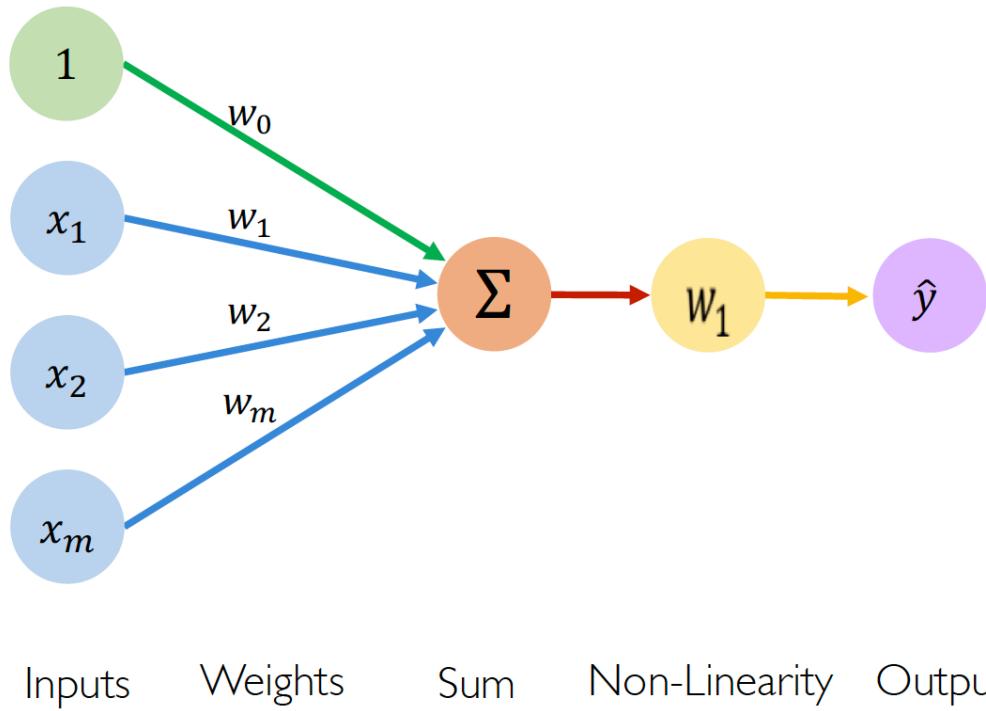


Lecture 24: Deep Neural Network

BIOS635

04/16/2020

The perceptron: the forward propagation



Linear combination of inputs \downarrow

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

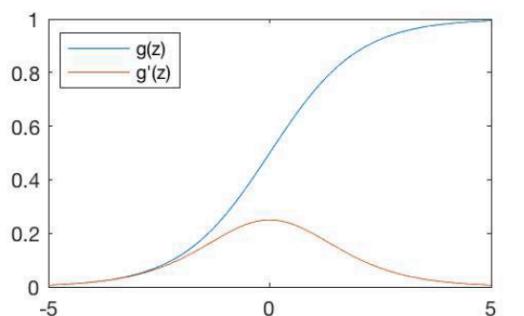
Non-linear activation function \uparrow

Bias \nearrow

Output \downarrow

Common activation function

Sigmoid Function



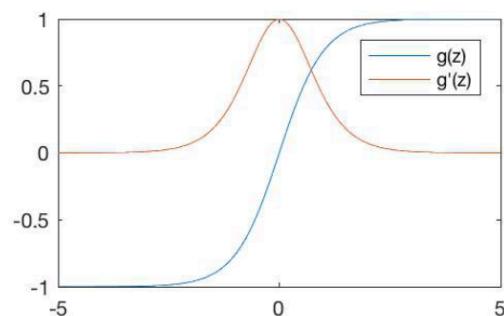
$$g(z) = \frac{1}{1 + e^{-z}}$$

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



`tf.nn.sigmoid(z)`

Hyperbolic Tangent



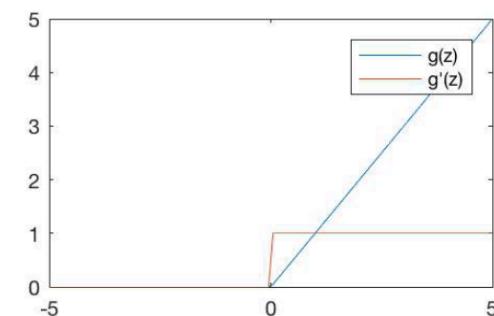
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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



`tf.nn.tanh(z)`

Rectified Linear Unit (ReLU)



