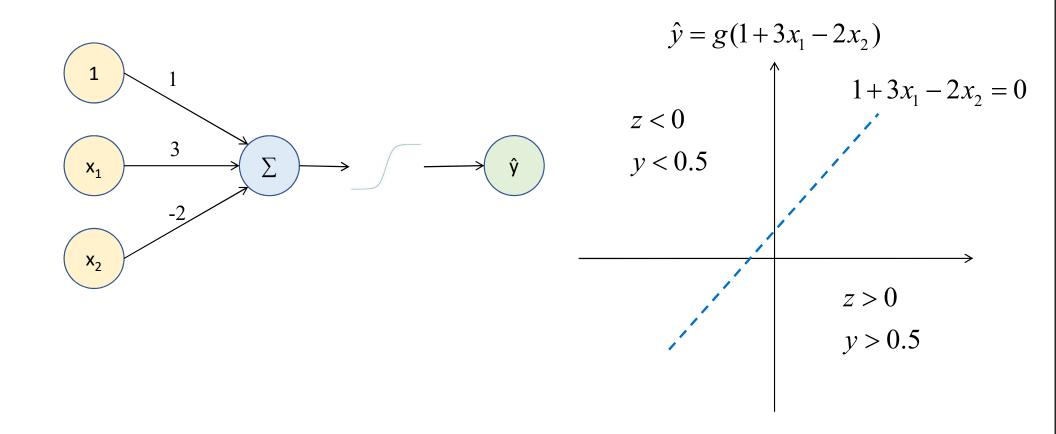
The Perceptron: Example



The Perceptron: Example

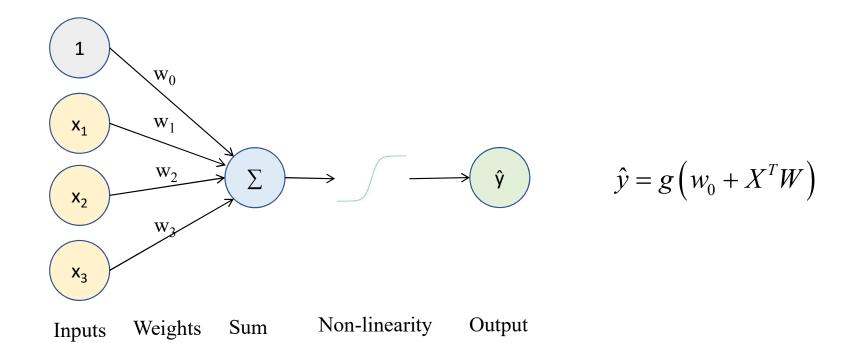


The Perceptron: Example

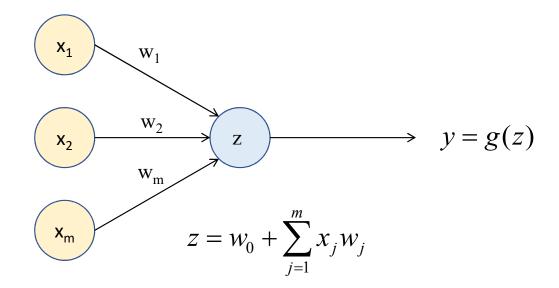


Building Neural Networks with Perceptrons

The Perceptron: Simplified

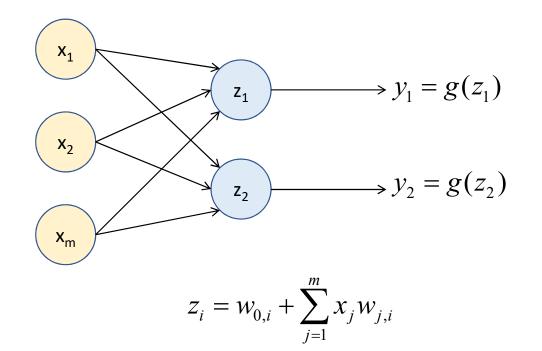


The Perceptron: Simplified

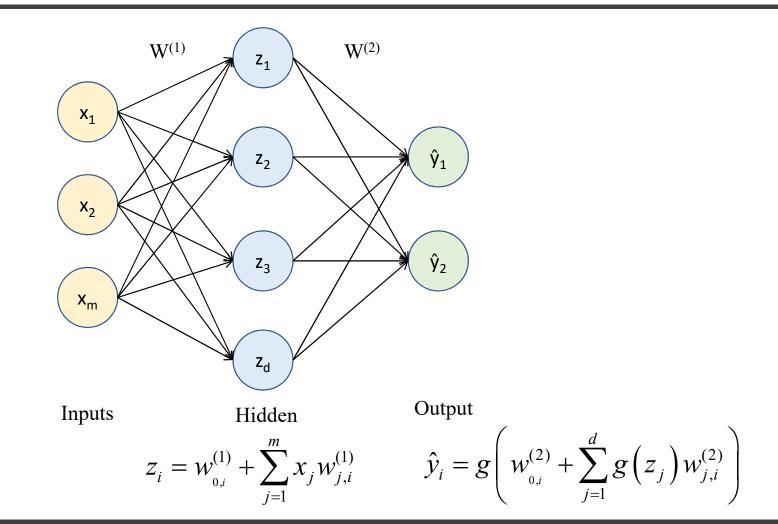


Multi Outputs

