

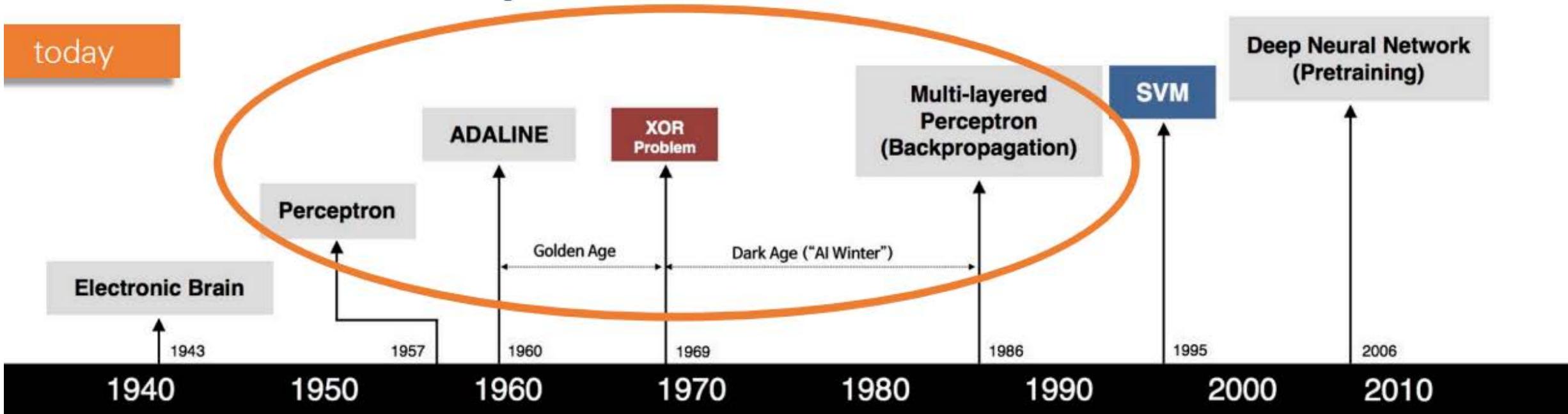
Lecture 6: Deep Neural Networks and Training

Zerrin Yumak
Utrecht University

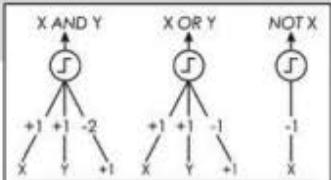
In this lecture

- Feedforward neural networks
- Activation functions
- Backpropagation
- Regularization
- Dropout
- Optimization algorithms
- Weight initialization
- Batch normalization
- Hyper parameter tuning

today



S. McCulloch - W. Pitts



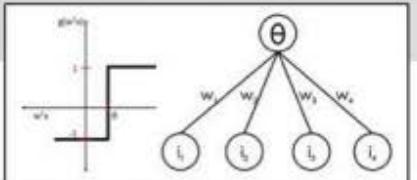
- Adjustable Weights
- Weights are not Learned



F. Rosenblatt



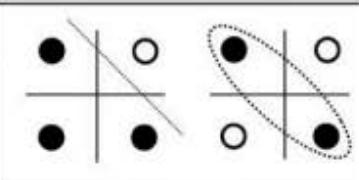
B. Widrow - M. Hoff



- Learnable Weights and Threshold



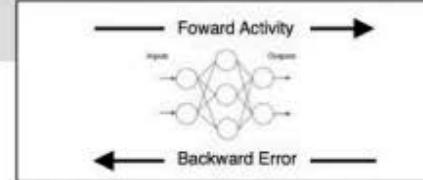
M. Minsky - S. Papert



- XOR Problem



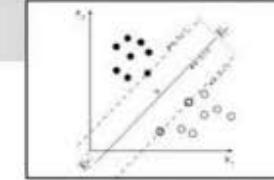
D. Rumelhart - G. Hinton - R. Williams



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



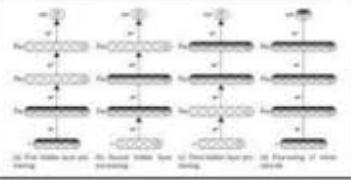
V. Vapnik - C. Cortes



- Limitations of learning prior knowledge
- Kernel function: Human Intervention



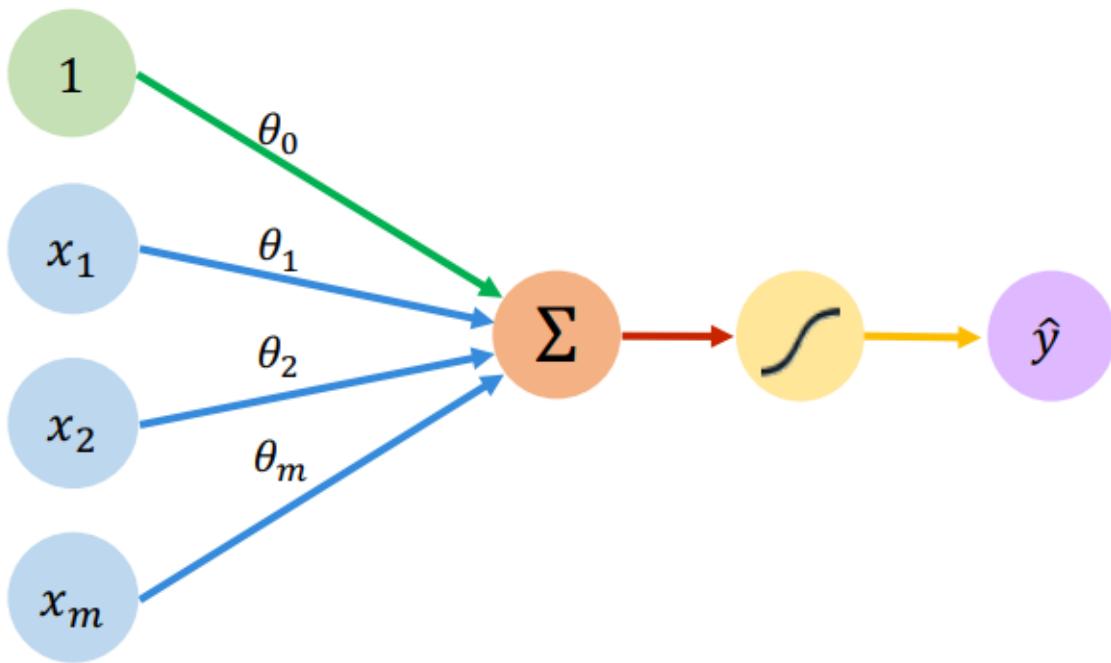
G. Hinton - S. Ruslan



- Hierarchical feature Learning

The Perceptron

- Building block of deep neural networks



Inputs Weights Sum Non-Linearity Output

Output \downarrow

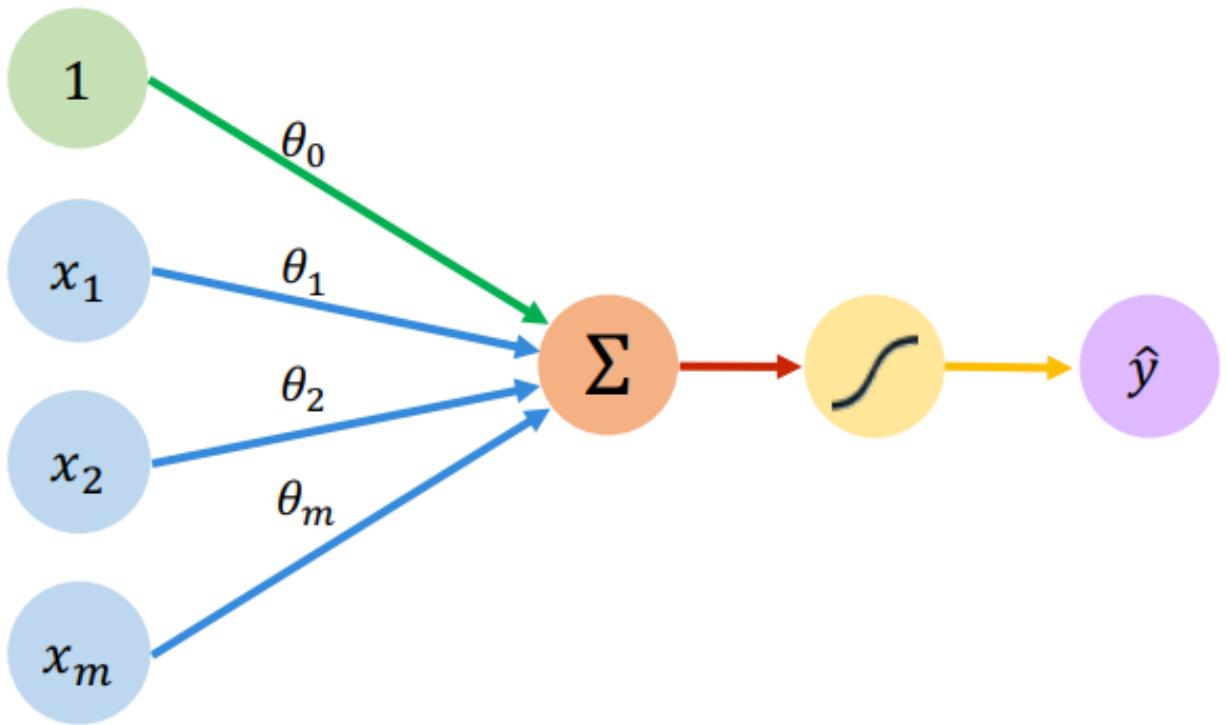
$$\hat{y} = g \left(\theta_0 + \sum_{i=1}^m x_i \theta_i \right)$$

Linear combination of inputs \downarrow

Non-linear activation function \uparrow

Bias \uparrow

The Perceptron



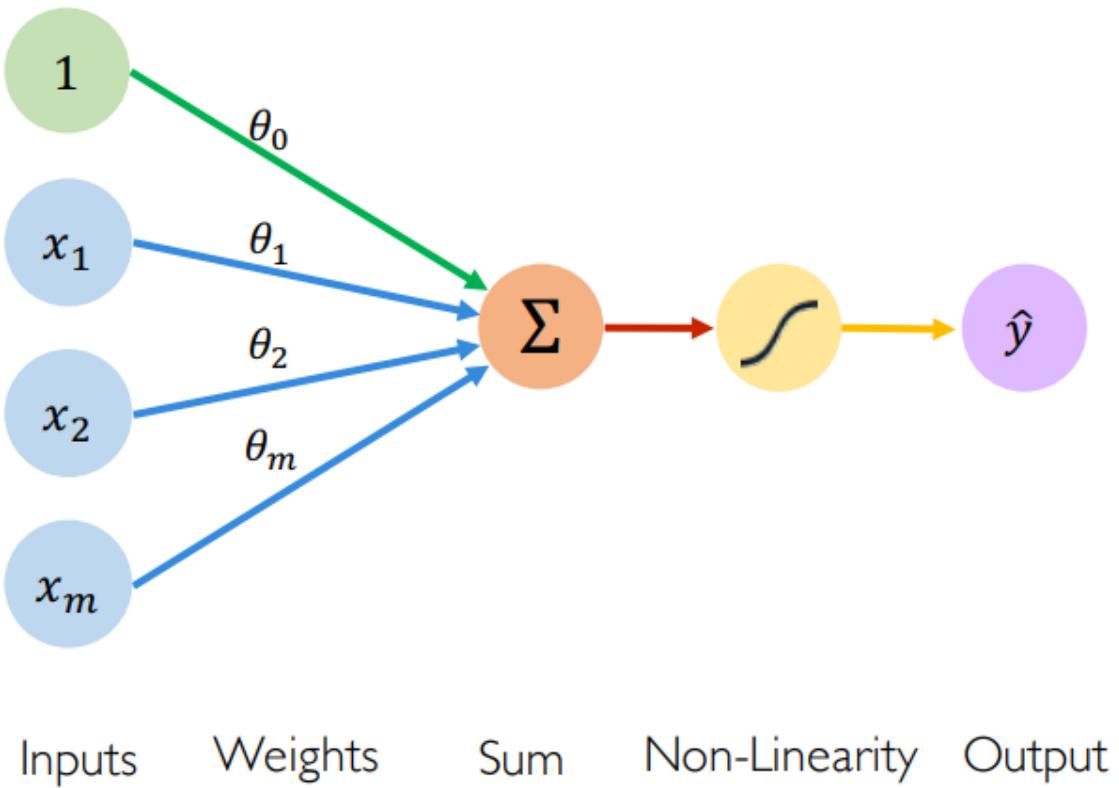
Inputs Weights Sum Non-Linearity Output

$$\hat{y} = g \left(\theta_0 + \sum_{i=1}^m x_i \theta_i \right)$$

$$\hat{y} = g(\theta_0 + \mathbf{X}^T \boldsymbol{\theta})$$

where: $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$

The Perceptron

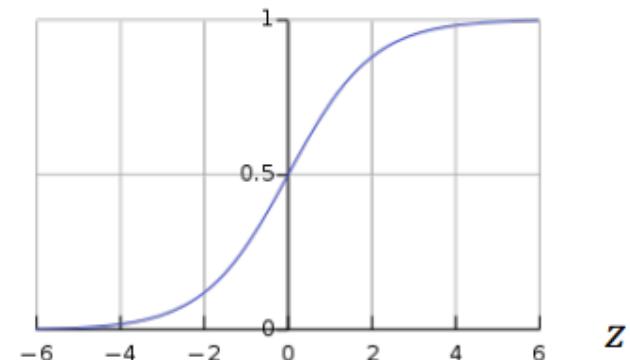


Activation Functions

$$\hat{y} = g(\theta_0 + X^T \theta)$$

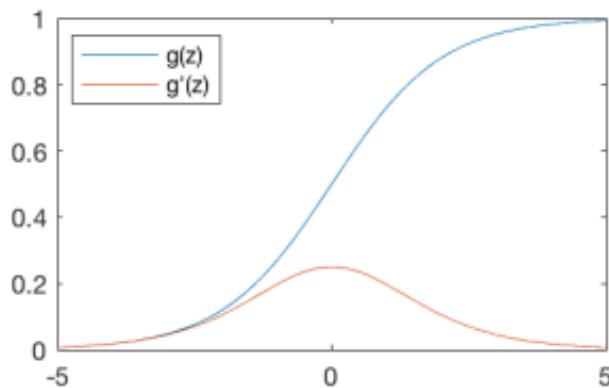
- Example: sigmoid function

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



Common Activation Functions

Sigmoid Function

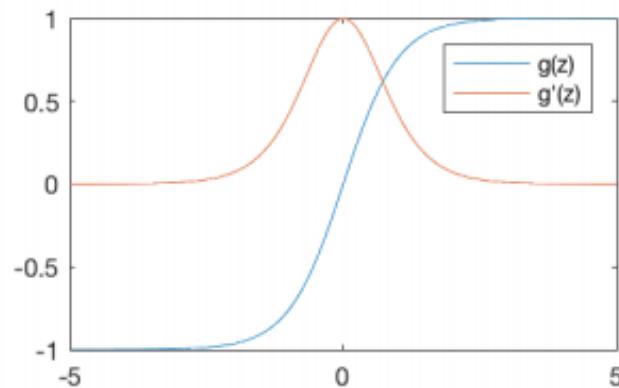


$$g(z) = \frac{1}{1 + e^{-z}}$$

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

 `tf.nn.sigmoid(z)`

Hyperbolic Tangent

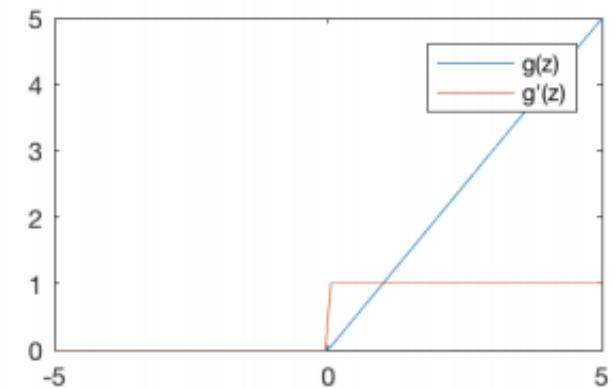


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

 `tf.nn.tanh(z)`

Rectified Linear Unit (ReLU)



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

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

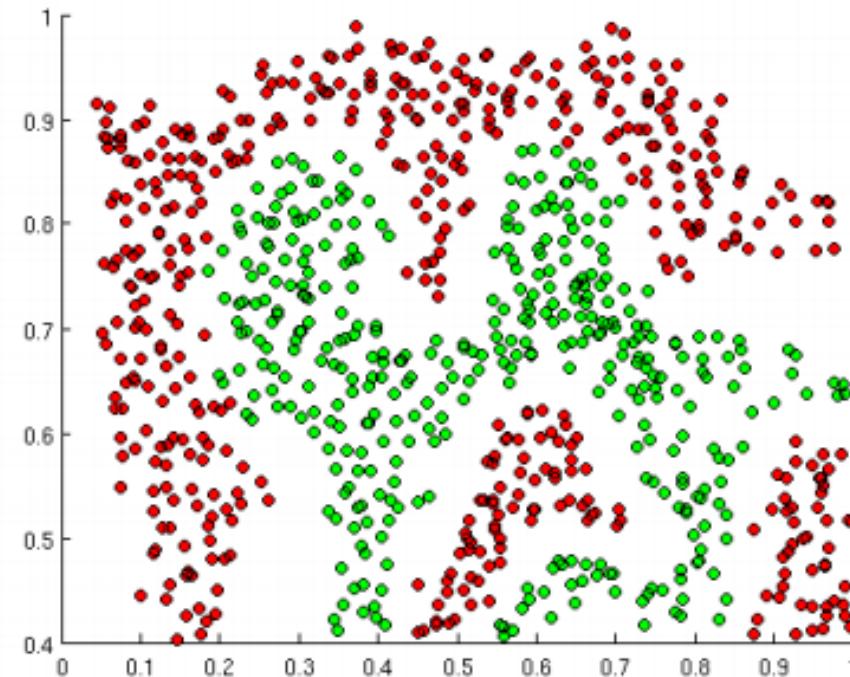
 `tf.nn.relu(z)`

NOTE: All activation functions are non-linear

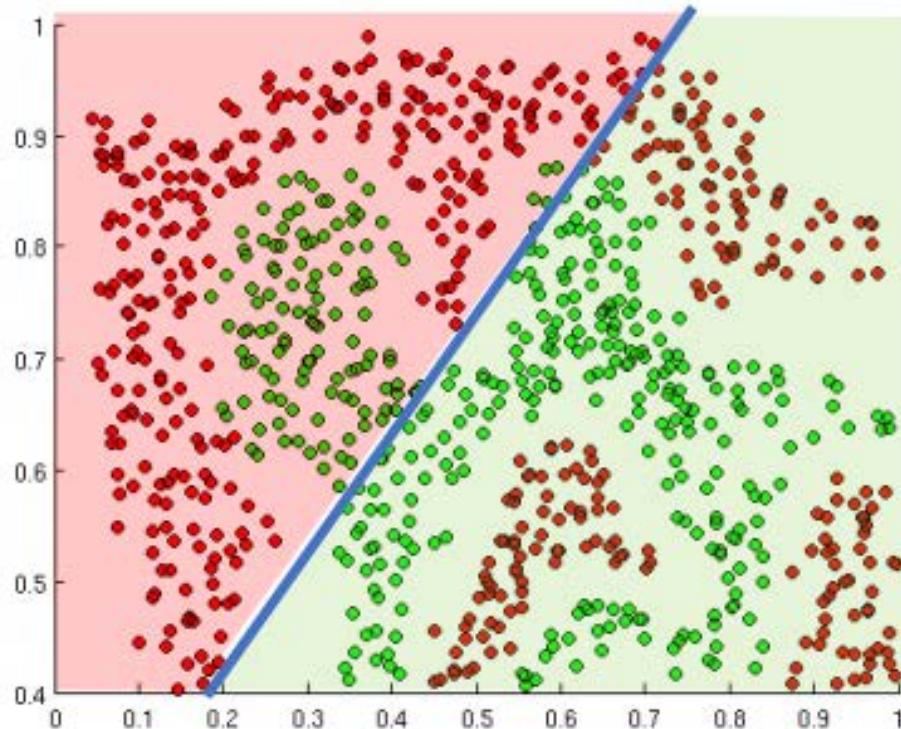
Why do we need activation functions?