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

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

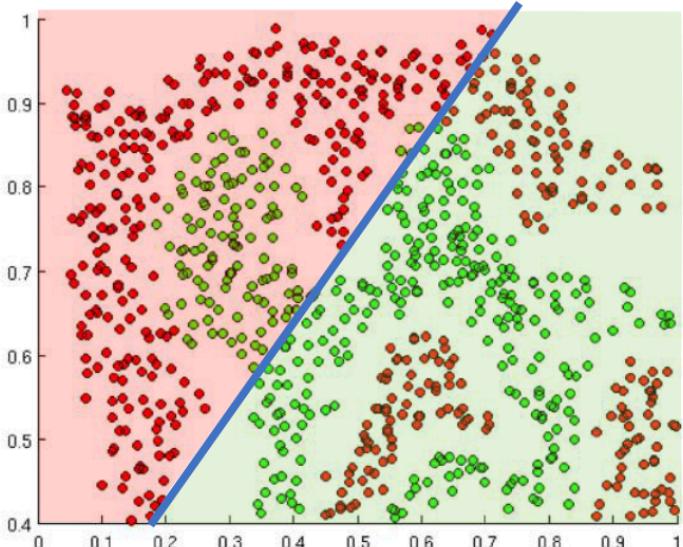


`tf.nn.relu(z)`

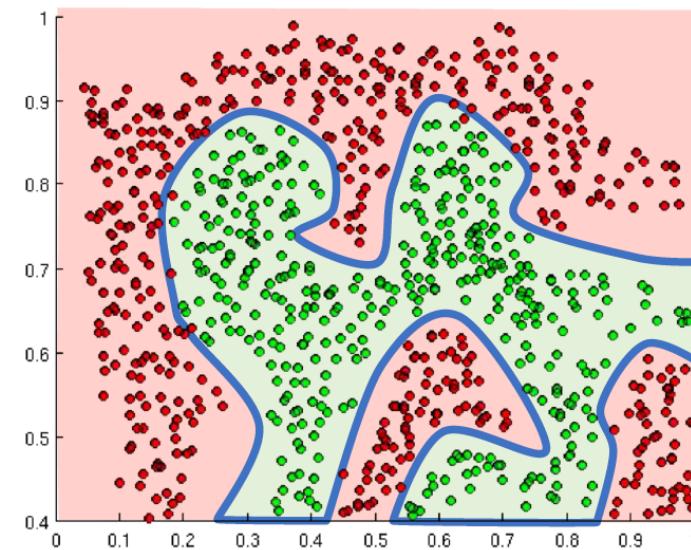
NOTE: All activation functions are non-linear