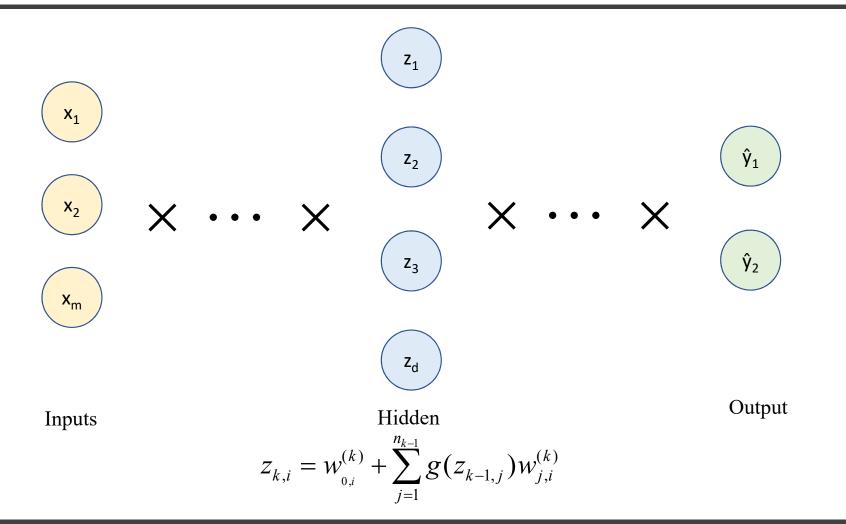
 Because all inputs are densely connected to all outputs, these layers are called Dense layers



Single Layer Neural Network



Deep Neural Network

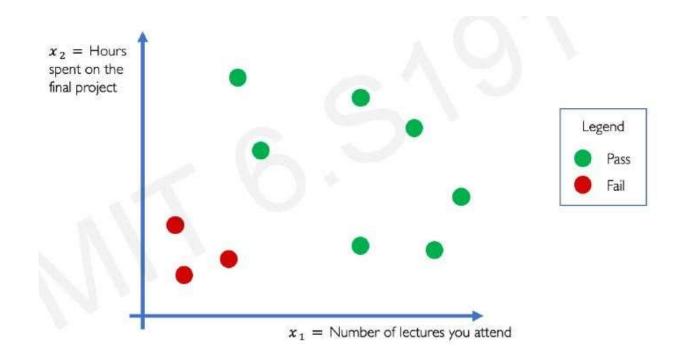


Applying Neural Networks

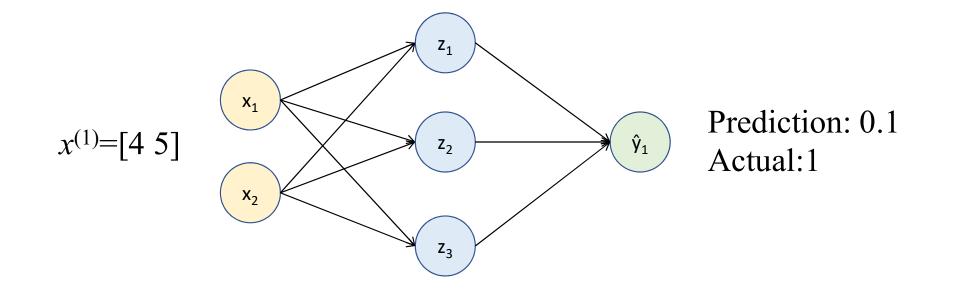
Example Pass the class

Let's start with a simple two feature model
 ✓ x₁ = # of lectures you attend
 ✓ x₂ = Hours spent on the final project