- To introduce non-linearities into the network

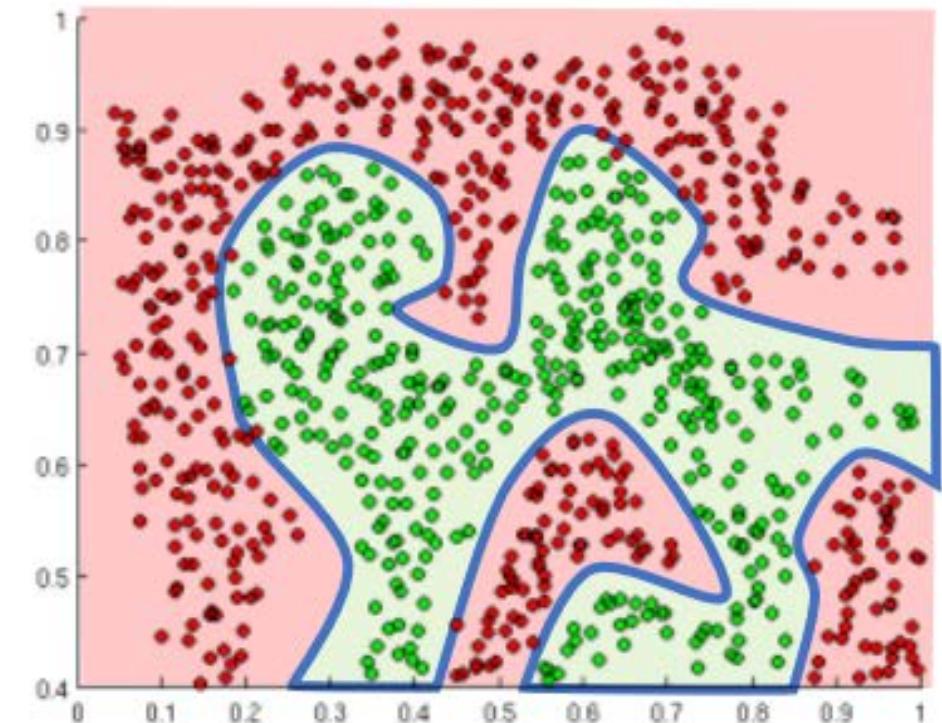
How to build a neural network to distinguish red and green points?



Linear vs Non-linear activation function

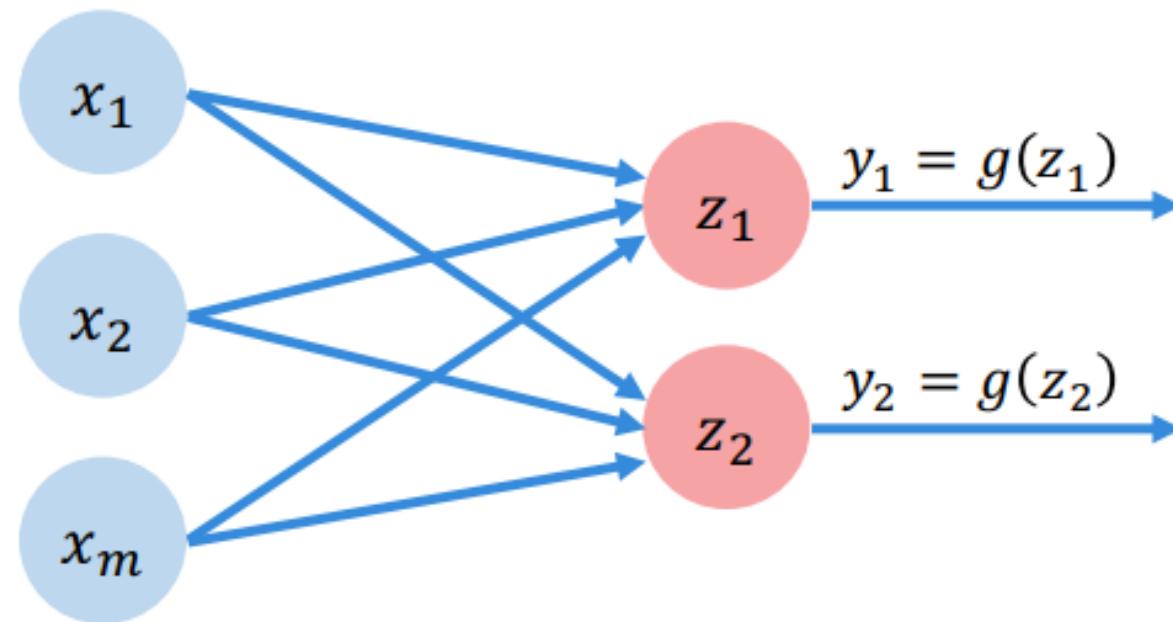


Linear activations produce linear decisions no matter the network size



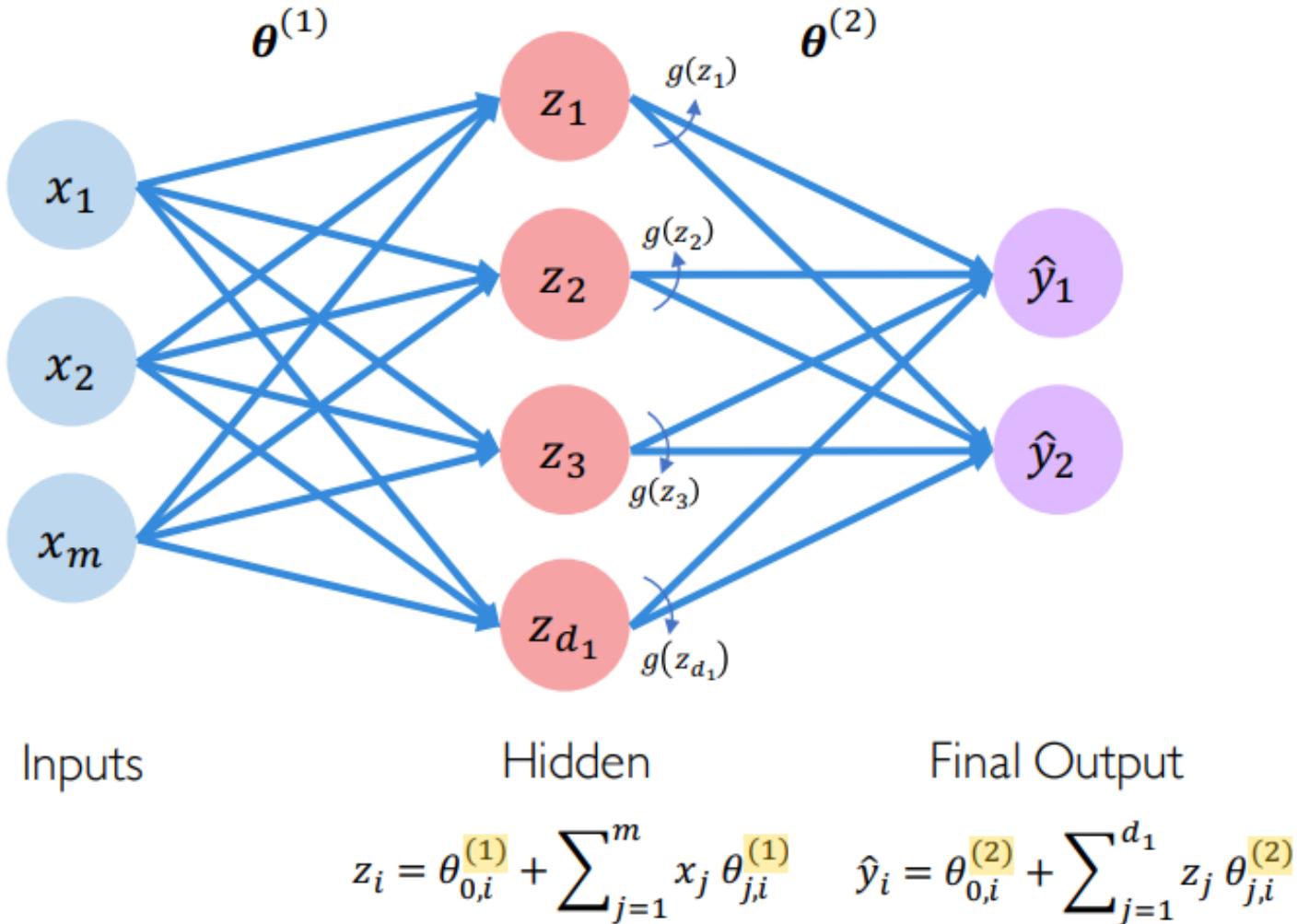
Non-linearities allow us to approximate arbitrarily complex functions

Multi-output perceptron

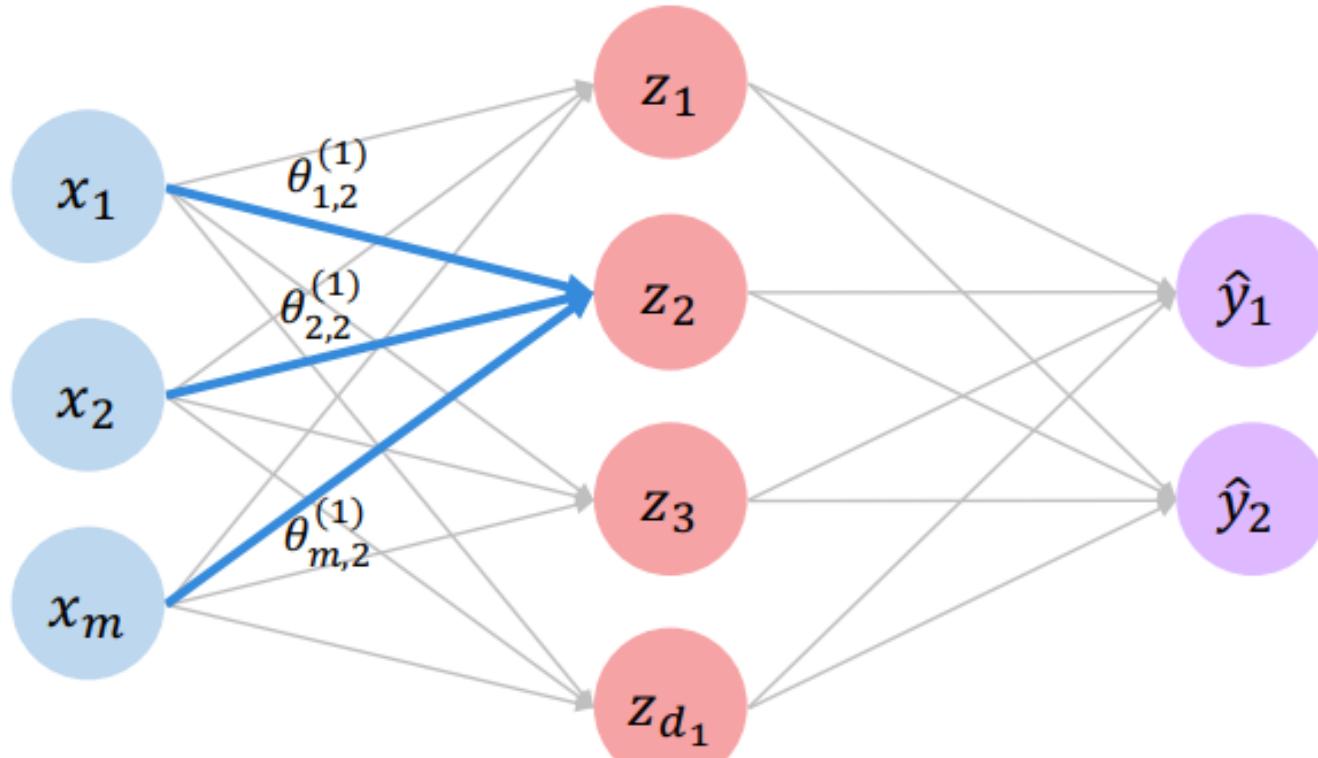


$$z_i = \theta_{0,i} + \sum_{j=1}^m x_j \theta_{j,i}$$

Single hidden layer neural network

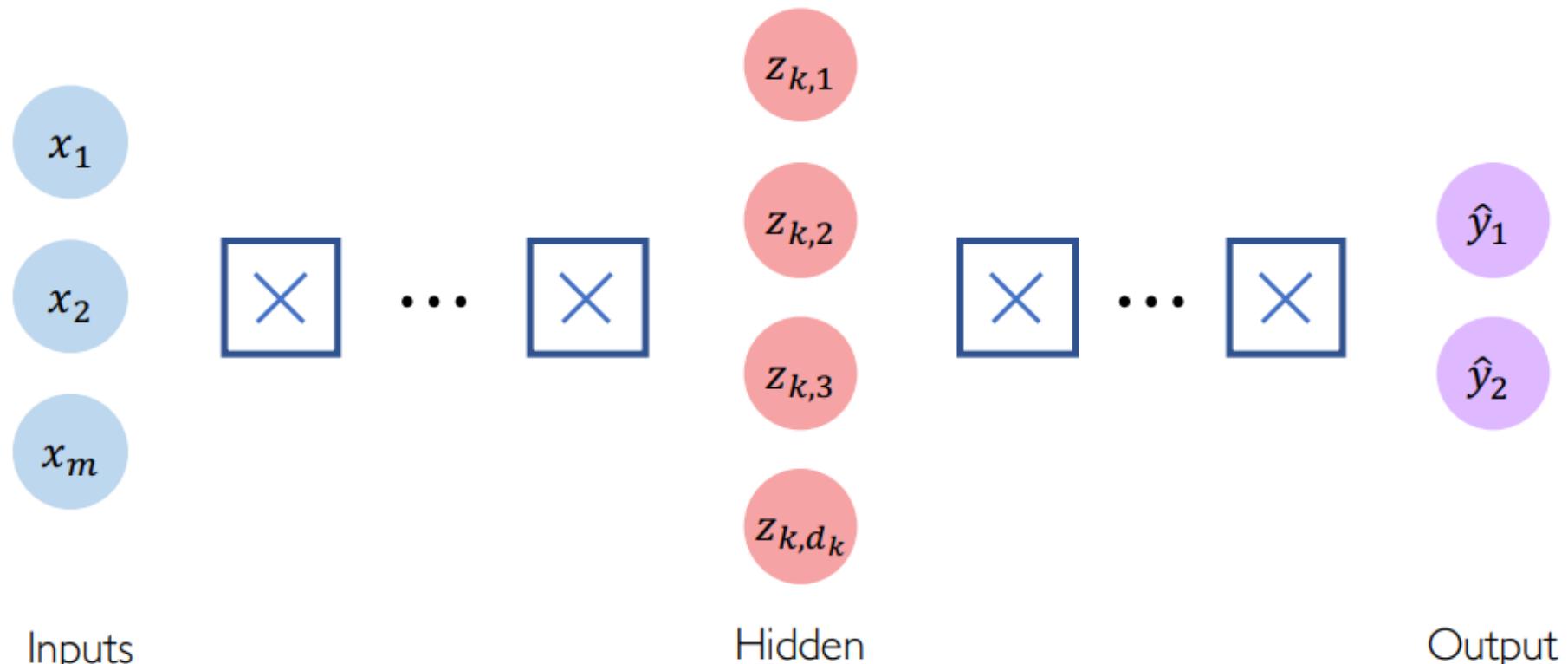


Single hidden layer neural network



$$\begin{aligned} z_2 &= \theta_{0,2}^{(1)} + \sum_{j=1}^m x_j \theta_{j,2}^{(1)} \\ &= \theta_{0,2}^{(1)} + x_1 \theta_{1,2}^{(1)} + x_2 \theta_{2,2}^{(1)} + x_m \theta_{m,2}^{(1)} \end{aligned}$$

Deep Neural Network



$$z_{k,i} = \theta_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) \theta_{j,i}^{(k)}$$

Example Problem

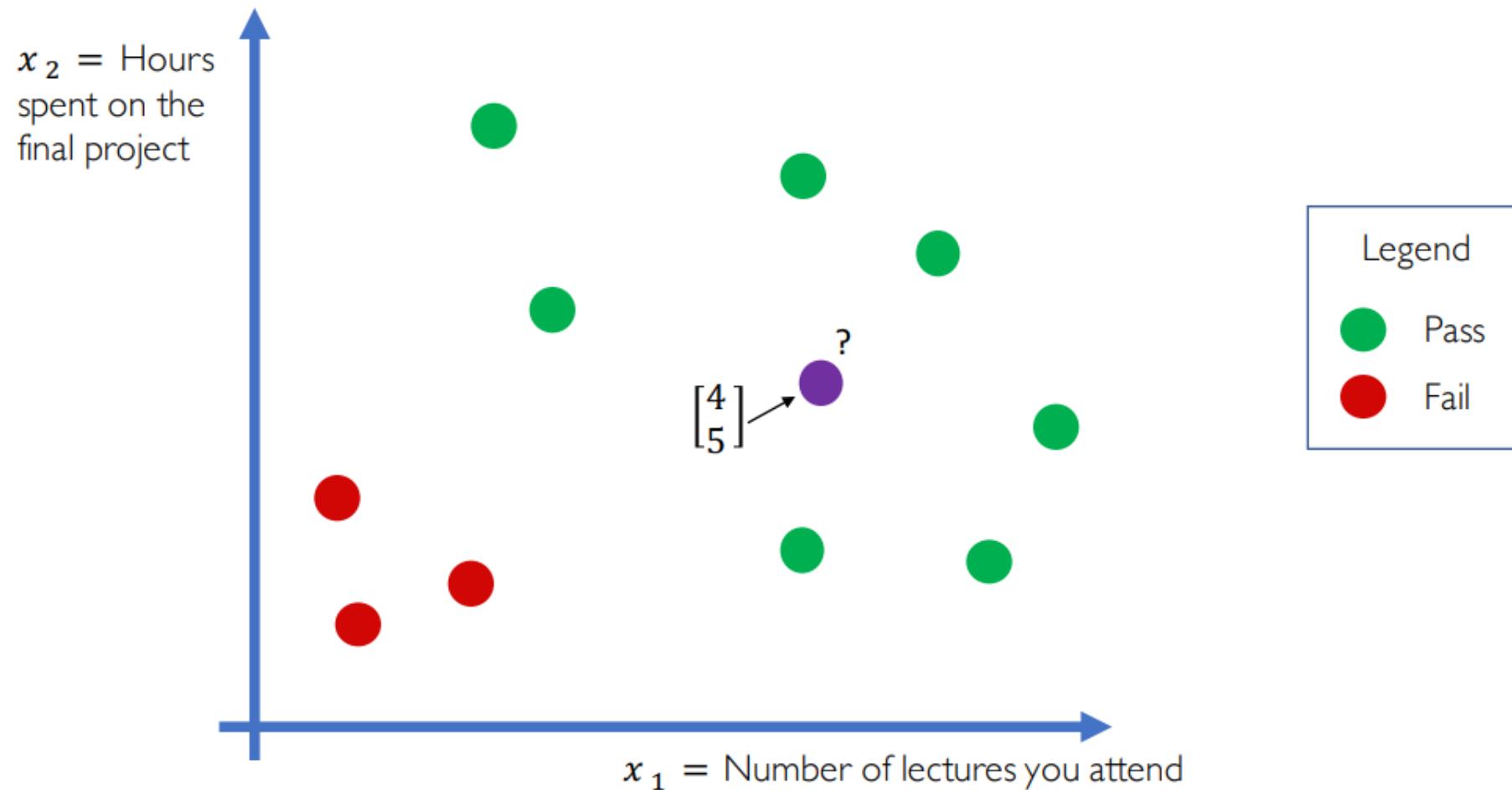
Will I pass this class?