Importance of activation functions

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

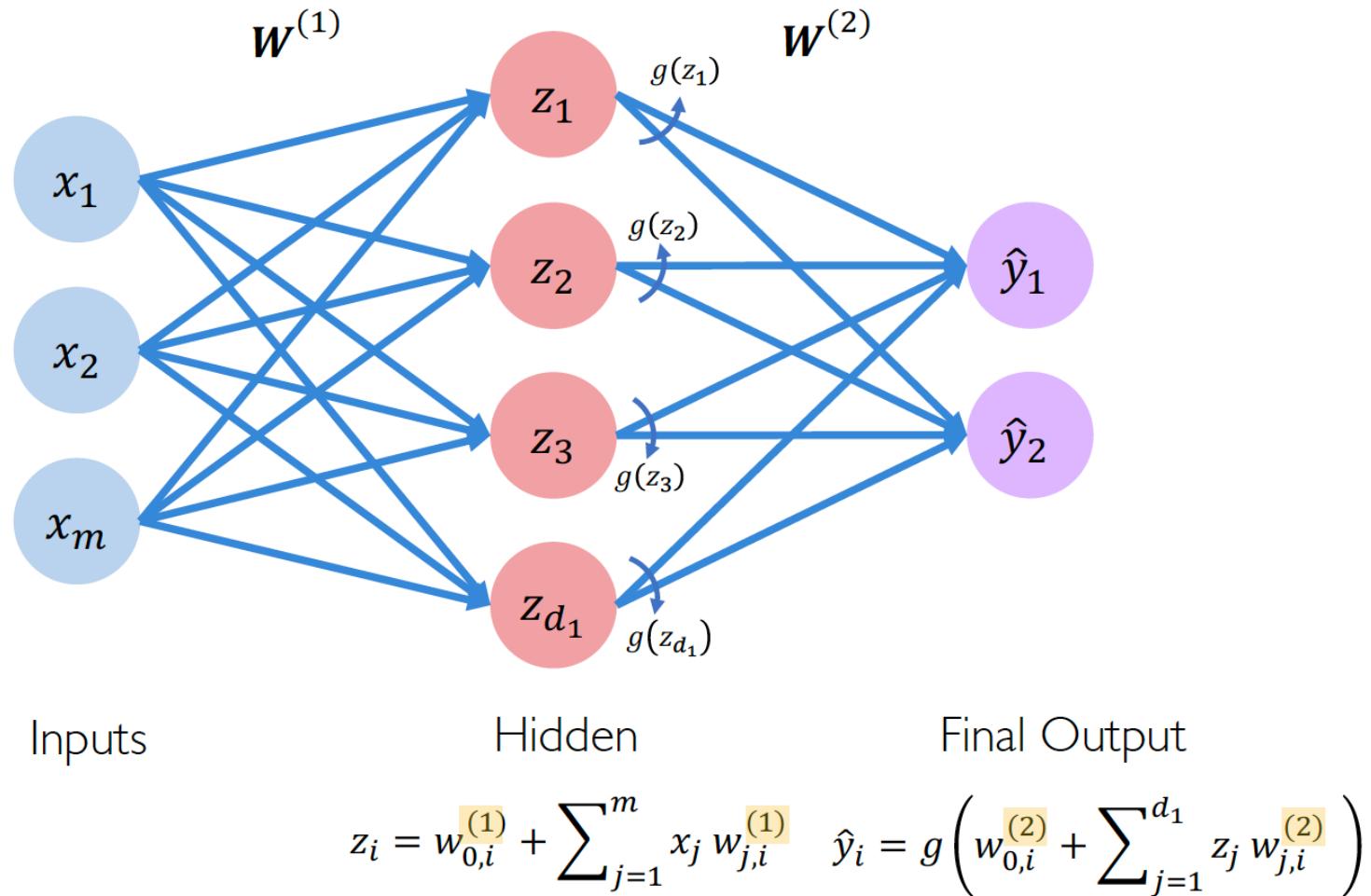


Linear Activation functions produce linear decisions no matter the network size

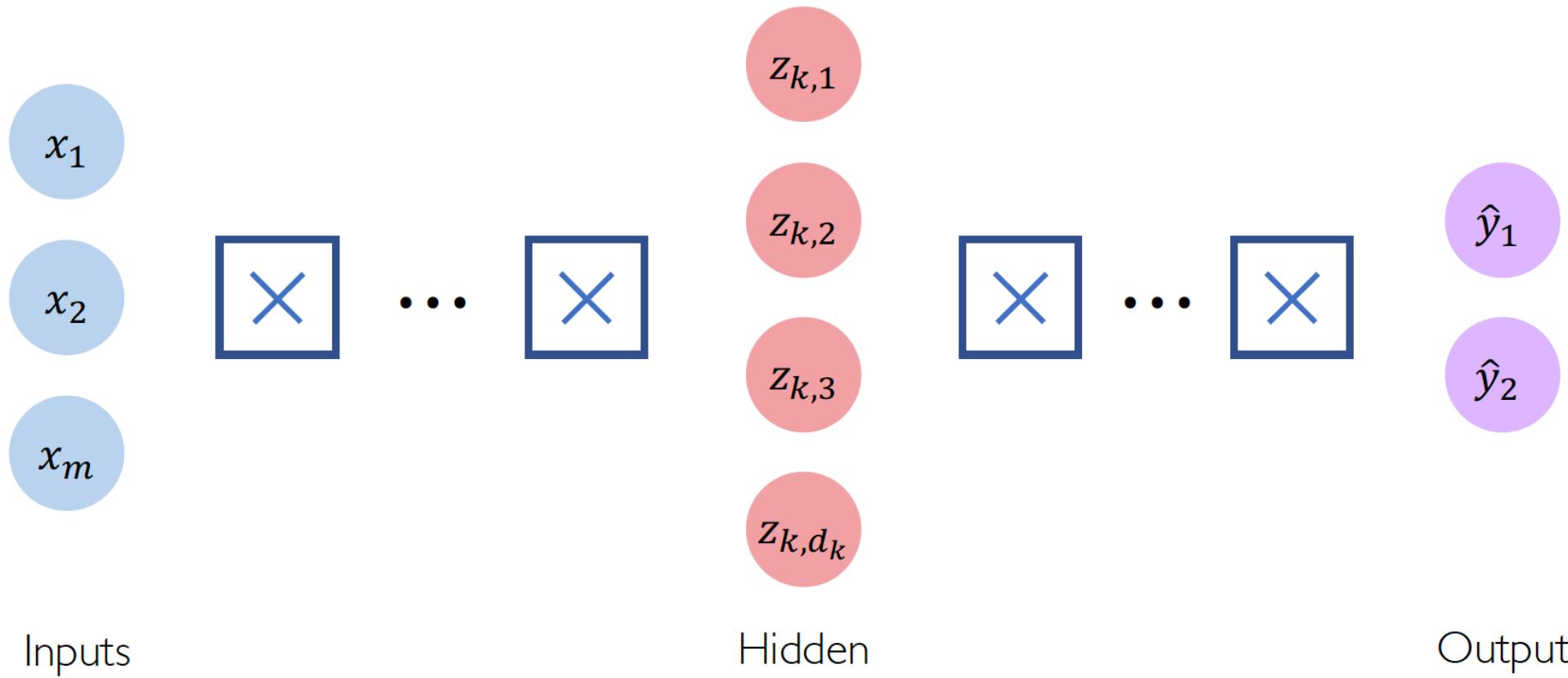


Non-linearities allow us to approximate arbitrarily complex functions

Perceptron: one hidden layer neural network



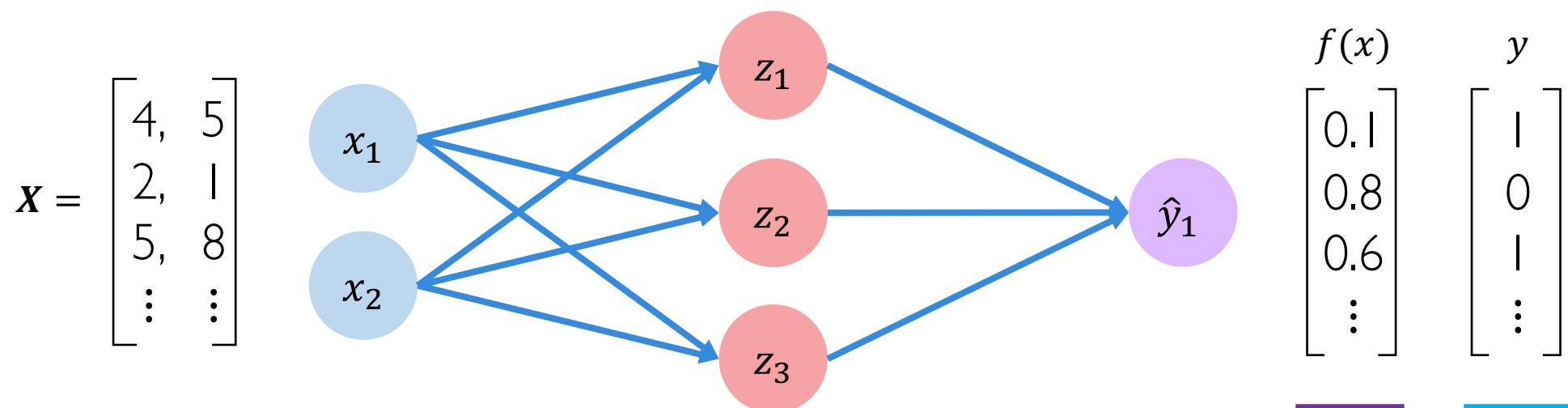
Deep neural network



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

Empirical loss

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



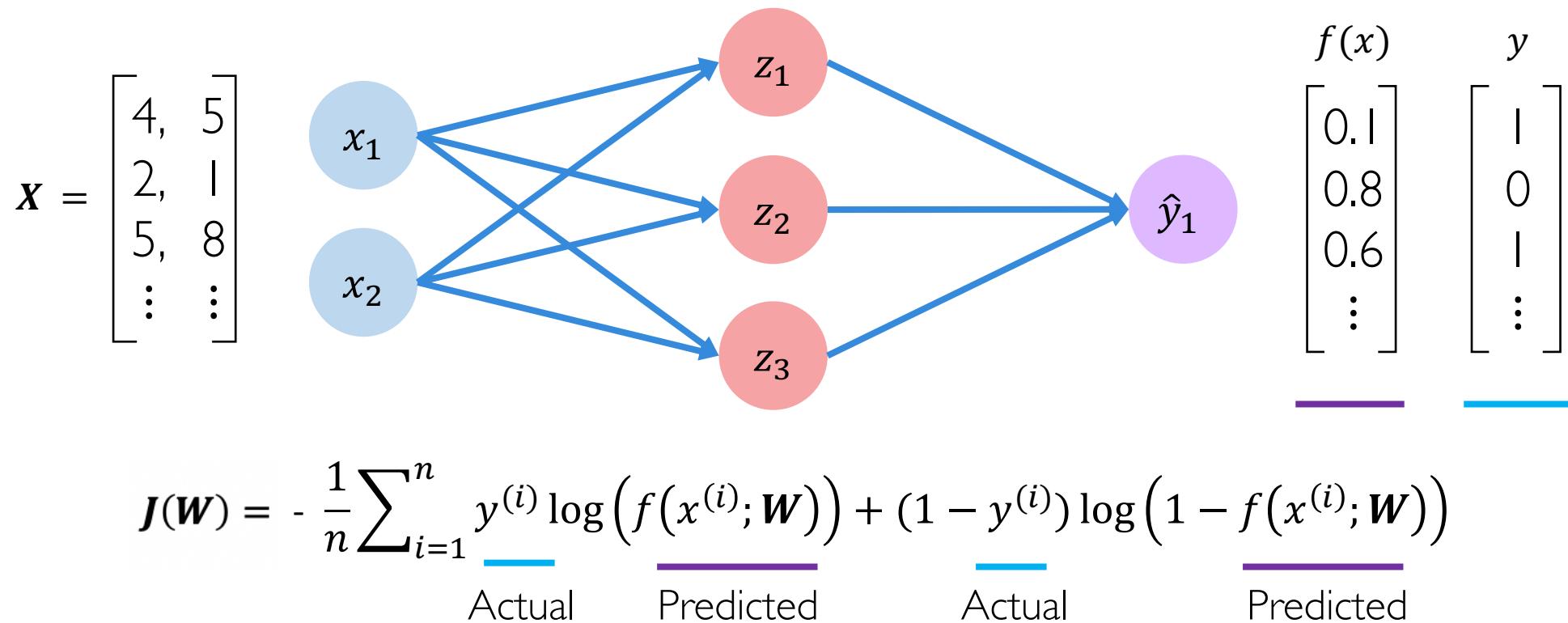
- Also known as:
- Objective function