Example

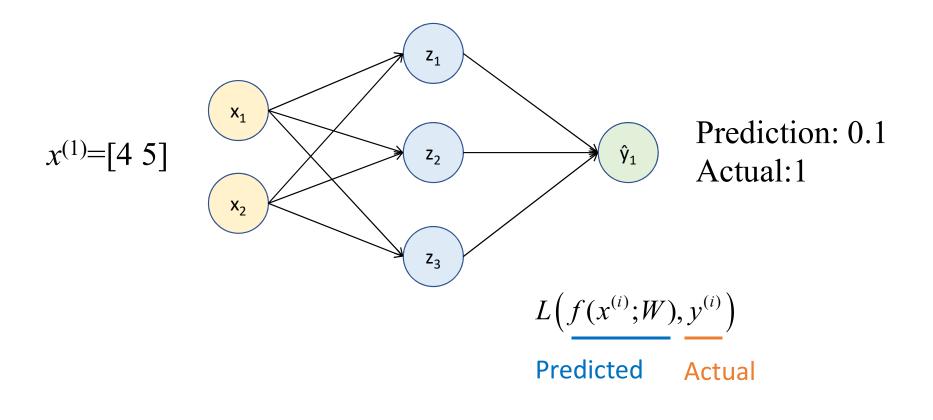


Example



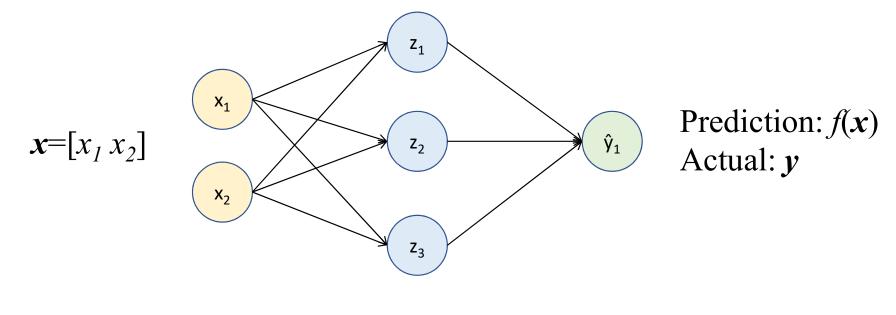
Quantifying Loss

 The loss of our network measures the cost incurred from incorrect predictions



Empirical Loss

• The empirical loss measures the total loss over our entire dataset



Objective function Cost function Empirical Risk

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

Mean Squared Error Loss

 Mean squared error loss can be used with regression models that output continuous real numbers

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} \left(y^{(i)} - f(x^{(i)}; W) \right)^2$$

Actual Predicted

Binary Cross Entropy Loss

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

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} y^{(i)} \log\left(f(x^{(i)}; W)\right) + (1 - y^{(i)}) \log\left(1 - f(x^{(i)}; W)\right)$$

Actual Predicted Actual Predicted

Loss Optimization

• We want to find the network weights that achieve the lowest loss

$$W^* = \underset{W}{\operatorname{arg\,min}} \frac{1}{n} \sum_{i=1}^n L(f(x^{(i)}; W), y^{(i)})$$
$$= \underset{W}{\operatorname{arg\,min}} J(W)$$

Gradient Descent Algorithm

- Loss is a function of the network weights
- Gradient Descent

1. Initialize weights randomly ~ $N(0, \sigma^2)$

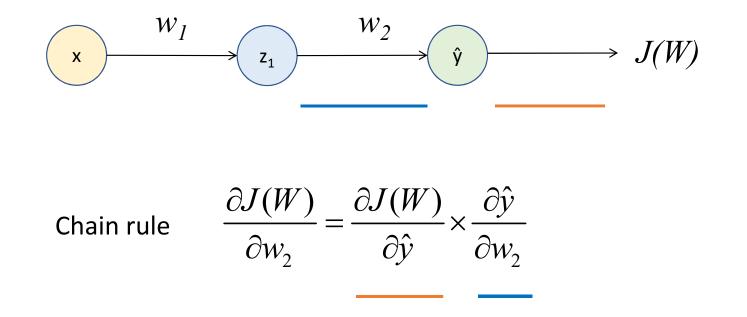
2. Loop until convergence:

3. Compute gradient,
$$\frac{\partial J(W)}{\partial W}$$

4. Update weights, $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$

5. Return weights

Computing Gradients: Backpropagation



Computing Gradients: Backpropagation

$$\begin{array}{c} x & \stackrel{W_{1}}{\longrightarrow} & \stackrel{V_{2}}{\longrightarrow} & \hat{y} & \stackrel{W_{2}}{\longrightarrow} & J(W) \\ \\ & \frac{\partial J(W)}{\partial w_{1}} = \frac{\partial J(W)}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_{1}} \times \frac{\partial z_{1}}{\partial w_{1}} \end{array}$$

Repeat this for every weight in the network in the network using gradients from later layers

Neural Networks in Practice: Optimization

Loss functions can be difficult to optimize

Optimization through gradient descent

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

Learning rate - How can we set ?

✓ Small learning rates converge 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

Adaptive Learning Rates

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
 - ✓ How large gradient is
 ✓ How fast learning is happening
 ✓ Size of particular weights
 - ✓ Etc.