Let's start with a simple two feature model

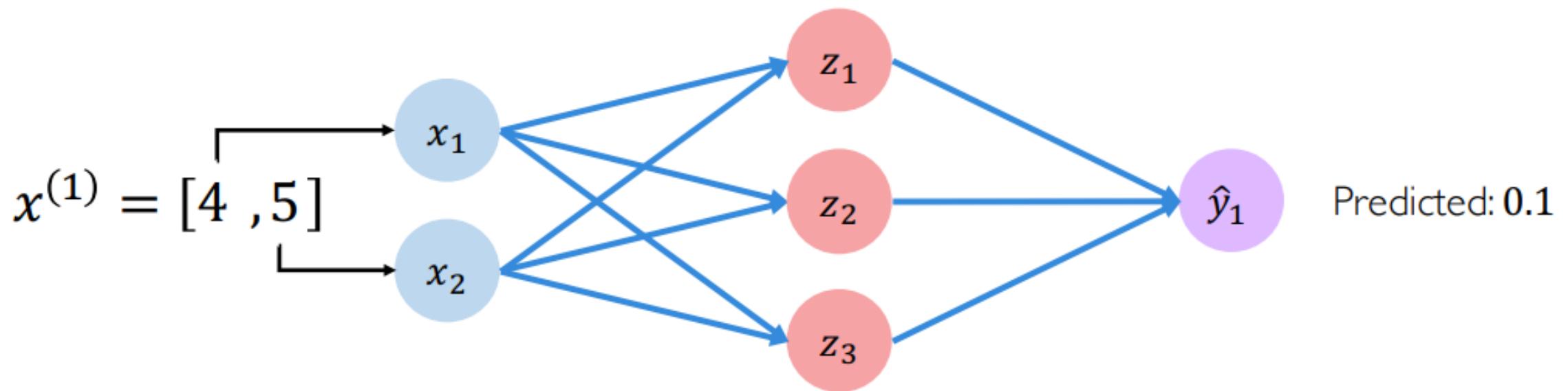
x_1 = Number of lectures you attend

x_2 = Hours spent on the final project

Example Problem

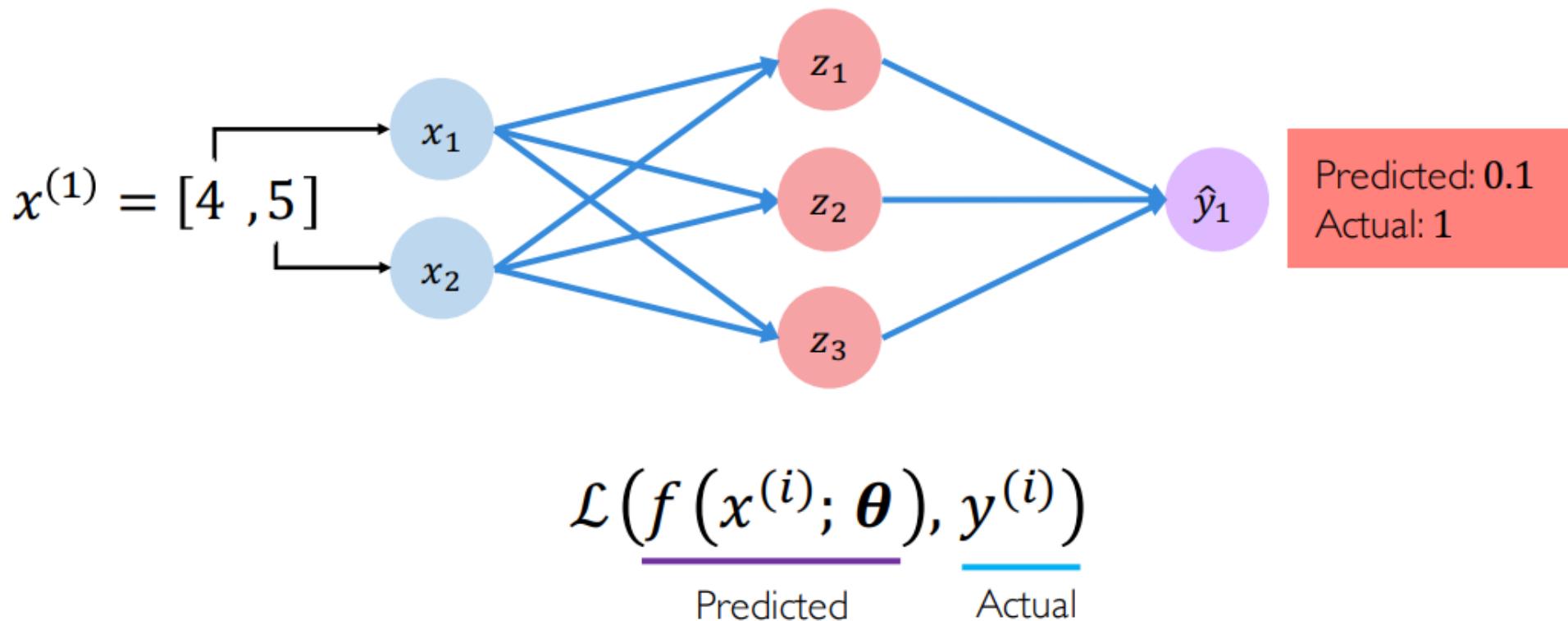


Example Problem



Quantifying loss

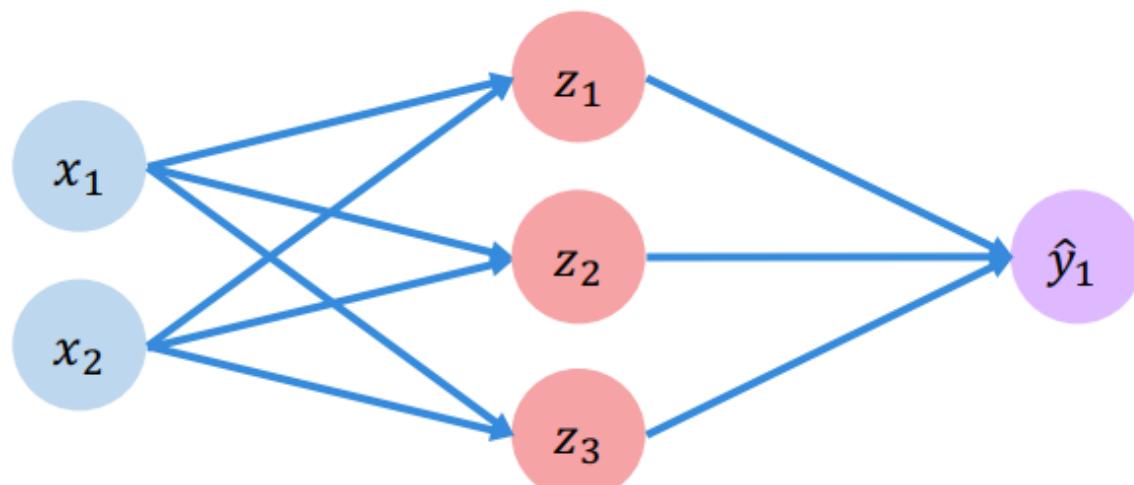
The *loss* of our network measures the cost incurred from incorrect predictions



Empirical Loss

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

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



$f(x)$	y
0.1	✗
0.8	✗
0.6	✓
\vdots	\vdots

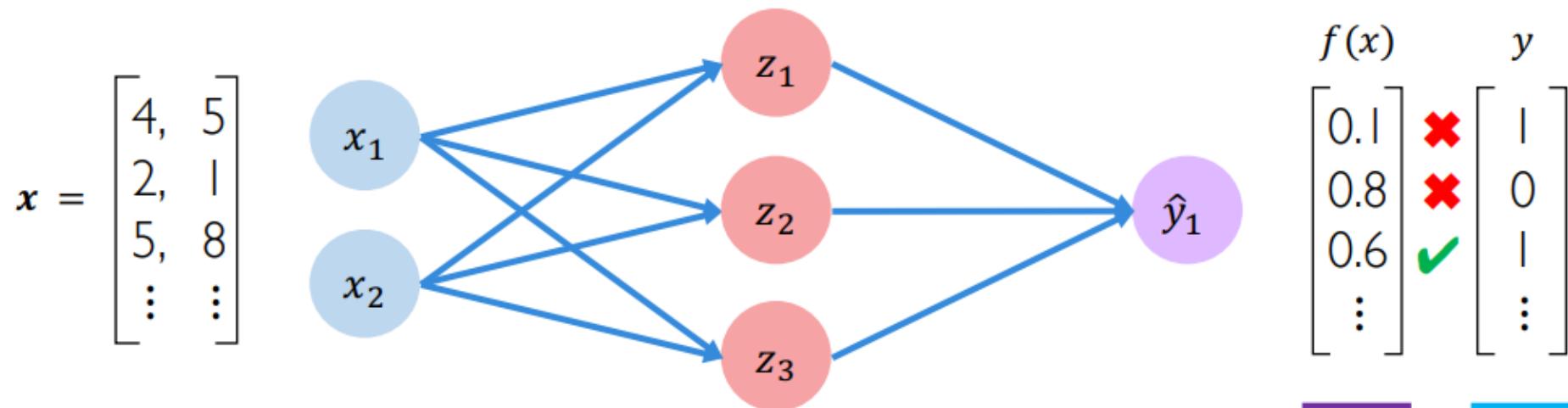
- Also known as:
- Objective function
 - Cost function
 - Empirical Risk

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

Predicted Actual

Binary Cross Entropy Loss

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



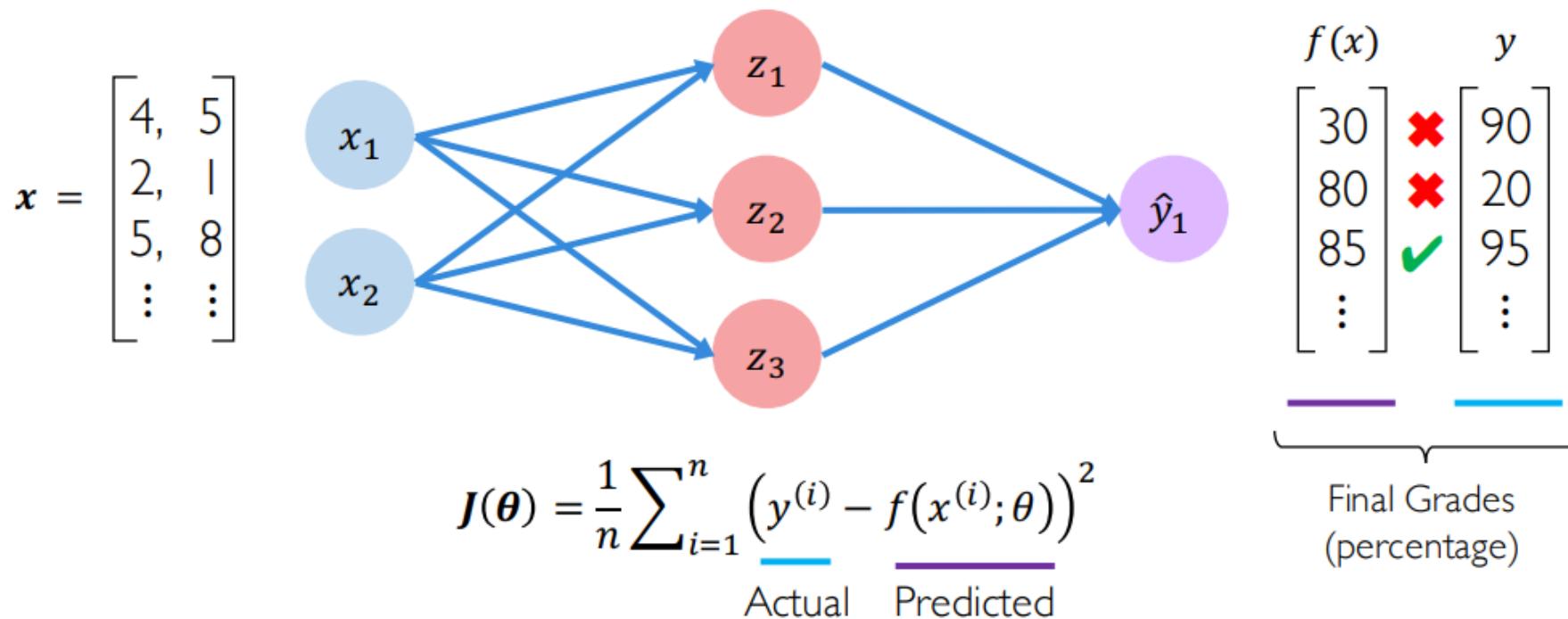
$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)} \log(f(x^{(i)}; \theta))}_{\text{Actual}} + \underbrace{(1 - y^{(i)}) \log(1 - f(x^{(i)}; \theta))}_{\text{Predicted}}$$



```
loss = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(model.y, model.pred))
```

Mean Squared Error Loss

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



```
loss = tf.reduce_mean( tf.square(tf.subtract(model.y, model.pred)) )
```

Loss Optimization

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

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \boldsymbol{\theta}), y^{(i)})$$

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

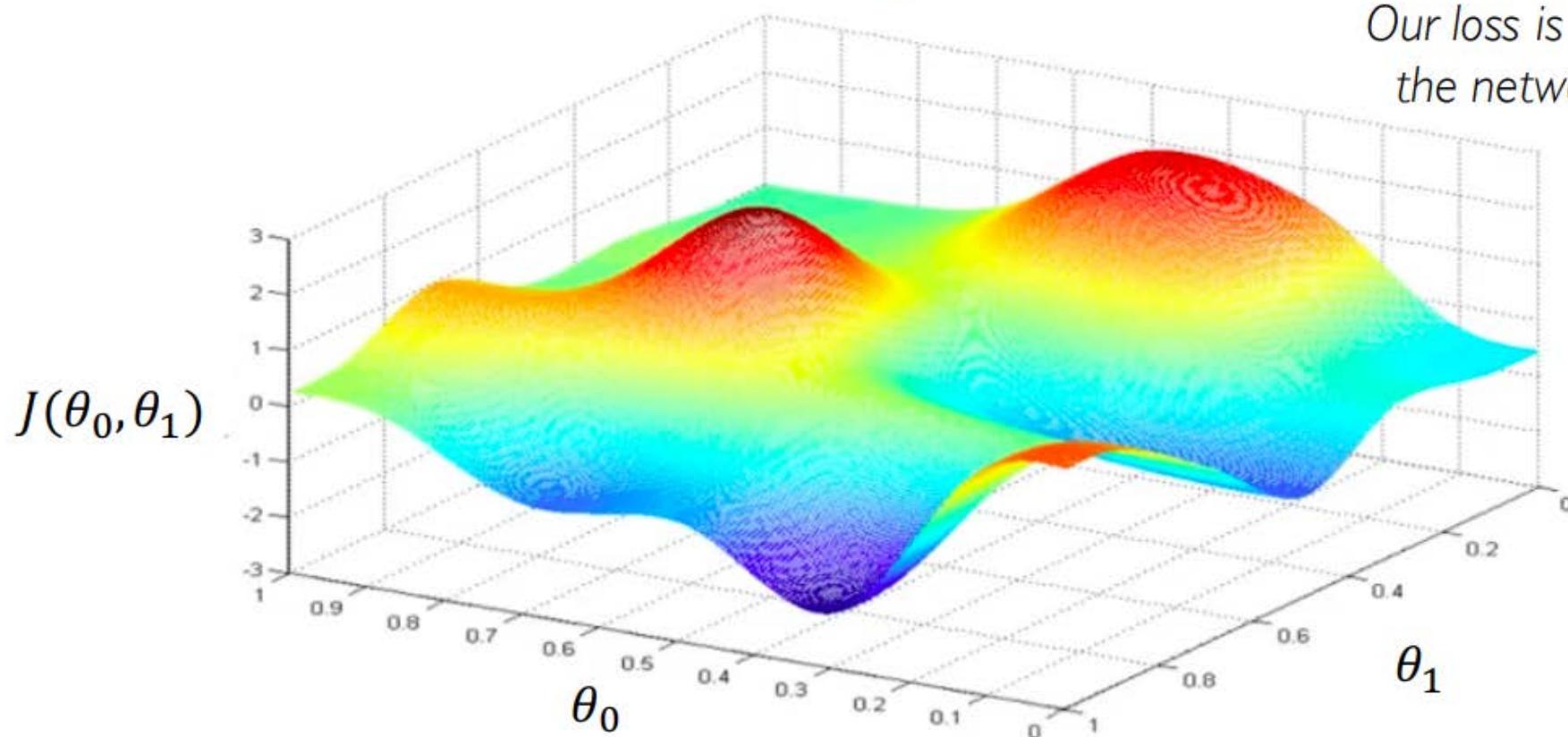


Remember:

$$\boldsymbol{\theta} = \{\theta^{(0)}, \theta^{(1)}, \dots\}$$

Loss Optimization

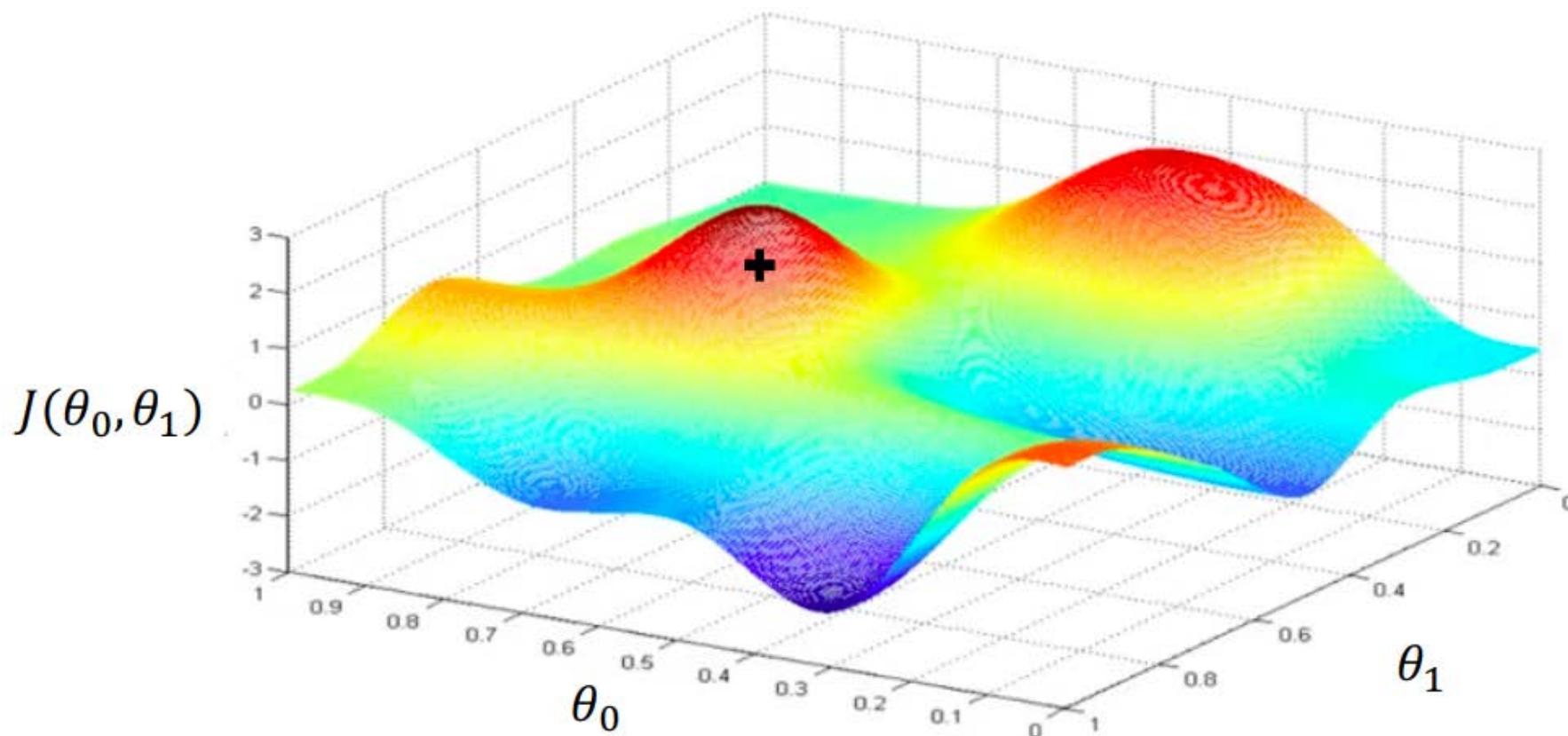
$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$$



Remember:
Our loss is a function of
the network weights!