 - Cost function
 - Empirical Risk

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

Predicted Actual

Binary cross entropy loss

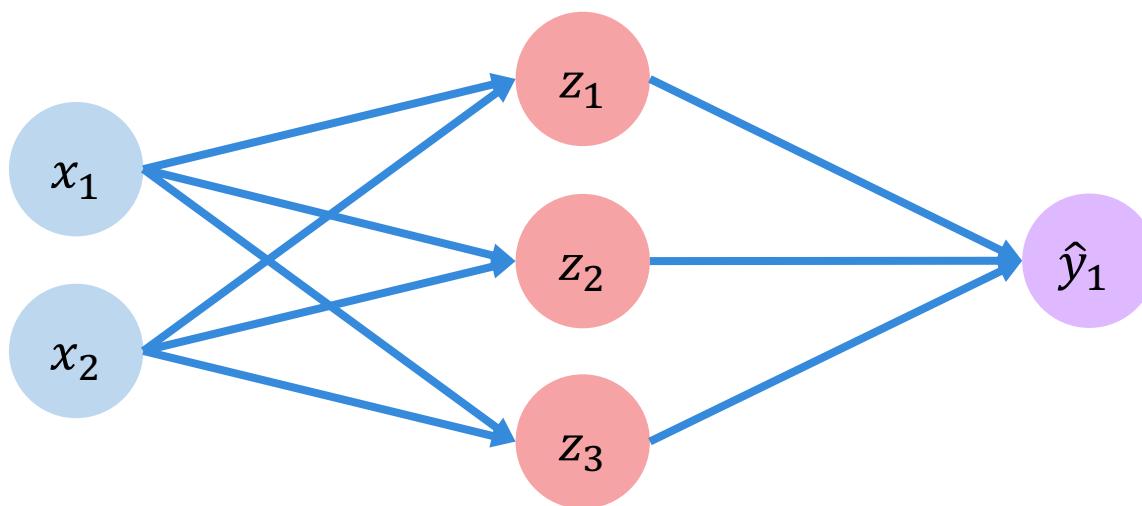
Cross entropy loss can be used with models that output a probability between 0 and 1



Mean squared error loss

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

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



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

Actual Predicted

$f(x)$	y
$\begin{bmatrix} 30 \\ 80 \\ 85 \\ \vdots \end{bmatrix}$	$\begin{bmatrix} 90 \\ 20 \\ 95 \\ \vdots \end{bmatrix}$



Final Grades
(percentage)