Optimizer - Adaptive Learning Rates

Algorithm

✓ SGD

✓ Adam

✓ Adadelta

 $\checkmark \mathsf{Adagrad}$

✓ RMSProp

Neural Networks in Practice: Mini-batches

Gradient Descent

- Loss is a function of the network weights
- Gradient Descent

1. Initialize weights randomly ~ $N(0, \sigma^2)$

2. Loop until convergence:

3. Compute gradient, $\frac{\partial J(W)}{\partial W}$ Can be very computationally intensive to compute 4. Update weights, $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$

5. Return weights

Stochastic Gradient Descent

- Stochastic Gradient Descent
 - 1. Initialize weights randomly ~ $N(0, \sigma^2)$
 - 2. Loop until convergence:
 - 3. Pick single data point *i*
 - 4. Compute gradient, $\frac{\partial J_i(W)}{\partial W}$

Easy to compute but very noisy (stochastic)

5. Update weights,
$$W \leftarrow W - \eta \frac{\partial J_i(W)}{\partial W}$$

6. Return weights

Stochastic Gradient Descent

- Stochastic Gradient Descent
 - 1. Initialize weights randomly ~ $N(0, \sigma^2)$
 - 2. Loop until convergence:
 - 3. Pick batch of B data points

4. Compute gradient,
$$\frac{\partial J(W)}{\partial W} = \frac{1}{B} \sum_{k=1}^{B} \frac{\partial J_k(W)}{\partial W}$$

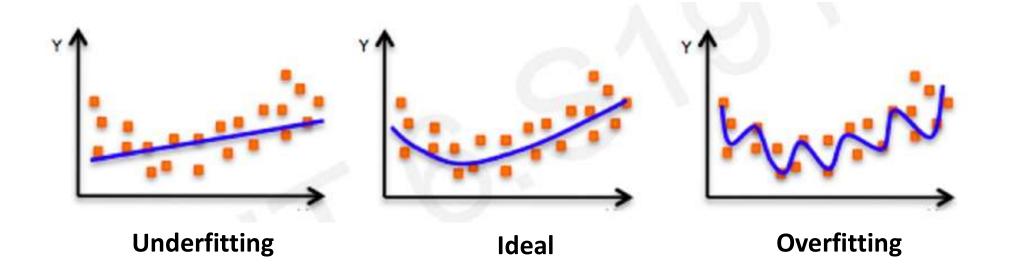
5. Update weights, $W \leftarrow W - \eta \frac{\partial J_i(W)}{\partial W}$

Fast to compute and a much Better estimate of the true gradient

6. Return weights

Neural Networks in Practice: Overfitting

The problem of overfitting



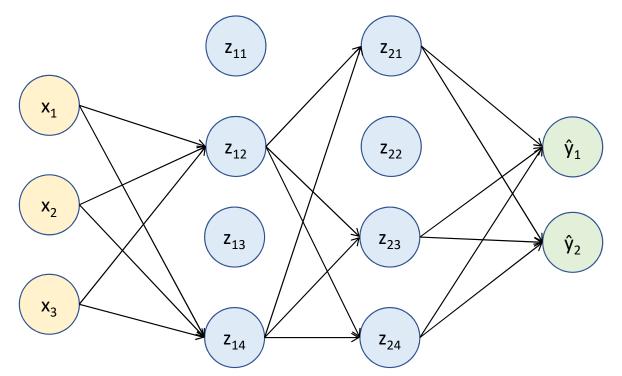
- Underfitting Model does not have capacity to fully learn the data
- Overfitting Too complex, extra parameters, does not generalize well

Regularization

- Technique that constrains our optimization problem to discourage complex models
- Improve generalization of our model on unseen data

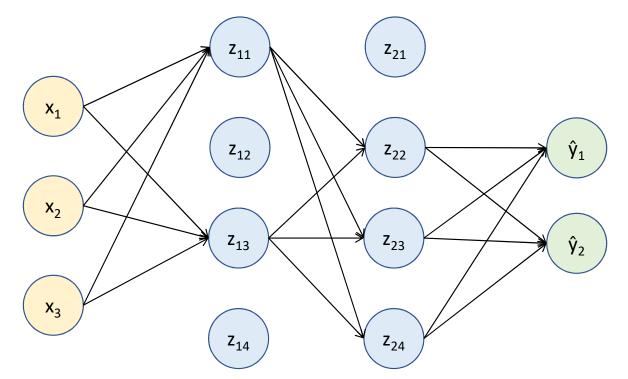
Regularization – Dropout

During training, randomly set some activations to 0
 ✓ Typically 'drop 50%' of activations in layer
 ✓ Forces network to not rely on any *i* node



Regularization – Dropout

During training, randomly set some activations to 0
 ✓ Typically 'drop 50%' of activations in layer
 ✓ Forces network to not rely on any *i* node



Regularization – Early Stopping

Stop training before we have a chance to overfit