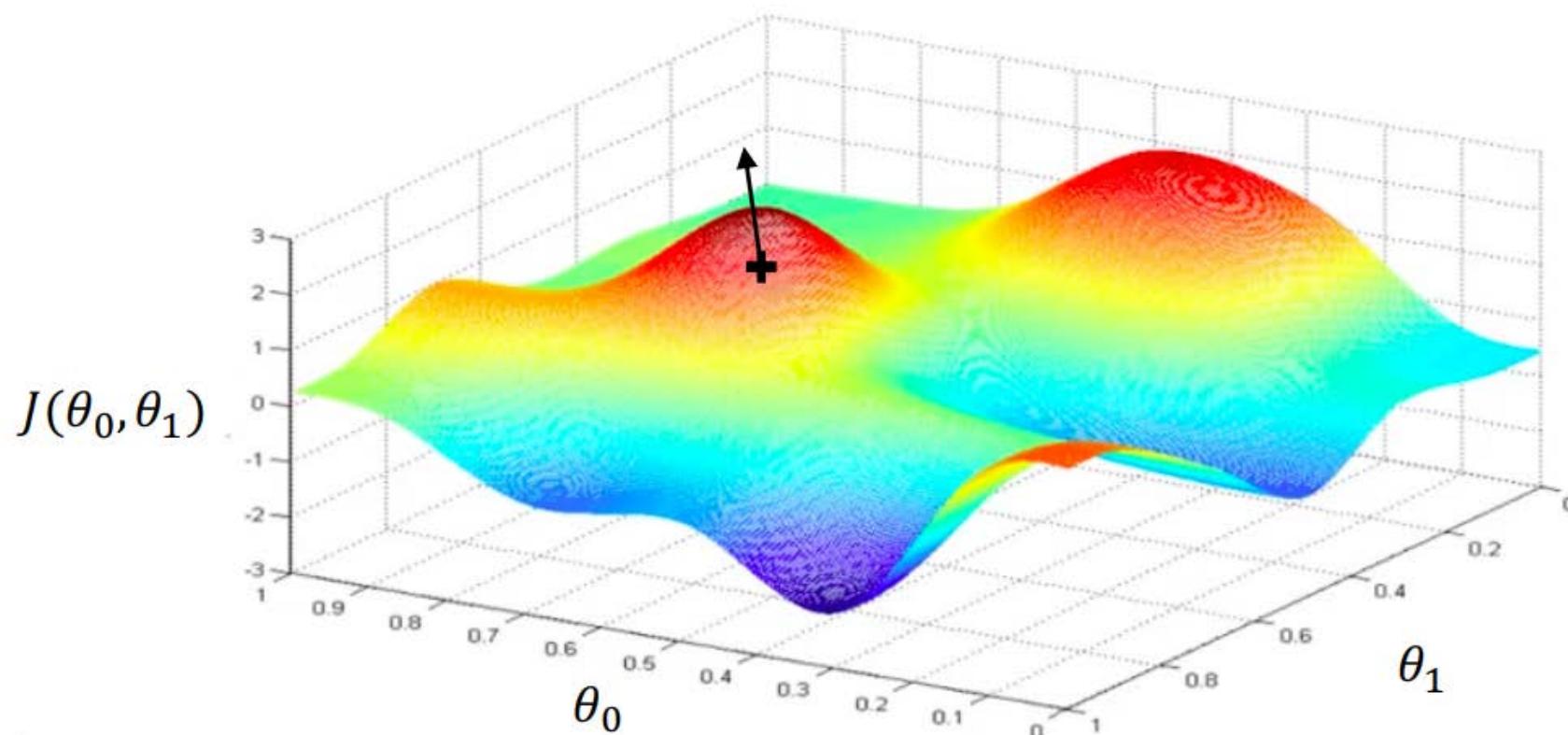
Loss Optimization

Randomly pick an initial (θ_0, θ_1)



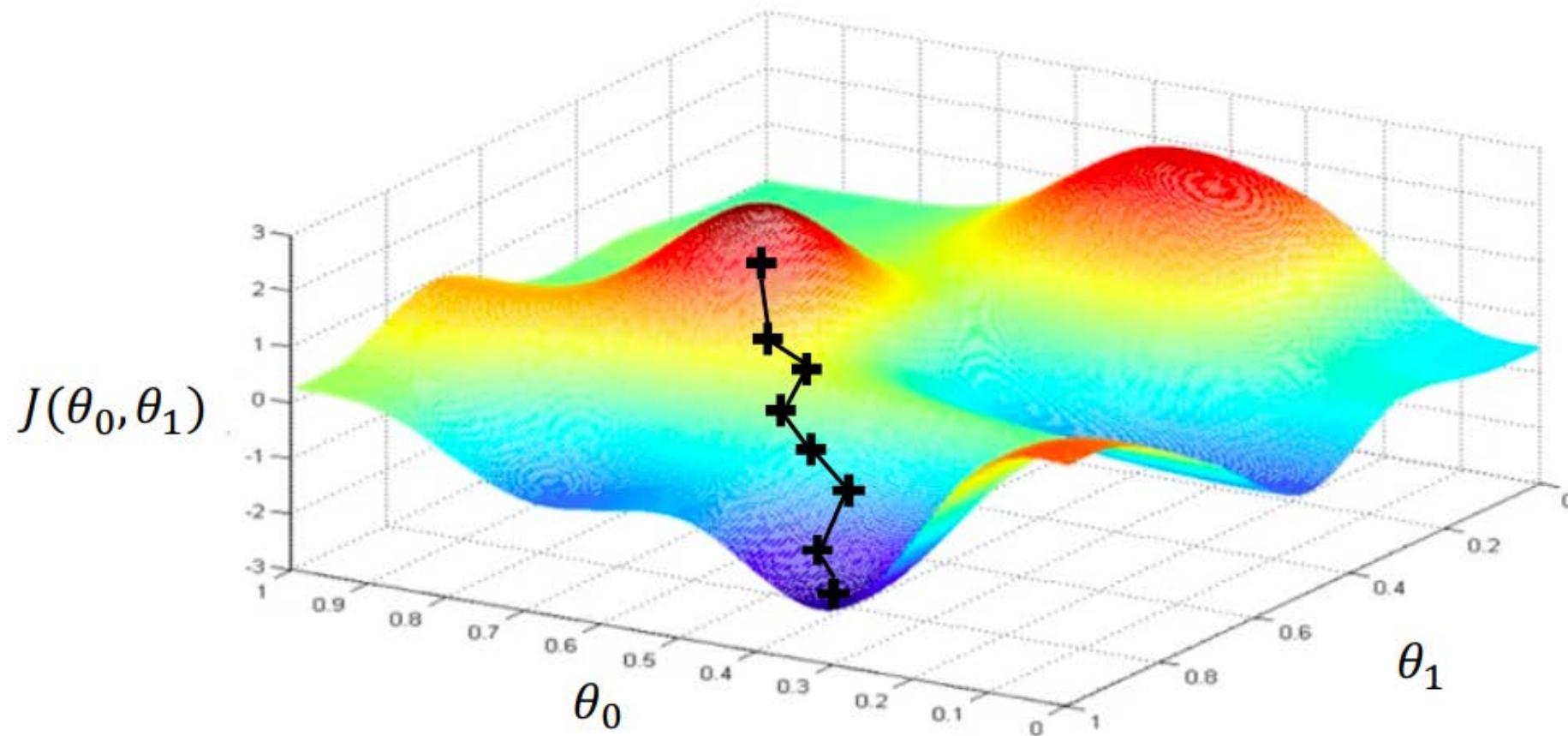
Loss Optimization

Compute gradient, $\frac{\partial J(\theta)}{\partial \theta}$



Gradient Descent

Repeat until convergence



Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\theta)}{\partial \theta}$
4. Update weights, $\theta \leftarrow \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$
5. Return weights

```
 weights = tf.random_normal(shape, stddev=sigma)
```

```
 grads = tf.gradients(ys=loss, xs=weights)
```

```
 weights_new = weights.assign(weights - lr * grads)
```

Gradient Descent

Algorithm

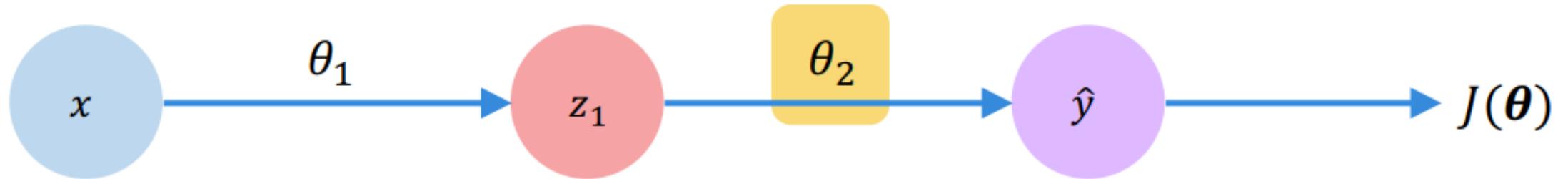
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\theta)}{\partial \theta}$
4. Update weights, $\theta \leftarrow \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$
5. Return weights

```
 weights = tf.random_normal(shape, stddev=sigma)
```

```
 grads = tf.gradients(ys=loss, xs=weights)
```

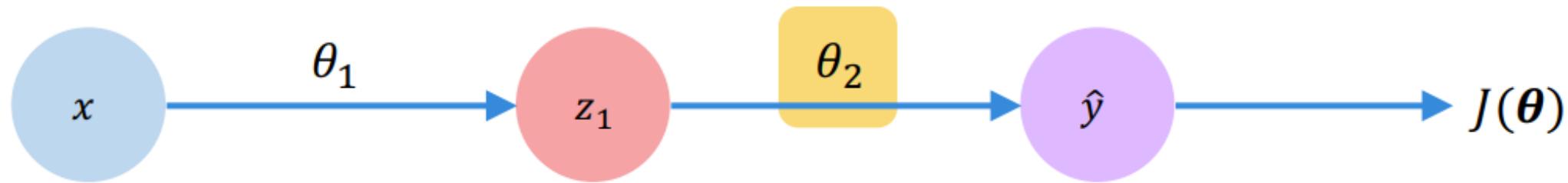
```
 weights_new = weights.assign(weights - lr * grads)
```

Computing Gradients: Backpropagation



How does a small change in one weight (ex. θ_2) affect the final loss $J(\boldsymbol{\theta})$?

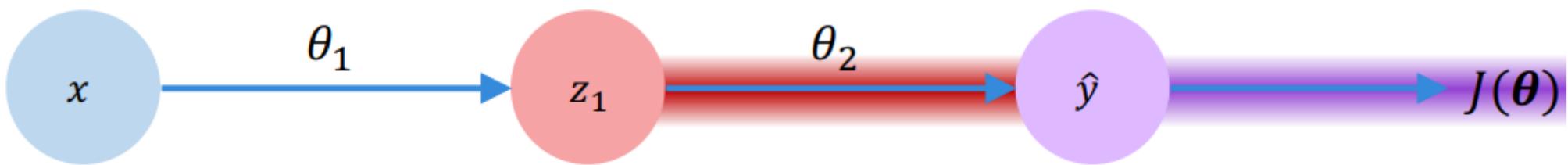
Computing Gradients: Backpropagation



$$\frac{\partial J(\theta)}{\partial \theta_2} =$$

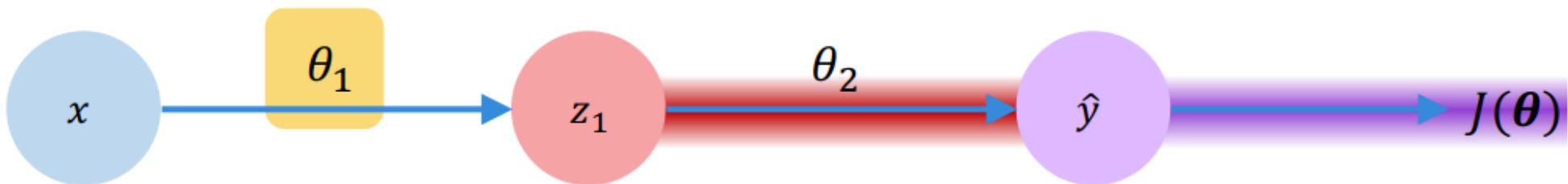
Let's use the chain rule!

Computing Gradients: Backpropagation



$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_2} = \underline{\frac{\partial J(\boldsymbol{\theta})}{\partial \hat{y}}} * \overline{\frac{\partial \hat{y}}{\partial \theta_2}}$$

Computing Gradients: Backpropagation

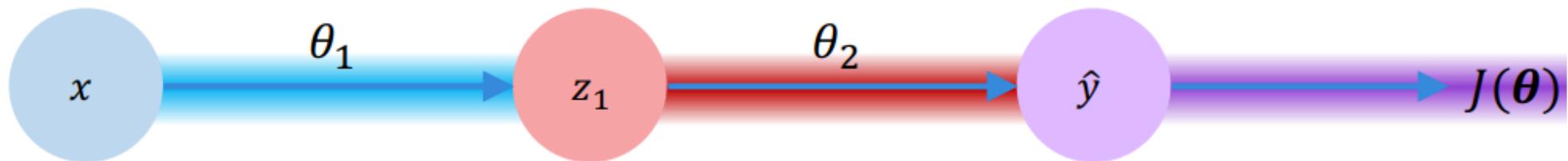


$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_1} = \frac{\partial J(\boldsymbol{\theta})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial \theta_1}$$

Apply chain rule!

Apply chain rule!

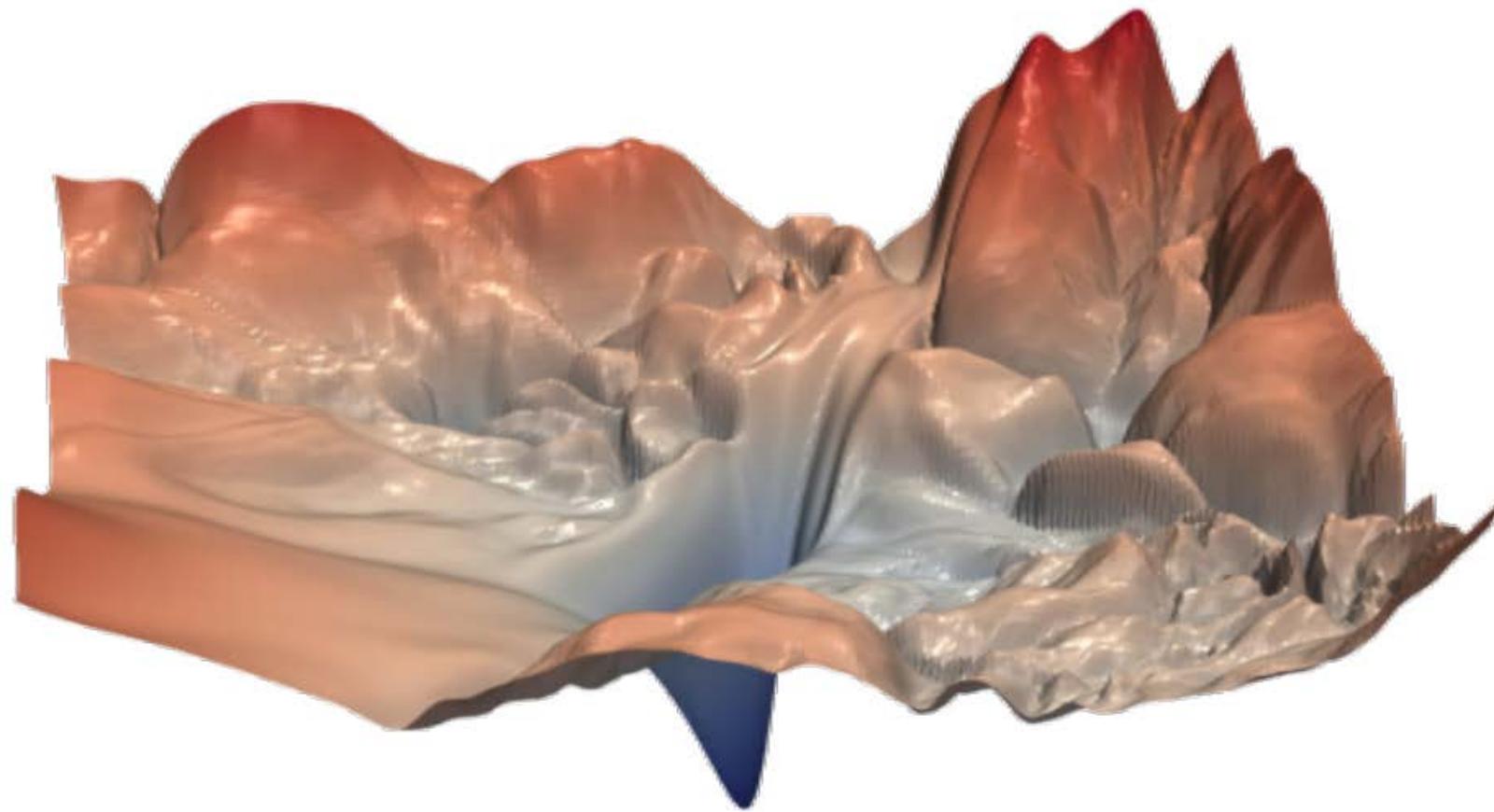
Computing Gradients: Backpropagation



$$\frac{\partial J(\theta)}{\partial \theta_1} = \underbrace{\frac{\partial J(\theta)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial \theta_1}}_{\text{blue}}$$

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

Training Neural Networks is Difficult



Loss functions can be difficult to optimize

Remember:

Optimization through gradient descent

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



How can we set the
learning rate?

Setting the learning rate

Small learning rates

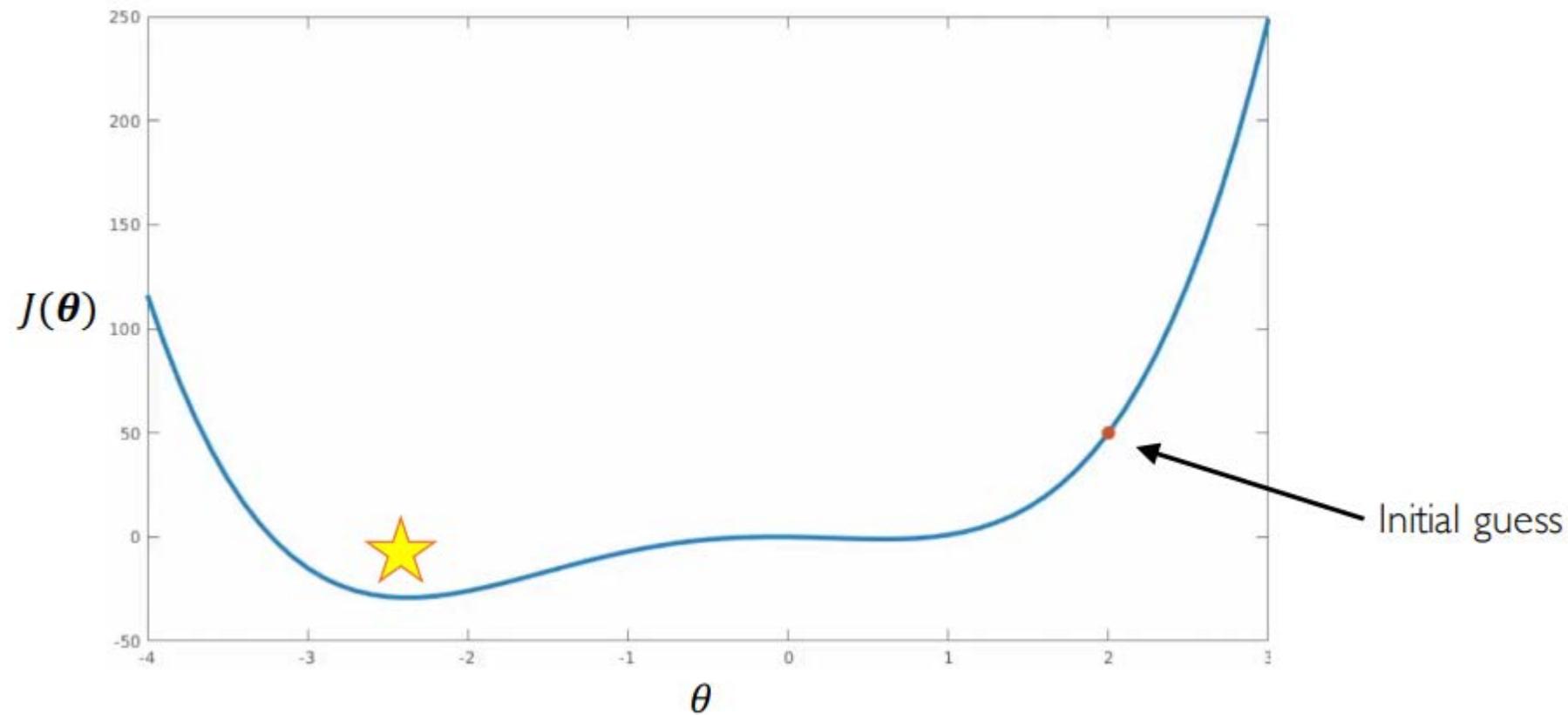
converges slowly and gets stuck in false local minima

Large learning rates

overshoot, become unstable and diverge

Stable learning rates

converge smoothly and avoid local minima



Adaptive Learning Rates

- Design an adaptive learning rate that adapts to the landscape
- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
 - How large the gradient is
 - How fast learning is happening
 - Etc..