Loss optimization

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

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

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

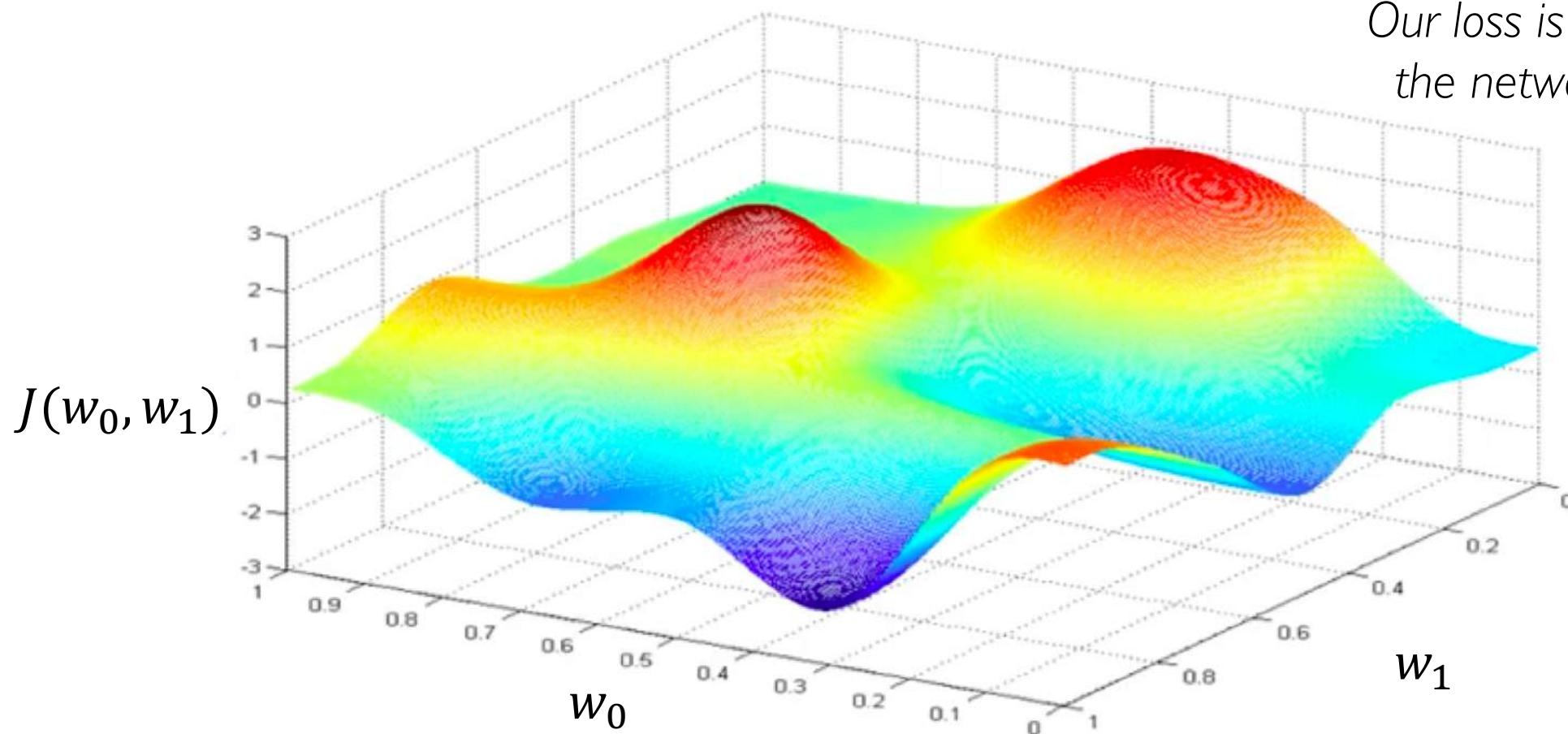


Remember:

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

Loss optimization

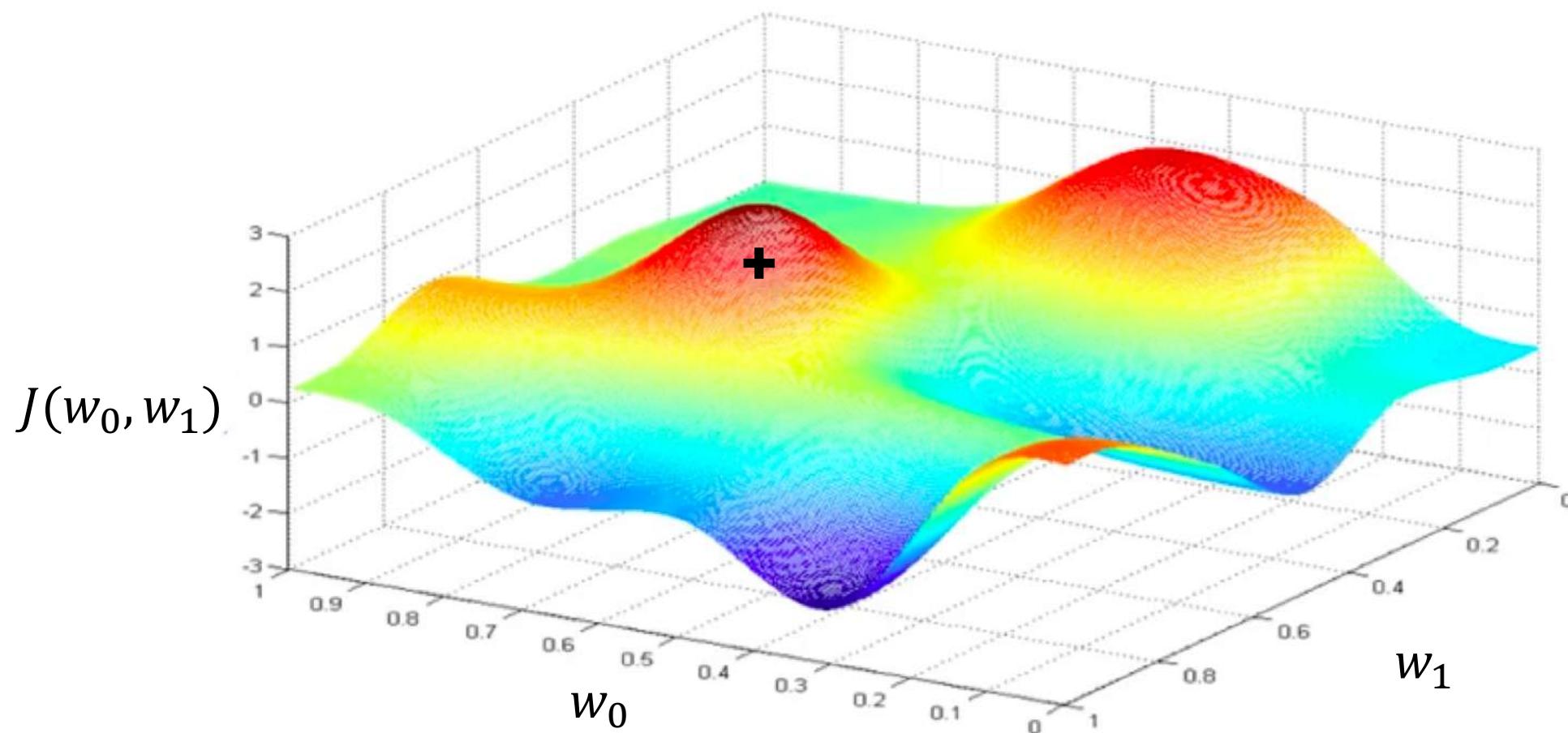
$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$