Adaptive Learning Rate Algorithms

- Momentum
- Adagrad
- Adadelta
- Adam
- RMSProp



`tf.train.MomentumOptimizer`



`tf.train.AdagradOptimizer`



`tf.train.AdadeltaOptimizer`



`tf.train.AdamOptimizer`



`tf.train.RMSPropOptimizer`

Qian et al. "On the momentum term in gradient descent learning algorithms." 1999.

Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011.

Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012.

Kingma et al. "Adam: A Method for Stochastic Optimization." 2014.

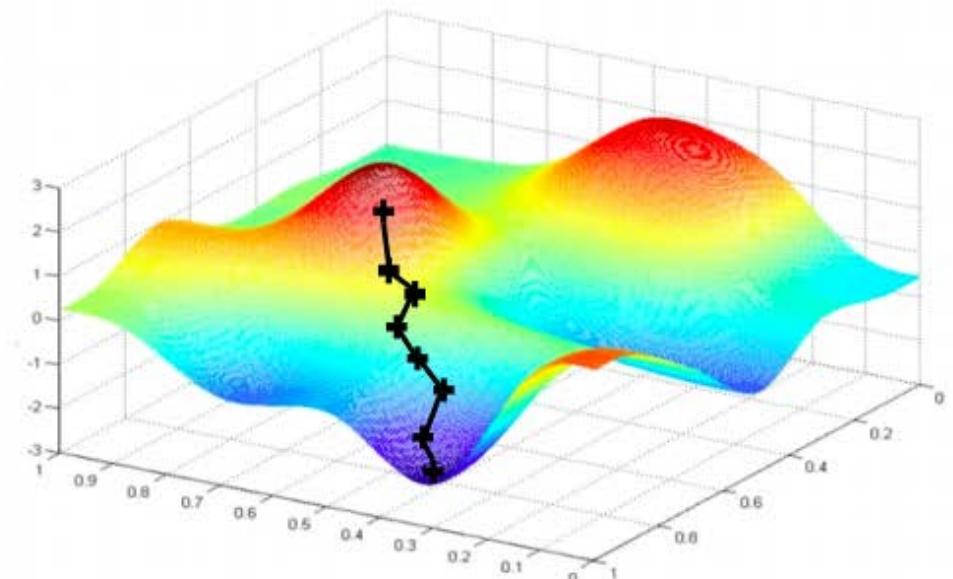
Hinton's Coursera lecture (unpublished)

Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\theta)}{\partial \theta}$
4. Update weights, $\theta \leftarrow \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$
5. Return weights

Can be very
computational to
compute!

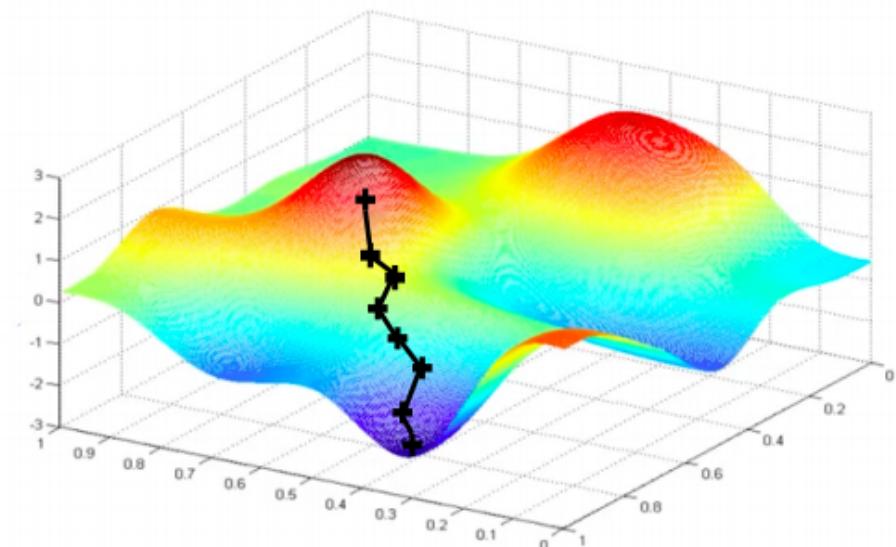


Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\theta)}{\partial \theta}$
5. Update weights, $\theta \leftarrow \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$
6. Return weights

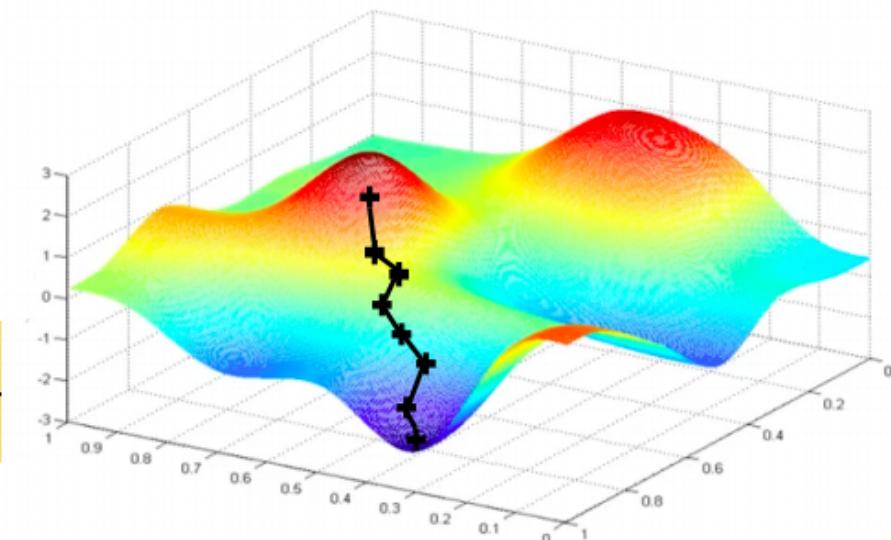
Easy to compute but
very noisy
(stochastic)!



Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient,
$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\theta)}{\partial \theta}$$
5. Update weights, $\theta \leftarrow \theta - \eta \frac{\partial J(\theta)}{\partial \theta}$
6. Return weights



Fast to compute and a much better estimate of the true gradient!

Mini-batches

- More accurate estimation of gradient
 - Smoother convergence
 - Allows for larger learning rates
- Mini-batches lead to fast training
 - Can parallelize computation + achieve significant speed increases on GPU's

Terminology

- **Number of iterations:** The number of times the gradient is estimated and the parameters of the neural network are updated using a batch of training instances
- **Batch size:** Number of training instances used in one iteration
- **Mini-batch:** When the total number of training instances N is large, a small number of training instances $B \ll N$ which constitute a mini-batch can be used in one iteration to estimate the gradient of the loss function and update the parameters of the network
- **Epoch:** It takes $n = N/B$ iterations to use the entire training data once. That is called an epoch. The total number of times the parameters get updates is $(N/B)*E$, where E is the number of epochs.

Three modes of gradient descent

- Batch mode: $N=B$, one epoch is same as one iteration.
- Mini-batch mode: $1 < B < N$, one epoch consists of N/B iterations.
- Stochastic mode: $B=1$, one epoch takes N iterations.

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: always works perfectly on training data

Your Dataset

Idea #2: Split data into **train** and **test**, choose hyperparameters that work best on test data

BAD: No idea how algorithm will perform on new data

train

test

Idea #3: Split data into **train**, **val**, and **test**; choose hyperparameters on val and evaluate on test

Better!

train

validation

test

Setting Hyperparameters

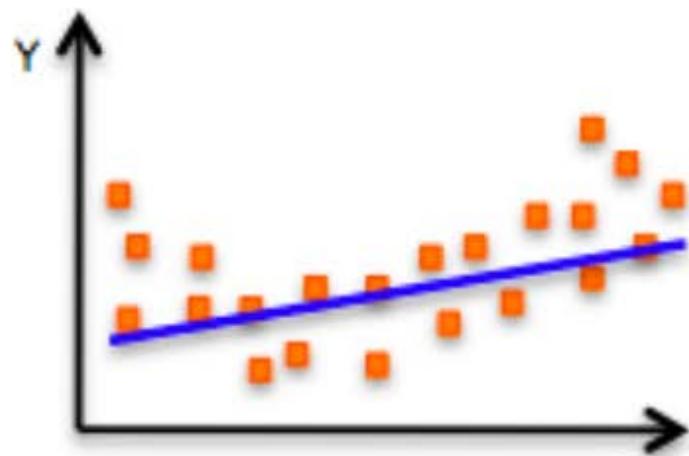
Your Dataset

Idea #4: Cross-Validation: Split data into **folds**,
try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test

Useful for small datasets, but not used too frequently in deep learning

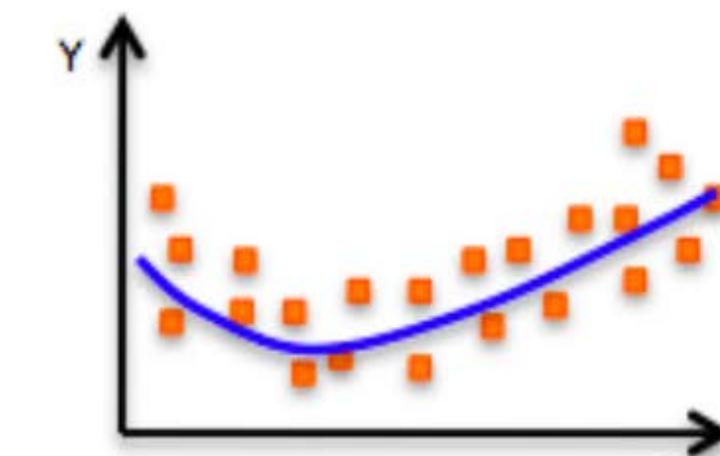
The Problem of Overfitting



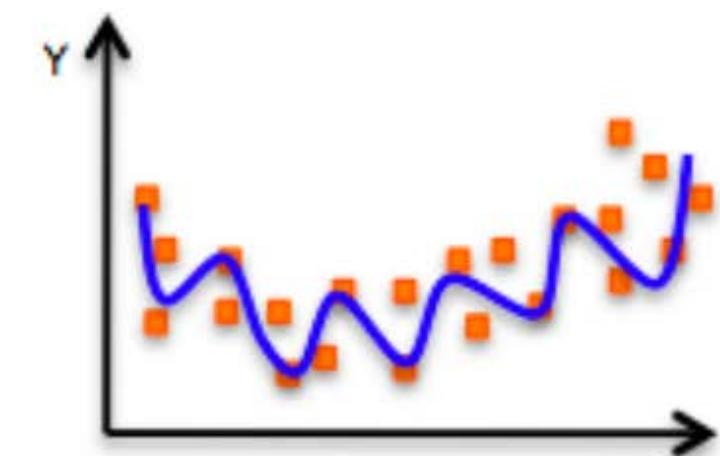
Underfitting

Model does not have capacity
to fully learn the data

High bias



Ideal fit



Overfitting

Too complex, extra parameters,
does not generalize well

High variance

High Bias vs High Variance

- High Bias (high training set error)
 - Use a bigger network
 - Try different optimization algorithms
 - Train longer
 - Try different architecture
- High Variance (high validation set error)
 - Collect more data
 - Use regularization
 - Try different NN architecture

Regularization

- What is it?
 - Technique that constrains our optimization problem to discourage complex models
- Why do we need it?
 - Improve generalization of our model on unseen data

Regularization 1: Penalizing weights

- Penalize large weights using penalties or constraints on their squared values (L2 penalty) or absolute values (L1 penalty)
- Neural networks have thousands (or millions of parameters)
 - Danger of overfitting

$$\theta^* \leftarrow \arg \min_{\theta} \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L(x; \theta_1, \dots, L)) + \lambda \Omega(\theta)$$

Regularization 1: L1 and L2 regularization

- L2 regularization (most popular)

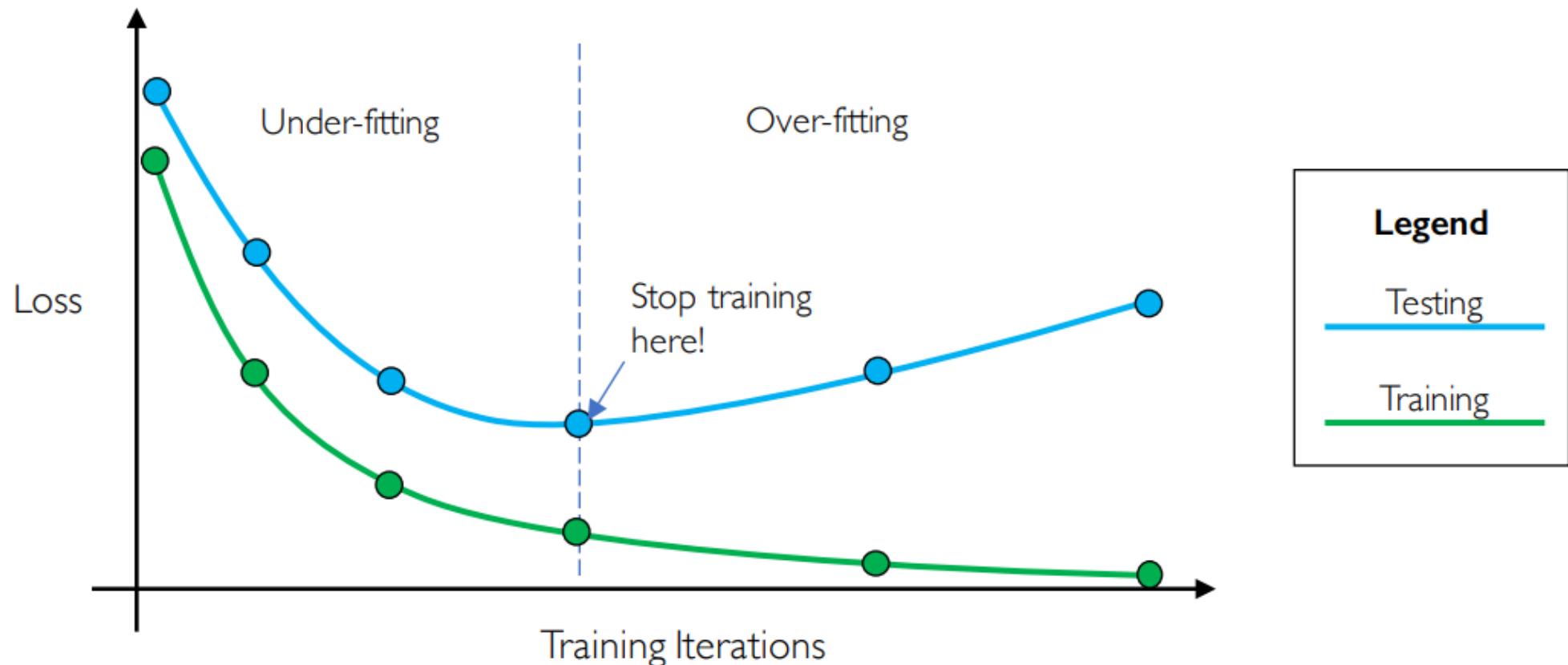
$$\theta^* \leftarrow \arg \min_{\theta} \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L(x; \theta_1, \dots, \theta_L)) + \frac{\lambda}{2} \sum_l \|\theta_l\|^2$$

- L1 regularization

$$\theta^* \leftarrow \arg \min_{\theta} \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L(x; \theta_1, \dots, \theta_L)) + \frac{\lambda}{2} \sum_l \|\theta_l\|$$

Regularization 2: Early Stopping

- Stop training before we have a chance to overfit

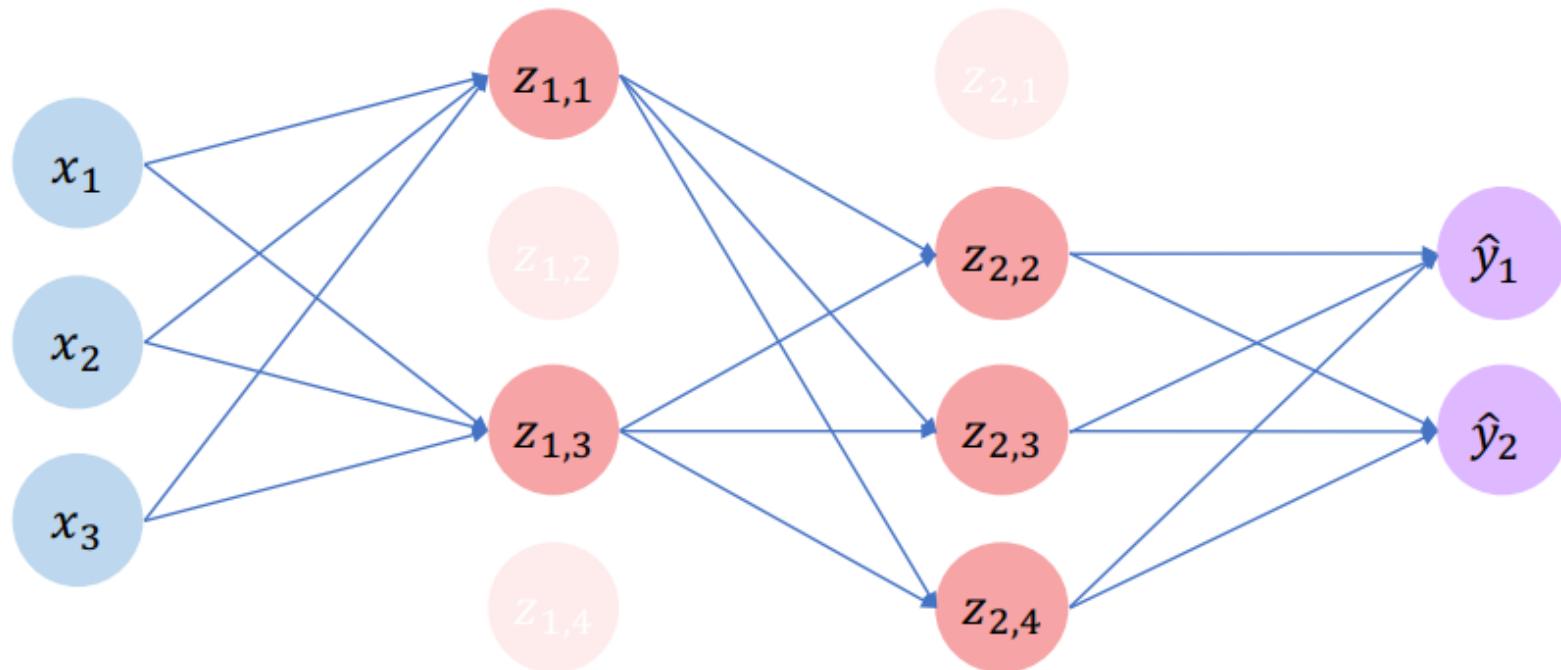


Regularization 3: Dropout

- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

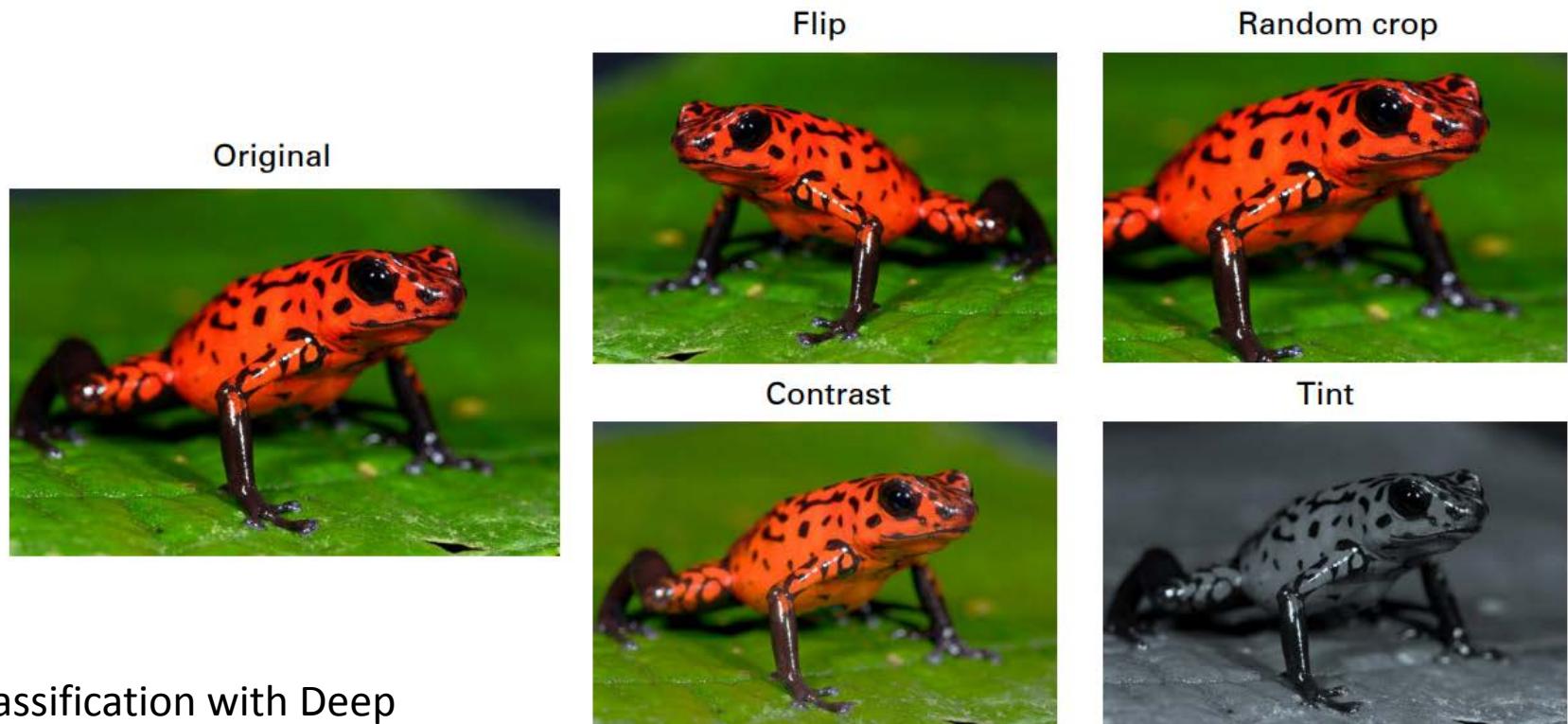


`tf.nn.dropout(hiddenLayer, p=0.5)`



Regularization 4: Data Augmentation

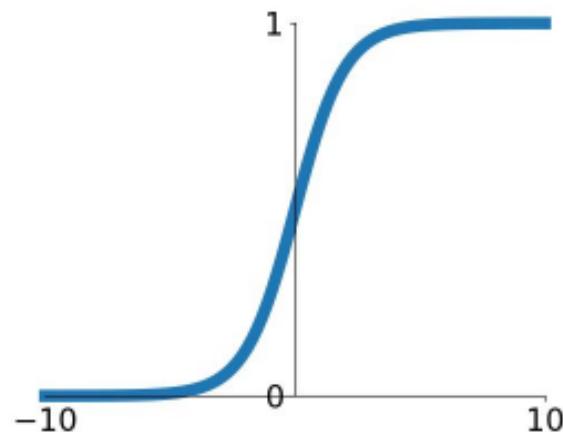
- Adding more data reduces overfitting
- Data collection and labelling is expensive
- Solution: Synthetically increase training dataset



Difference between Activation Functions

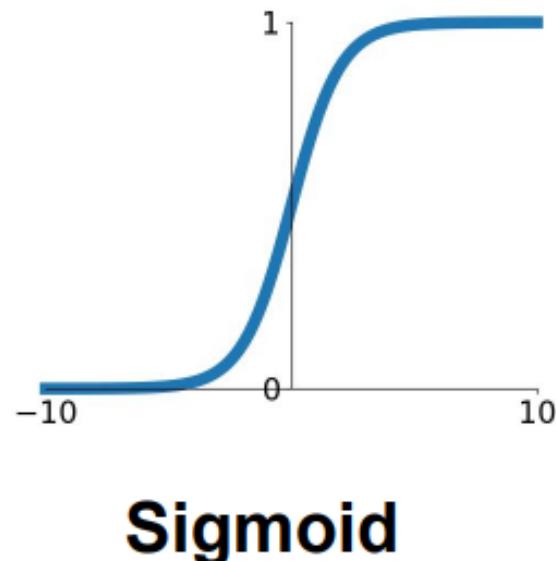
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron



Sigmoid

Difference between Activation Functions



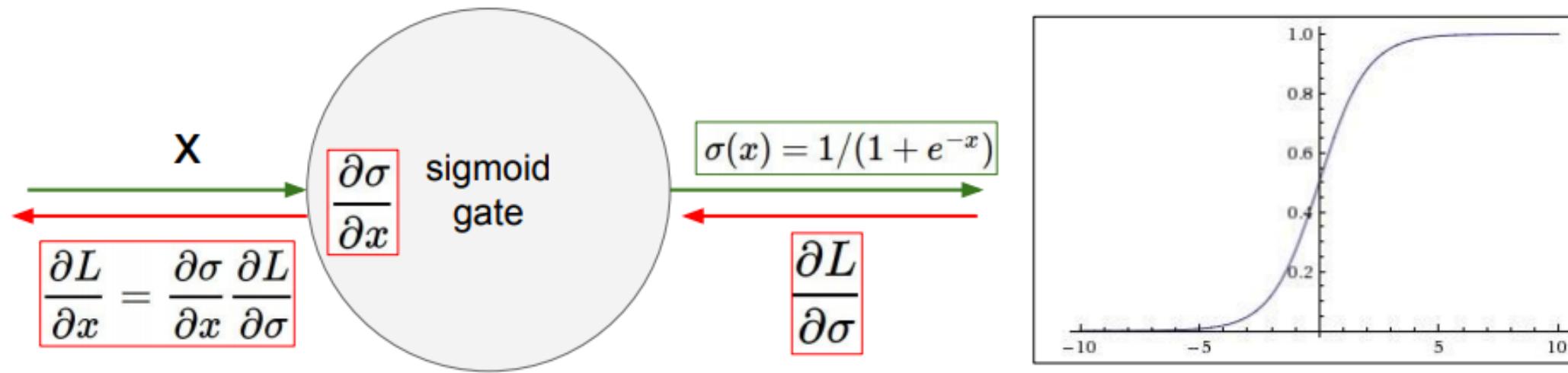
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0, 1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients

Difference between Activation Functions



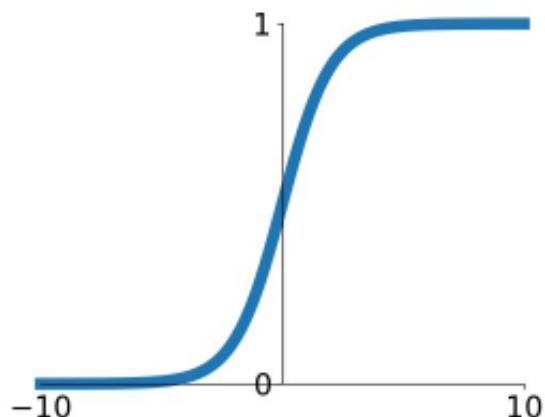
What happens when $x = -10$?

What happens when $x = 0$?

What happens when $x = 10$?

Difference between Activation Functions

$$\sigma(x) = 1/(1 + e^{-x})$$



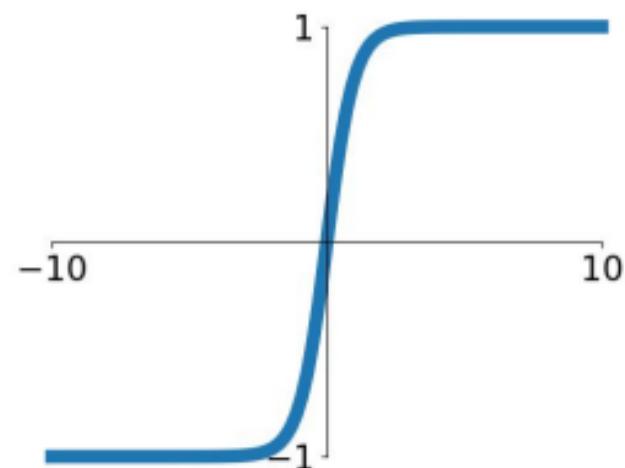
Sigmoid

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3. $\exp()$ is a bit compute expensive

Difference between Activation Functions

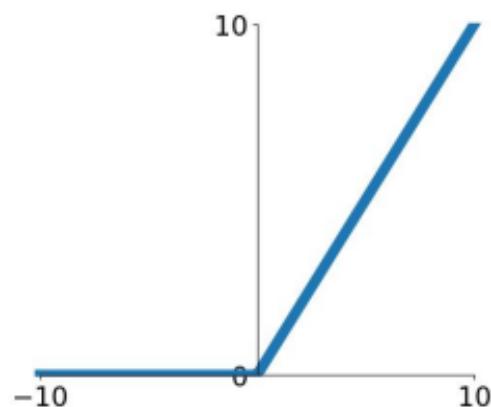


$\tanh(x)$

- Squashes numbers to range [-1,1]
- zero centered (nice)
- still kills gradients when saturated :(

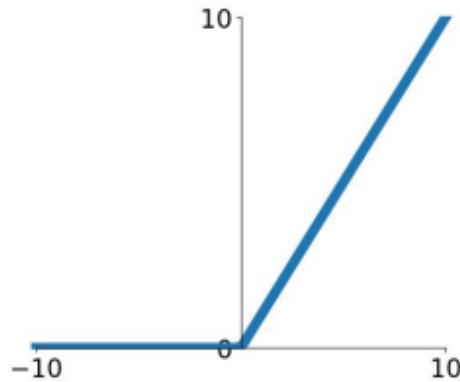
Difference between Activation Functions

- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid



ReLU
(Rectified Linear Unit)

Difference between Activation Functions



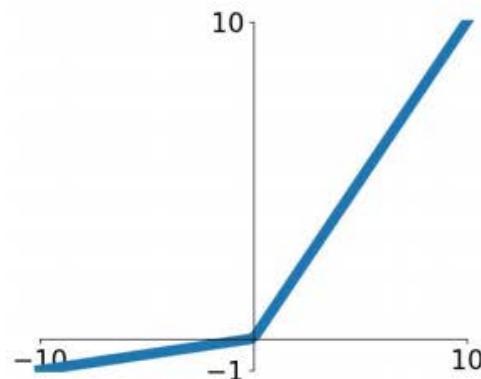
ReLU
(Rectified Linear Unit)

- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid

- Not zero-centered output
- An annoyance:

hint: what is the gradient when $x < 0$?

Difference between Activation Functions



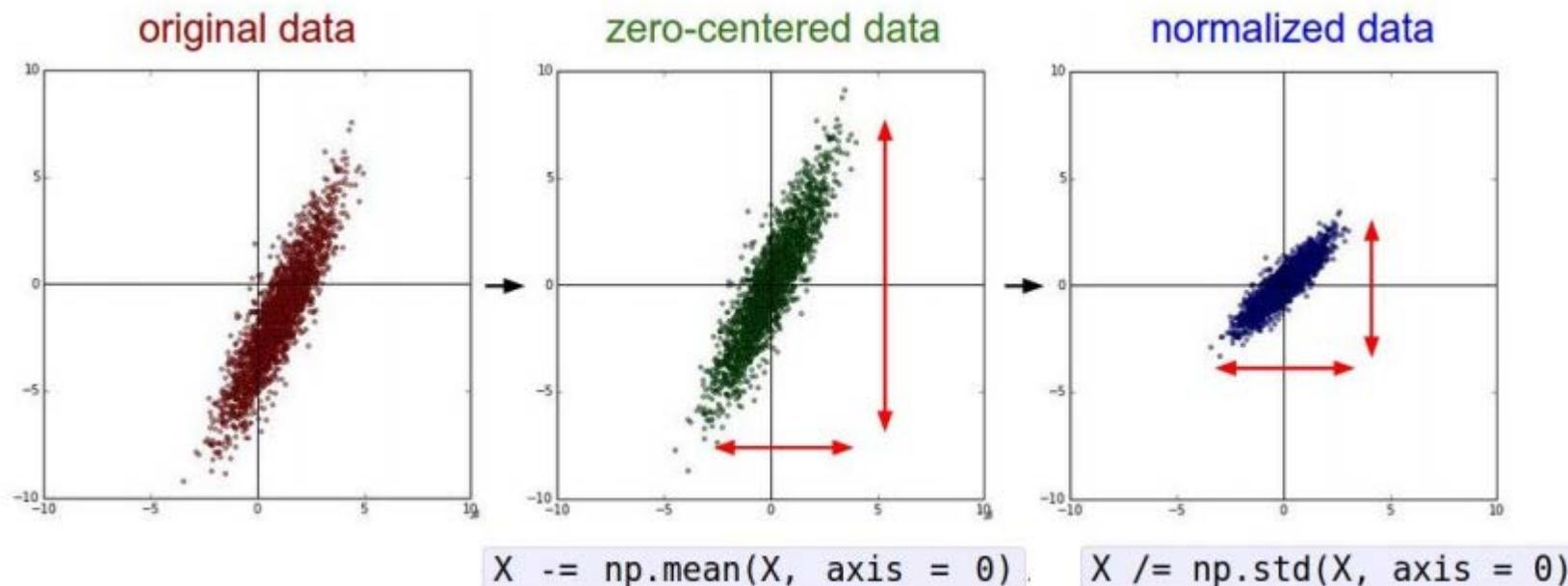
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Leaky ReLU

$$f(x) = \max(0.01x, x)$$

Normalizing inputs

- Normalized inputs helps for the learning process
- Subtract mean and normalize variances
- Use the same mean and variance to normalize the test (you want them to go through the same transitions)



Batch Normalization

- Similar to input normalization, you can normalize the values in the hidden layer
- Two additional parameters to be trained

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

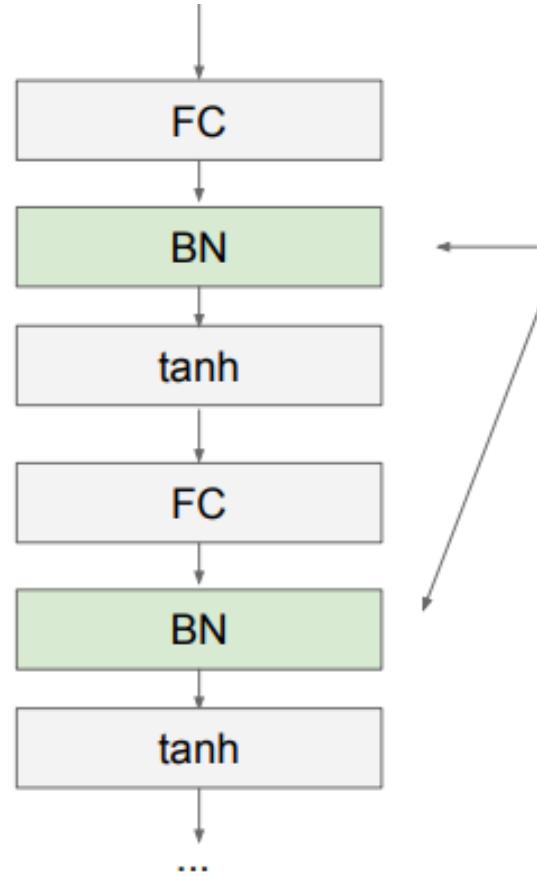
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Batch Normalization



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

Vanishing/exploding gradients

- Vanishing gradients: As we get back deep in the neural network, gradient tends to get smaller through hidden layers
 - In other words, neurons in the earlier layers learn much more slowly than neurons in later layers
- Exploding gradients: Gradients get much larger in earlier layers, unstable gradient
- How you initialize the network weights is important!!

Weight initialization

- Initialize with all 0s or 1s?
 - Behaves like a linear model, hidden units become symmetric
- Traditionally weights of a neural network were set to small random numbers
- Weight initialization is a whole field of study, careful weight initialization can speed up the learning process

<https://machinelearningmastery.com/why-initialize-a-neural-network-with-random-weights/>

<https://medium.com/usf-msds/deep-learning-best-practices-1-weight-initialization-14e5c0295b94>

Weight Initialization (Best practices)

- For $\tanh(z)$ (also called Xavier initialization)

$$\sqrt{\frac{1}{\text{size}^{[l-1]}}}$$

$$W^{[l]} = \text{np.random.randn}(\text{size_}l, \text{size_}l-1) * \text{np.sqrt}(1/\text{size_}l-1)$$

- For $\text{ReLU}(z)$

$$\sqrt{\frac{2}{\text{size}^{[l-1]}}}$$

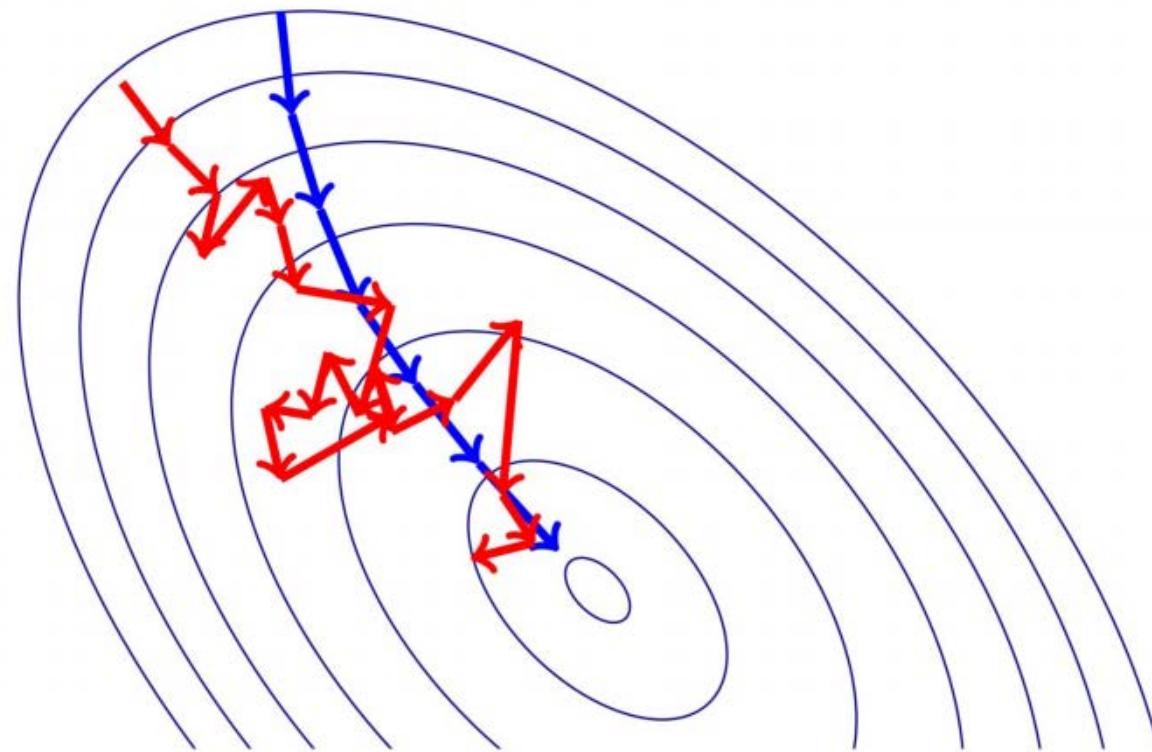
$$W^{[l]} = \text{np.random.randn}(\text{size_}l, \text{size_}l-1) * \text{np.sqrt}(2/\text{size_}l-1)$$

Proper initialization is an active area of research...

- **Understanding the difficulty of training deep feedforward neural networks** Glorot and Bengio, 2010
- **Exact solutions to the nonlinear dynamics of learning in deep linear neural networks** Saxe et al, 2013
- **Random walk initialization for training very deep feedforward networks** Sussillo and Abbott, 2014
- **Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification** He et al., 2015
- **Data-dependent Initializations of Convolutional Neural Networks** Krähenbühl et al., 2015
- **All you need is a good init**, Mishkin and Matas, 2015

...