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

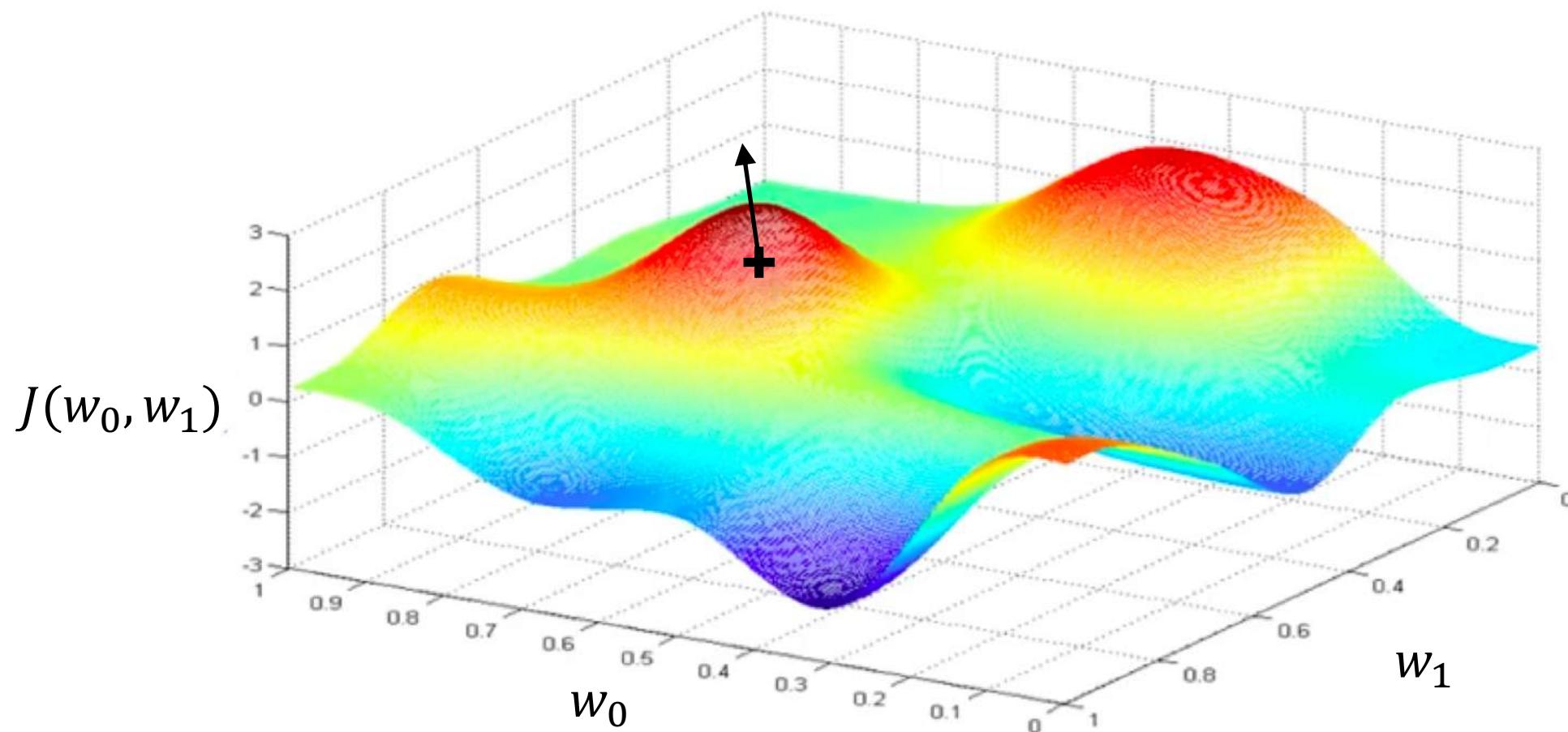
Loss optimization

Randomly pick an initial (w_0, w_1)