Stochastic Gradient Descent vs Gradient Descent

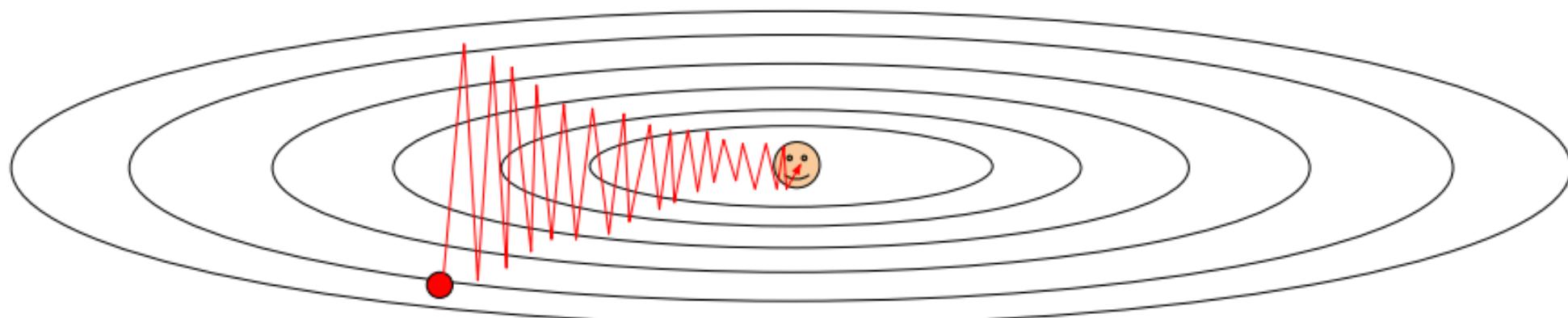


Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?

What does gradient descent do?

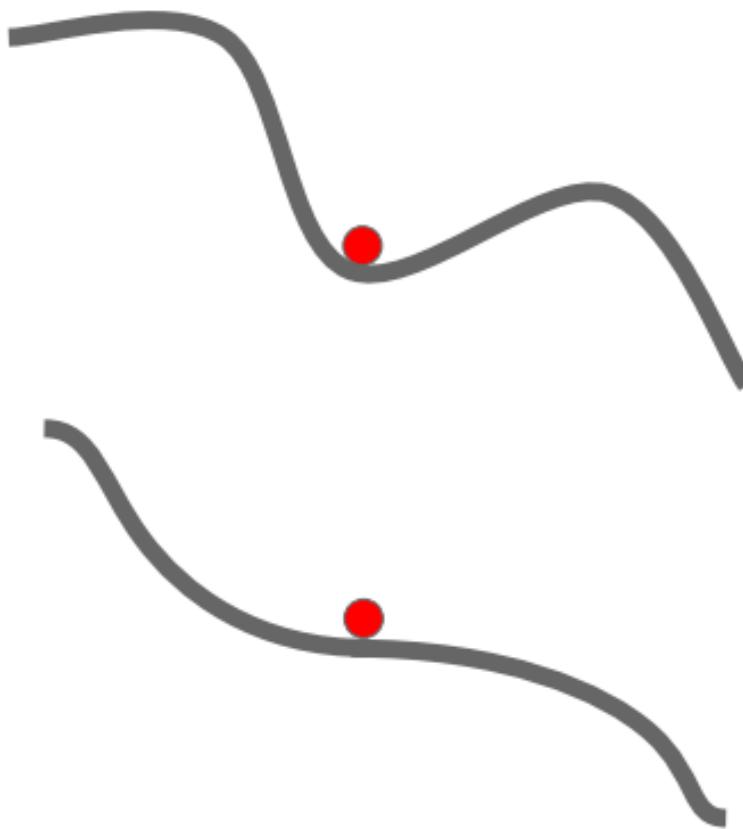
Very slow progress along shallow dimension, jitter along steep direction



Optimization: Problems with SGD

What if the loss
function has a
local minima or
saddle point?

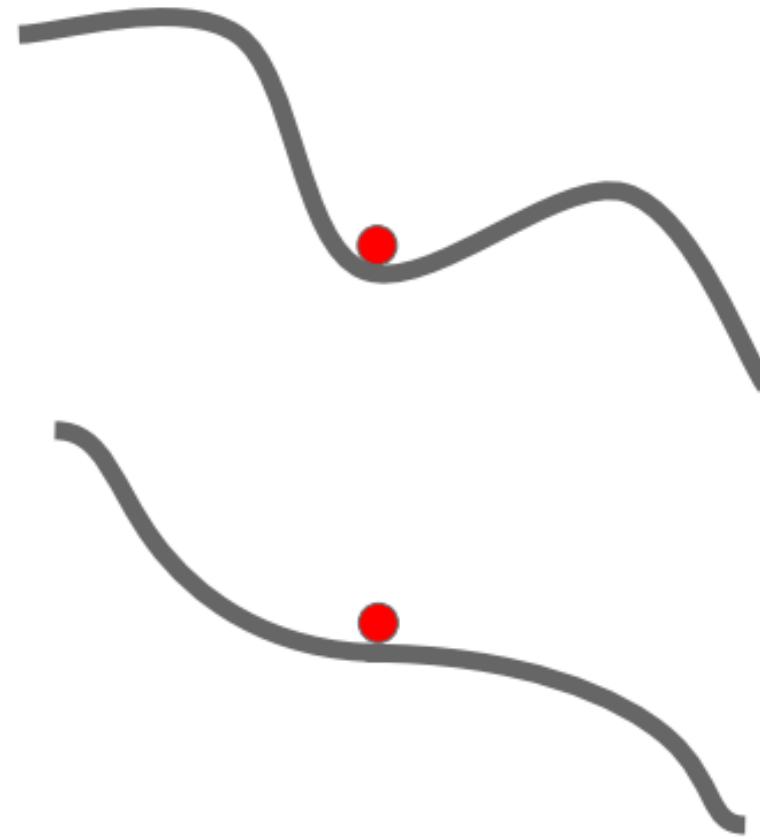
Zero gradient,
gradient descent
gets stuck



Optimization: Problems with SGD

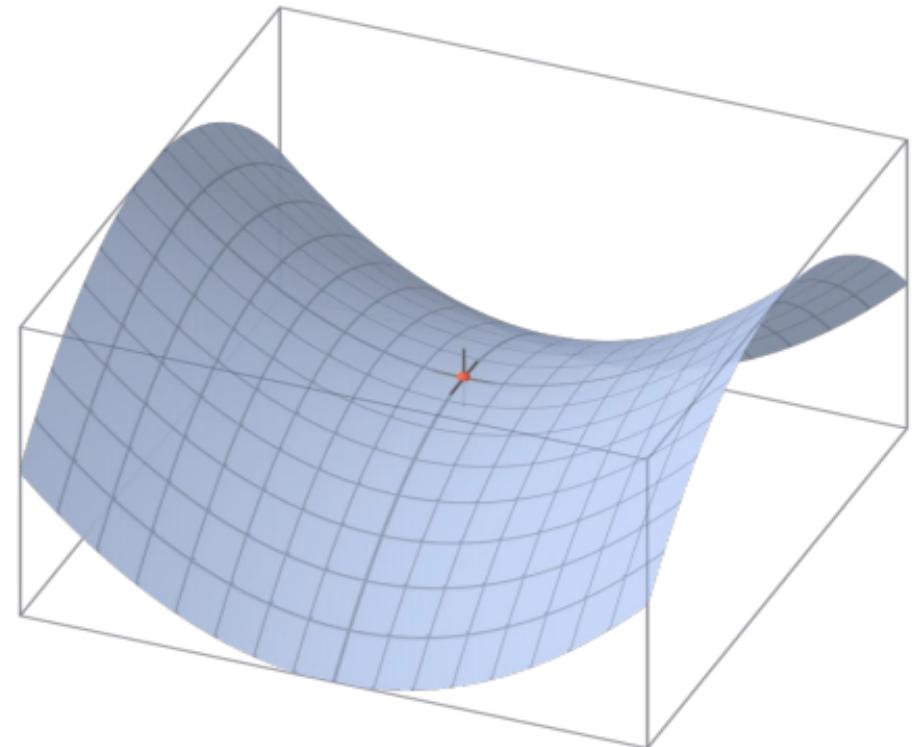
What if the loss
function has a
local minima or
saddle point?

Saddle points much
more common in
high dimension



Saddle points

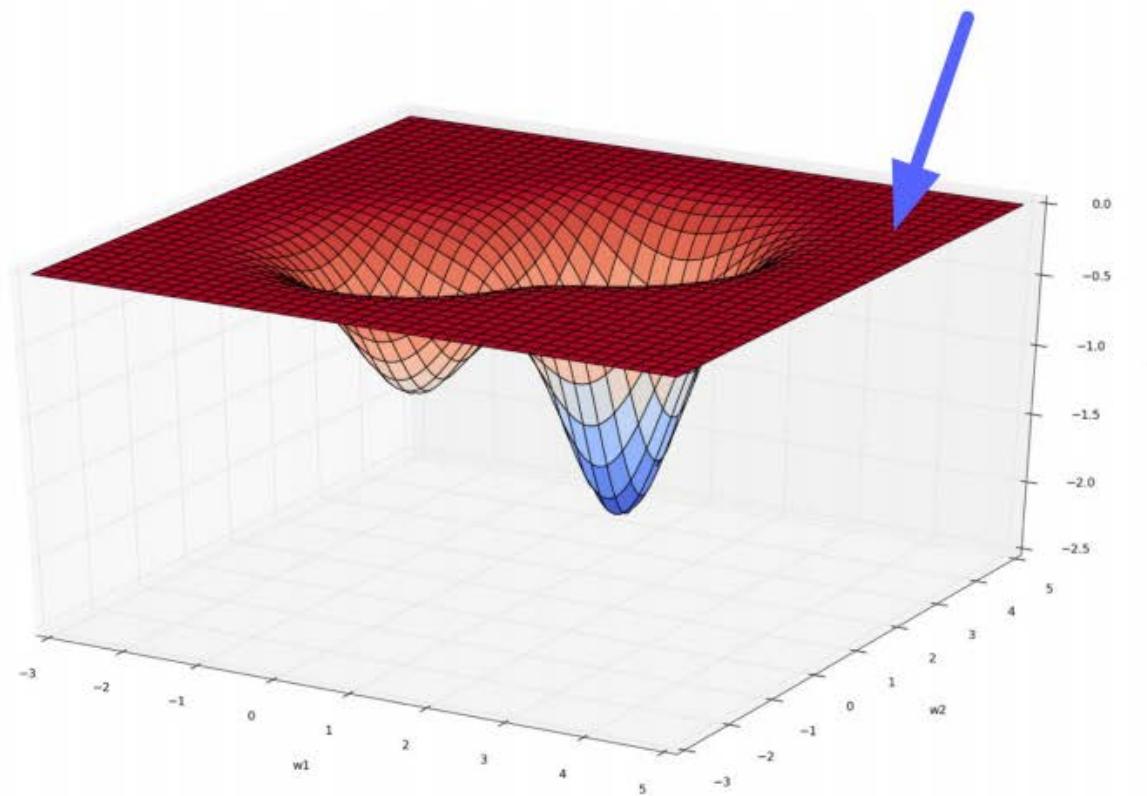
- At a **saddle point**, $\frac{\partial L}{\partial W} = 0$ even though we are not at a minimum. Some directions curve upwards, and others curve downwards.
- When would saddle points be a problem?
 - If we're exactly on the saddle point, then we're stuck.
 - If we're slightly to the side, then we can get unstuck.



Saddle points much more common in high dimensions!

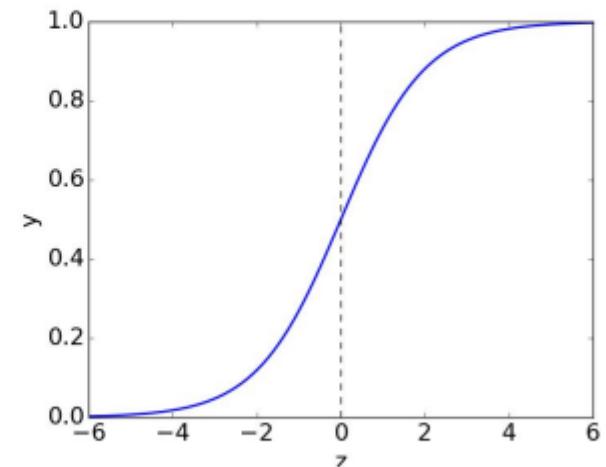
Plateaux

- A flat region is called a **plateau**. (Plural: plateaux)



Plateaux

- An important example of a plateau is a **saturated unit**. This is when it is in the flat region of its activation function.
- If $\phi'(z_i)$ is always close to zero, then the weights will get stuck.
- If there is a ReLU unit whose input z_i is always negative, the weight derivatives will be exactly 0. We call this a **dead unit**.



SGD + Momentum

SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:  
    dx = compute_gradient(x)  
    x += learning_rate * dx
```

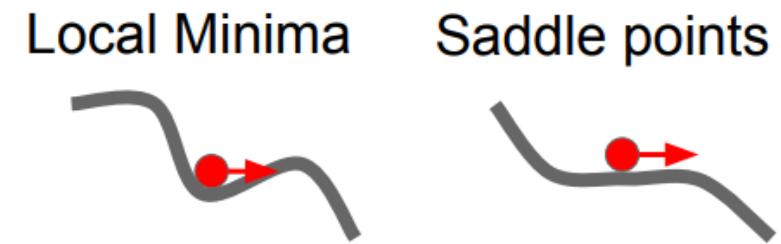
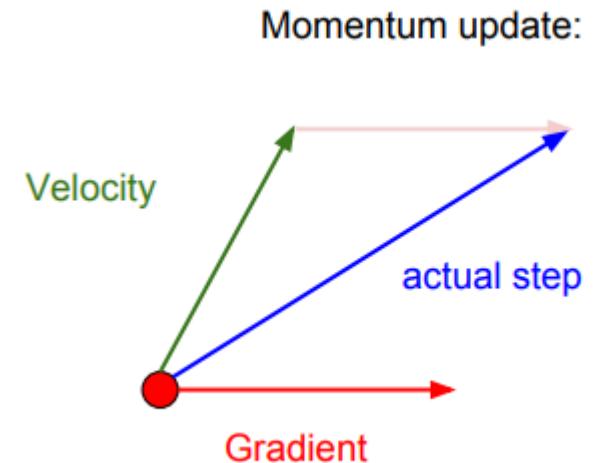
SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0  
while True:  
    dx = compute_gradient(x)  
    vx = rho * vx + dx  
    x += learning_rate * vx
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho=0.9 or 0.99



AdaGrad and RMSProp (Root Mean square prop)

AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



RMSProp

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

Adam (Adaptive Moment Estimation)