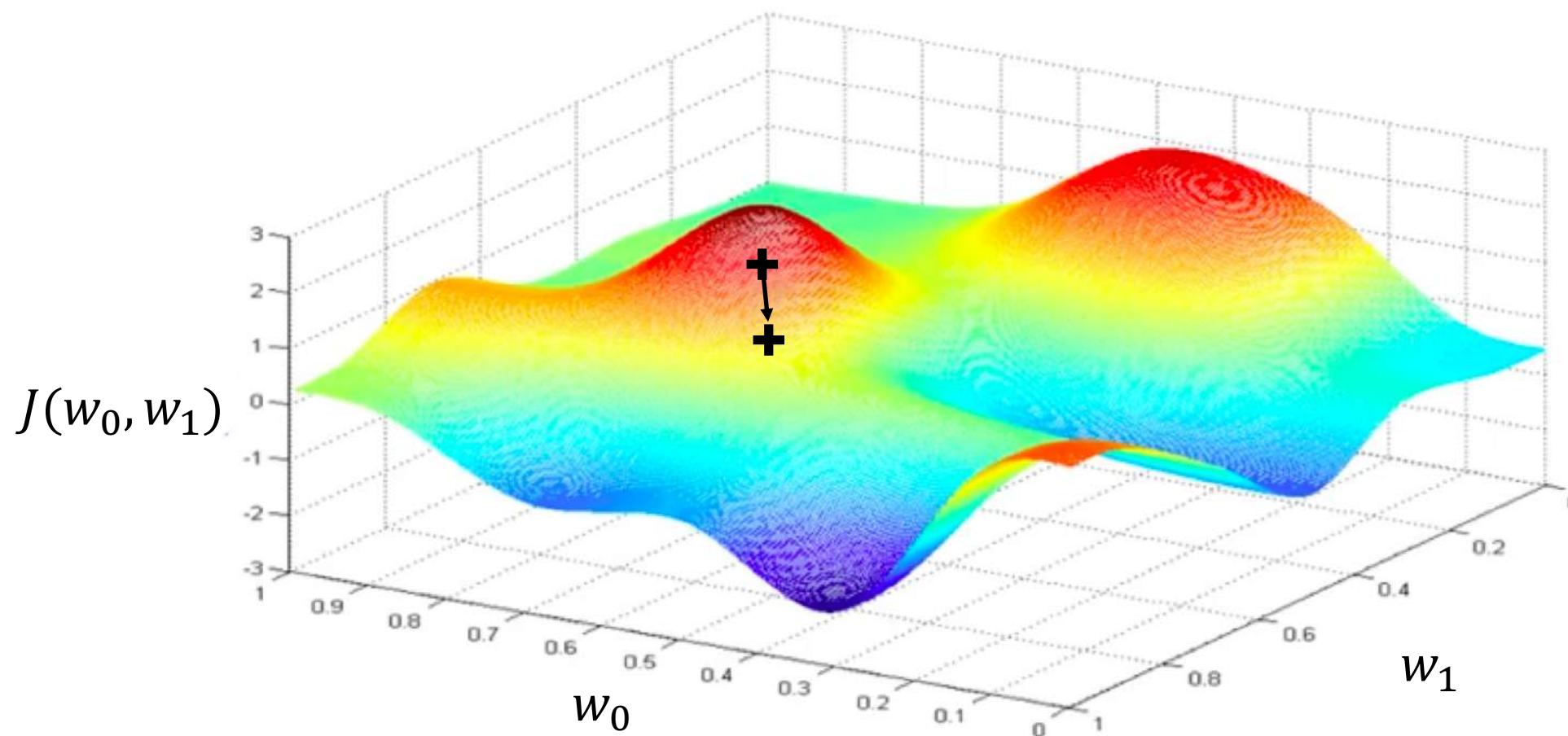
Loss optimization

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



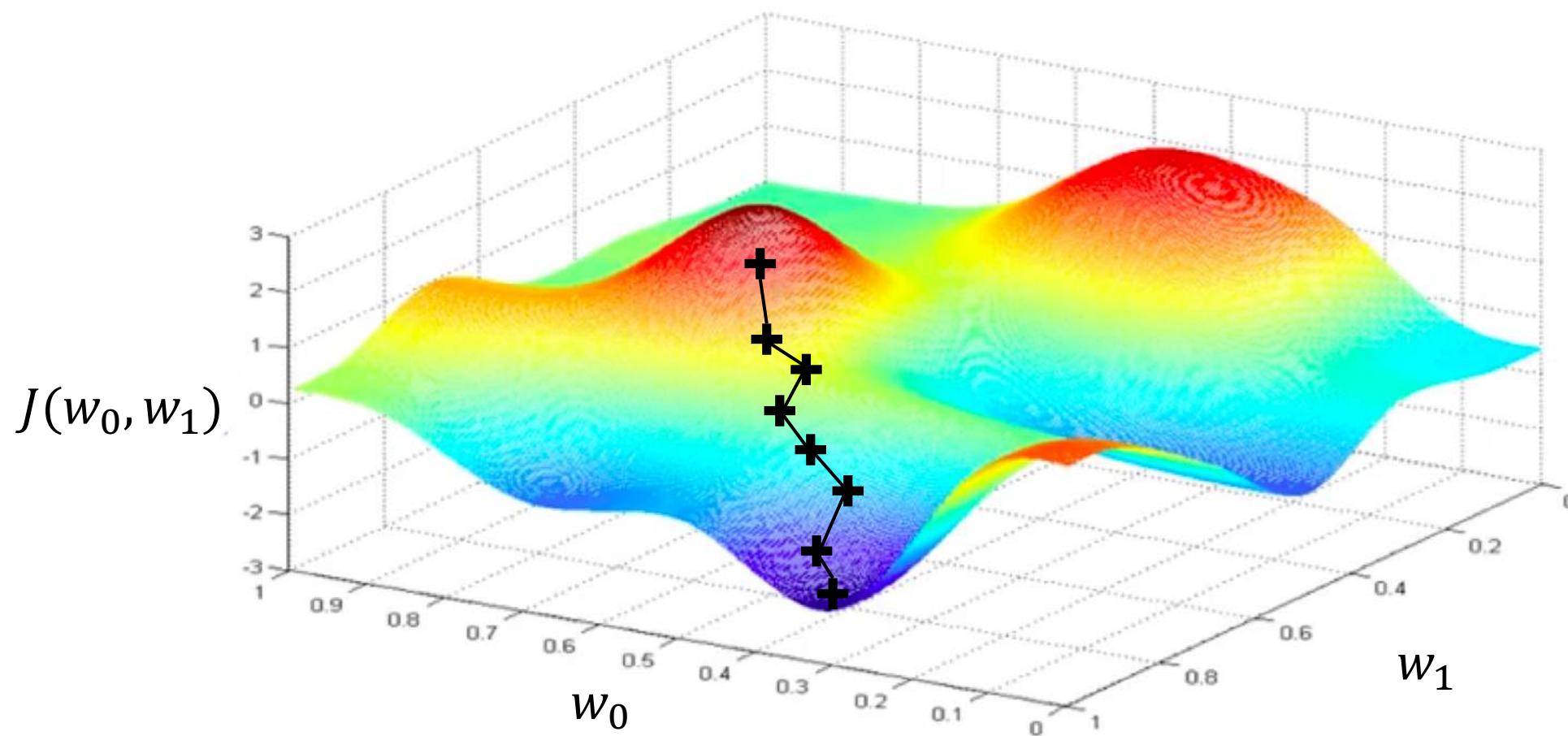
Loss optimization

Take small step in opposite direction of gradient



Loss optimization

Repeat until convergence

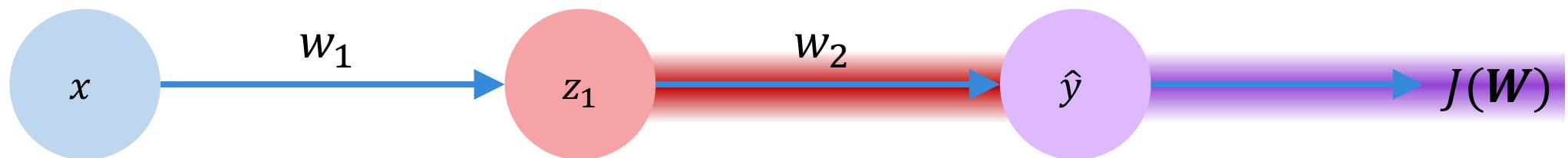


Gradient descent

Algorithm

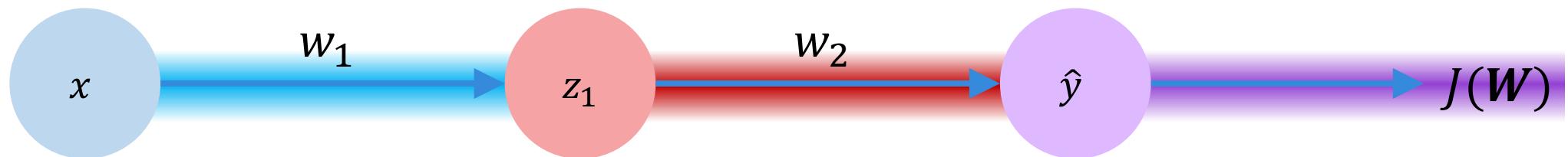
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Computing gradients: back propagation



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

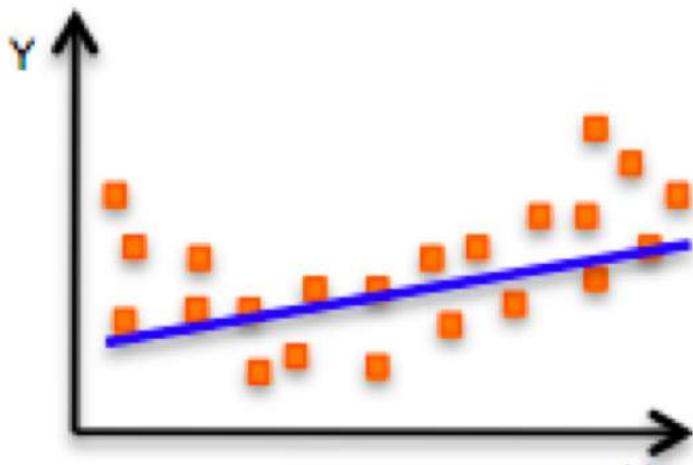