```
first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7))
```

Momentum

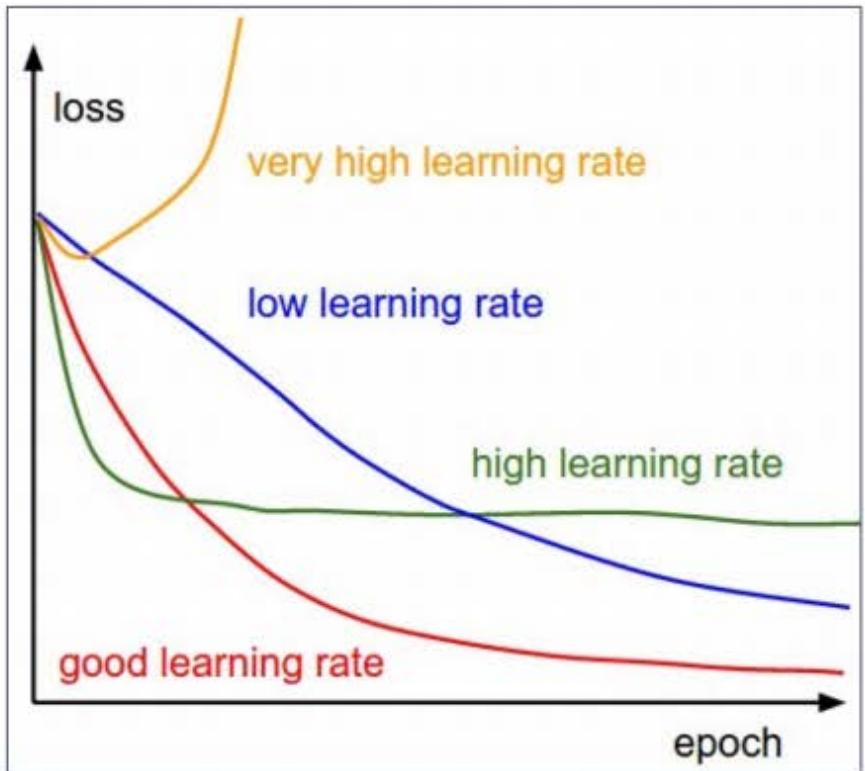
Bias correction

AdaGrad / RMSProp

Bias correction for the fact that
first and second moment
estimates start at zero

Adam with $\beta_1 = 0.9$,
 $\beta_2 = 0.999$, and $\text{learning_rate} = 1e-3$ or $5e-4$
is a great starting point for many models!

SGD, SGD+Momentum, Adagrad, RMSProp, Adam
all have learning rate as a hyperparameter



=> Learning rate decay over time!

step decay:

e.g. decay learning rate by half every few epochs.

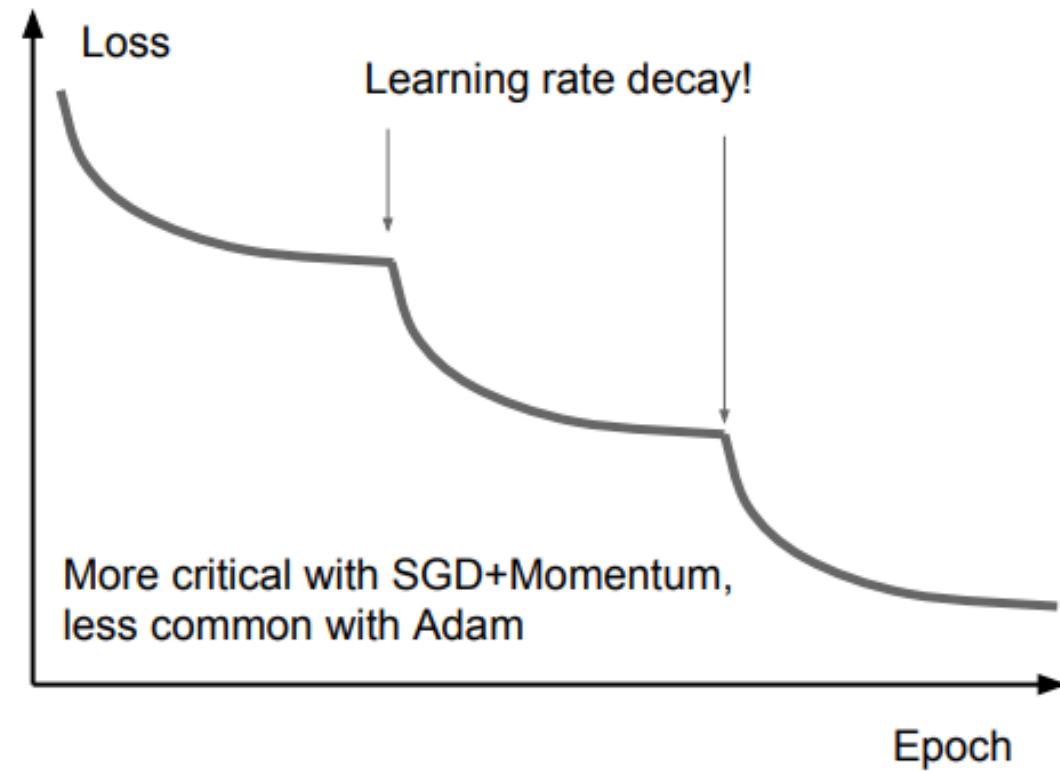
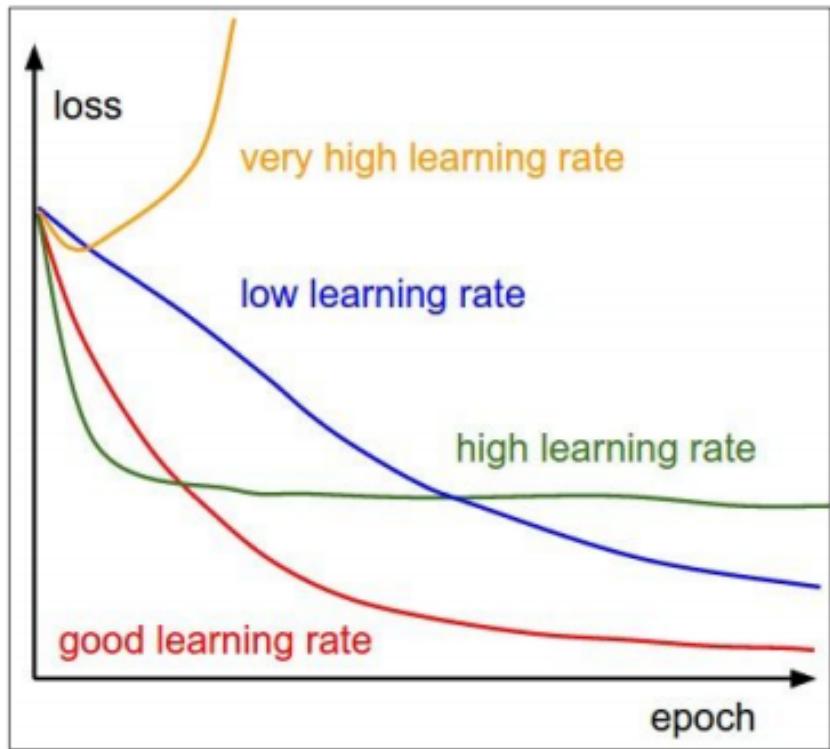
exponential decay:

$$\alpha = \alpha_0 e^{-kt}$$

1/t decay:

$$\alpha = \alpha_0 / (1 + kt)$$

SGD, SGD+Momentum, Adagrad, RMSProp, Adam
all have learning rate as a hyperparameter



Hyperparameters tuning

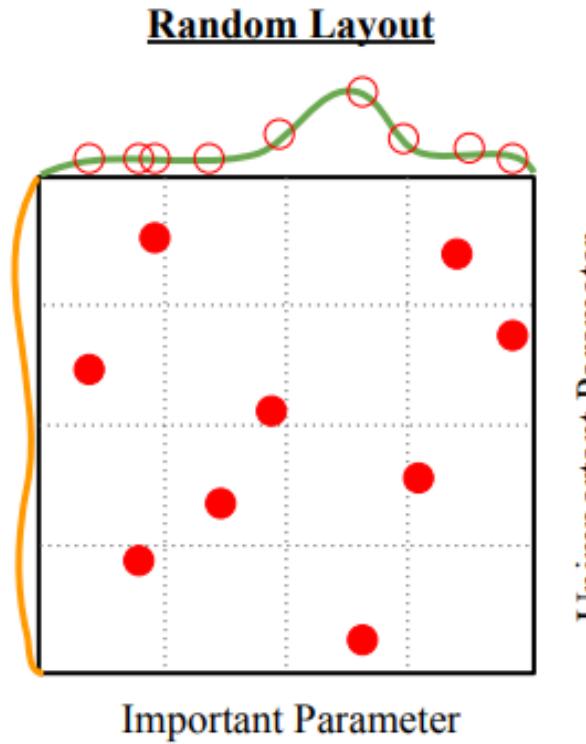
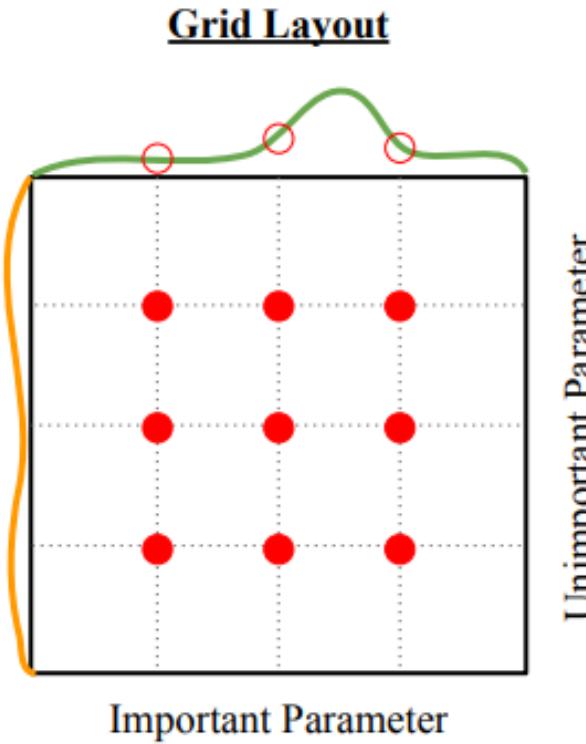


Illustration of Bergstra et al., 2012 by Shayne Longpre, copyright CS231n 2017

Hyperparameters tuning

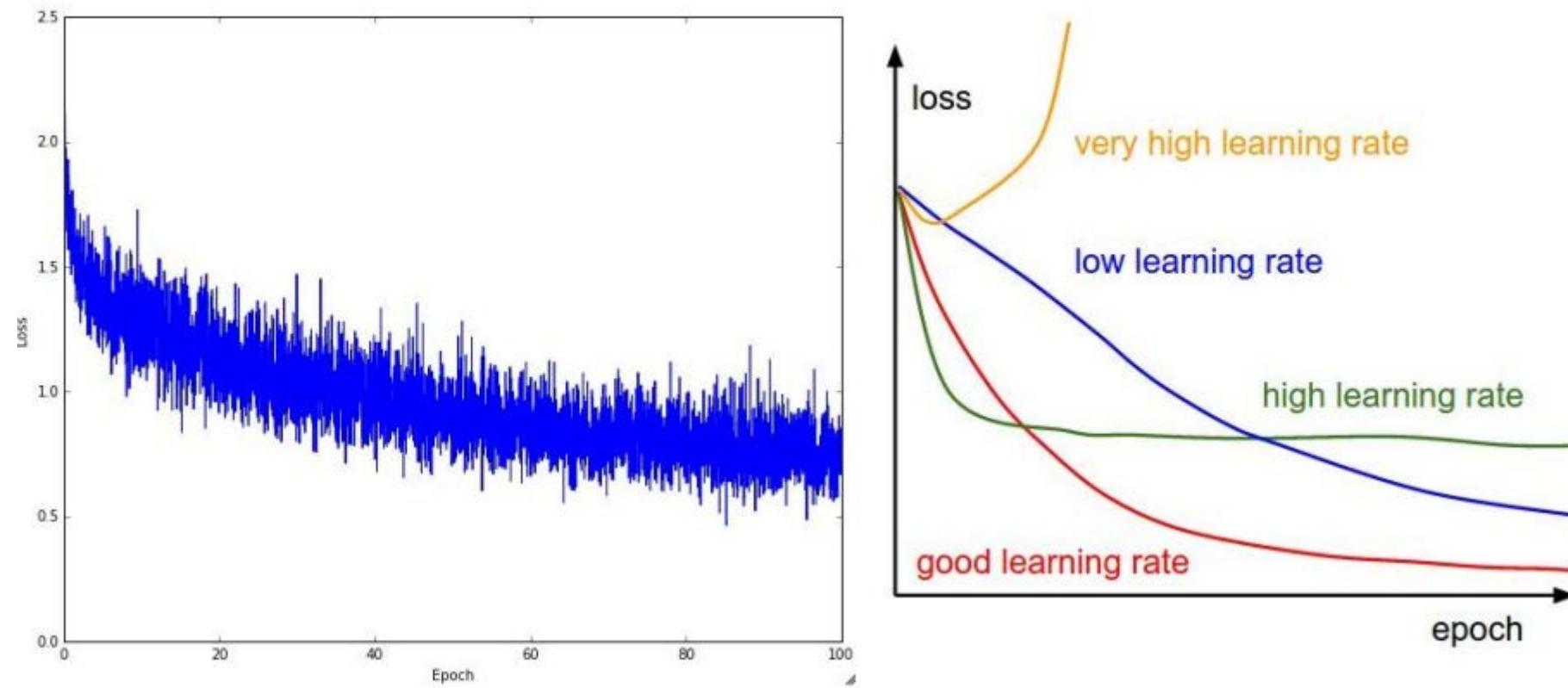
- network architecture
- learning rate, its decay schedule, update type
- regularization (L2/Dropout strength)

neural networks practitioner
music = loss function

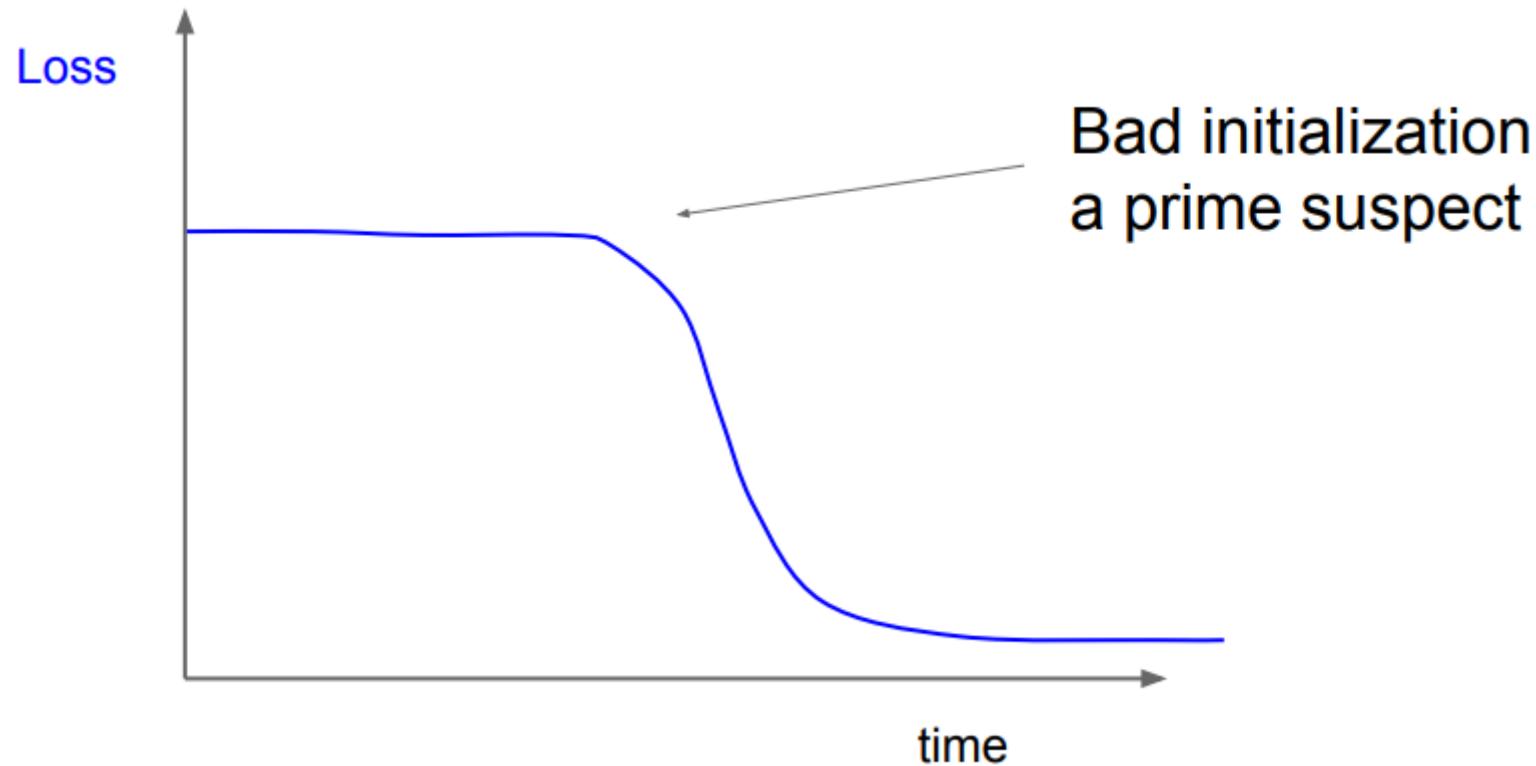


[This image](#) by Paolo Guereta is licensed under [CC-BY 2.0](#)

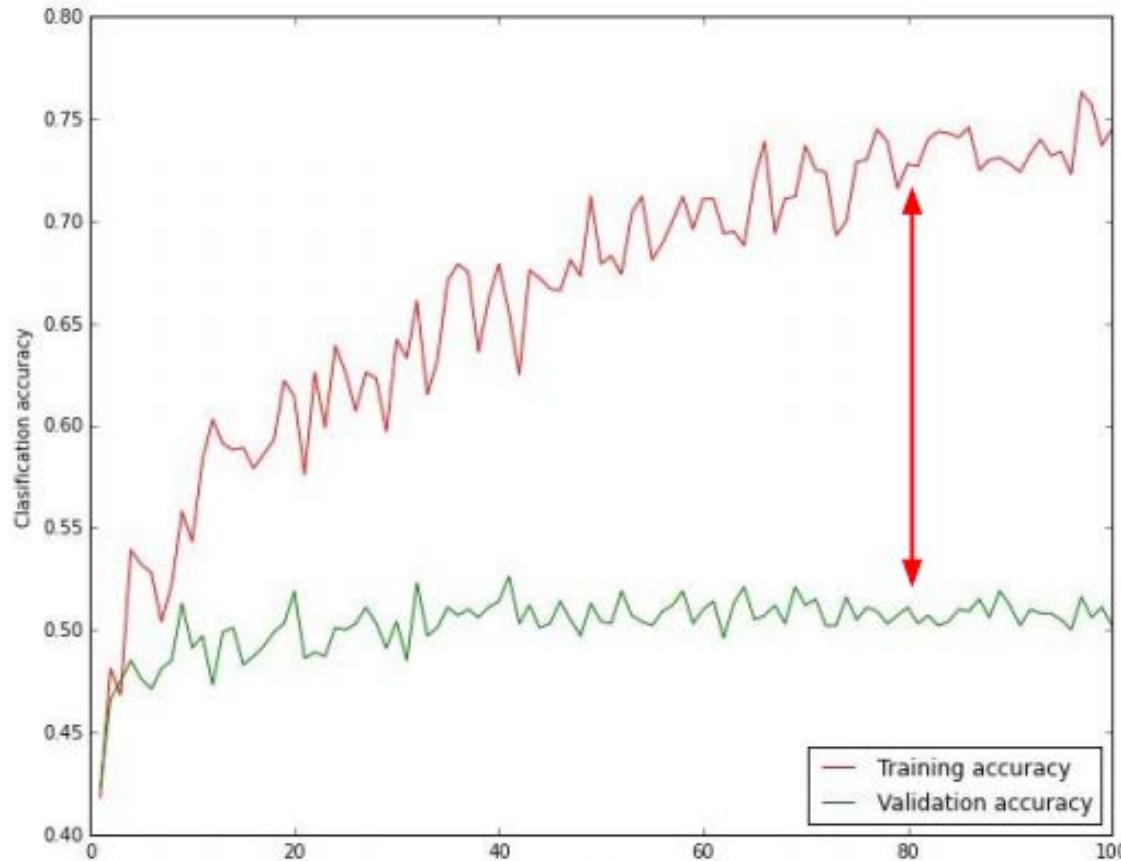
Monitor and visualize the loss curve



Monitor and visualize the loss curve



Monitor and visualize the accuracy



big gap = overfitting
=> increase regularization strength?

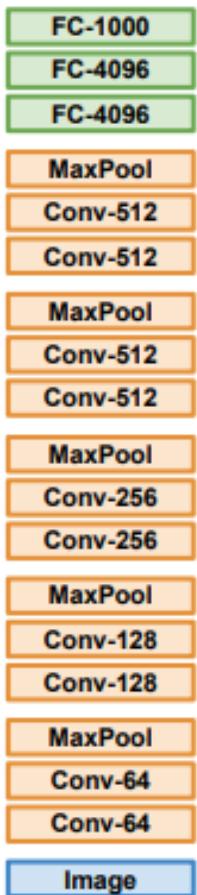
no gap
=> increase model capacity?

Babysitting one model vs training many models

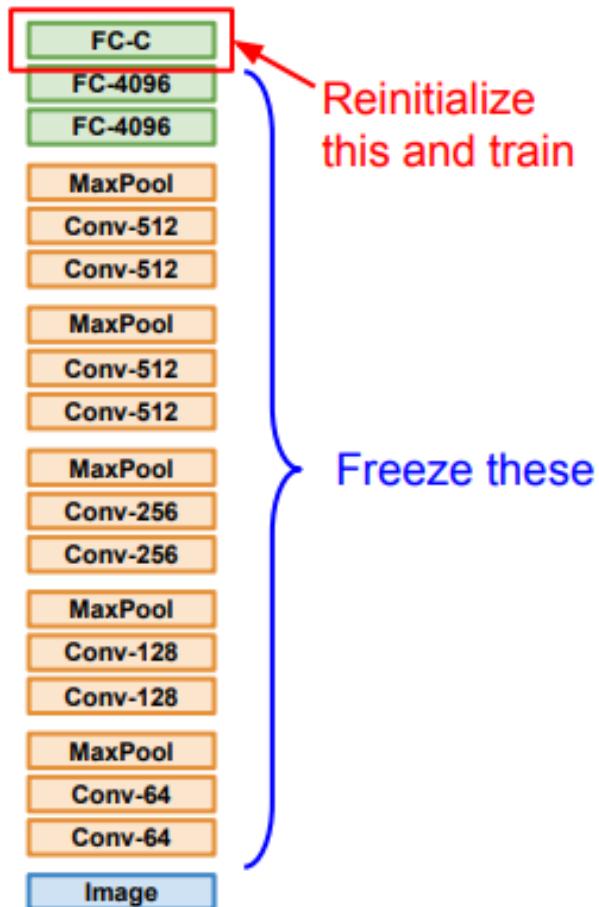
- Model Ensembles
- 1. Train multiple independent models
- 2. At test time average their results
- Enjoy 2% extra performance

Transfer learning

1. Train on Imagenet



2. Small Dataset (C classes)



Freeze these

Reinitialize
this and train

Deep learning frameworks provide models of pretrained models so you might not need to train your own:

Caffe: <https://github.com/BVLC/caffe/wiki/Model-Zoo>

TensorFlow: <https://github.com/tensorflow/models>

PyTorch: <https://github.com/pytorch/vision>

Summary

- Many steps and parameters
 - Normalization
 - Weight initialization
 - Learning rate
 - Number of hidden units
 - Mini-batch size
 - Number of layers
 - Batch normalization
 - Optimization algorithms
 - Learning rate decay

In your projects..

- Describe the steps you went through? e.g.
 - What is the training, validation, test set? Why did you split the data like this?
 - Which hyperparameters did you test first, why?
- Compare and reason about the results by looking at the loss curve and accuracy, e.g.
 - Compare different weight initialization methods
 - Compare different activation functions
 - Compare different optimization functions
 - Try different learning rates
 - Compare with and without batch normalization
 - Etc..
- Give also performance metrics
 - How much time it took for training?
 - How much time it took for testing?
 - On CPU, GPU? What are the machine specs?

Reading the research papers, critical thinking and in-depth analysis results into higher grades!
Avoid saying “We applied this and it worked well”. Try to explain why it worked!

Thoughts on research

- Scientific truth does not follow the fashion
 - Do not hesitate being a contrarian if you have good reasons
- Experiments are crucial
 - Do not aim at beating the state-of-the-art, aim at understanding the phenomena
- On the proper use of mathematics
 - A theorem is not like a subroutine that one can apply blindly
 - Theorems should not limit creativity

Supplementary reading and video

- [Deep Learning book](#), Chapter 6, 7 and 8
- <http://neuralnetworksanddeeplearning.com/>, Michael Nielsen
- <https://www.youtube.com/playlist?list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH>, Hugo Larochelle's video lectures (1.1 to 2.7)
- <https://webcolleges.uva.nl/Mediasite/Play/947ccbc9b11940c0ad5ab39ebb154c461d>, Efstratios Gavves' [Lecture 3](#)
- Machine Learning and Deep Learning courses on Coursera by Andrew Ng

References

- MIT 6. S191 Introduction to Deep Learning
- CS231n: Convolutional Neural Networks
- CMP8784: Deep Learning, Hacettepe University
- (Slides mainly adopted from the above courses)

Tensorflow tutorial

- <https://www.tensorflow.org/tutorials/>