Computing gradients: back propagation



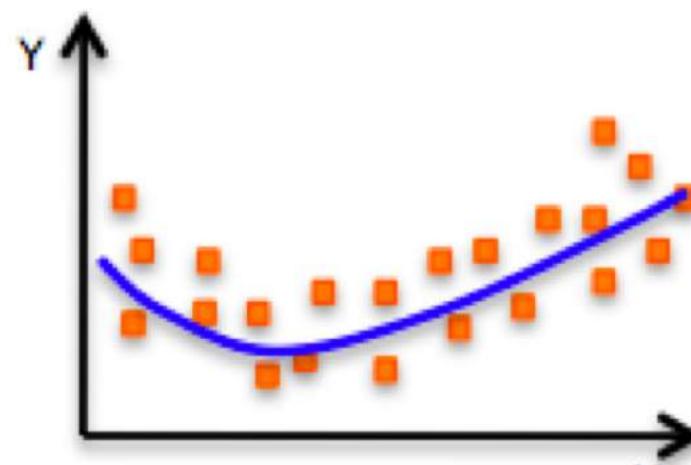
$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underbrace{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}}_{\text{purple bar}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red bar}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue bar}}$$

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

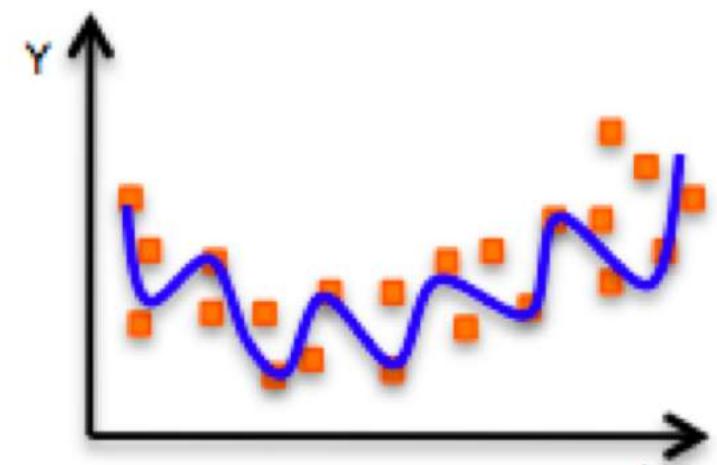
The problem of overfitting



Underfitting
Model does not have capacity
to fully learn the data



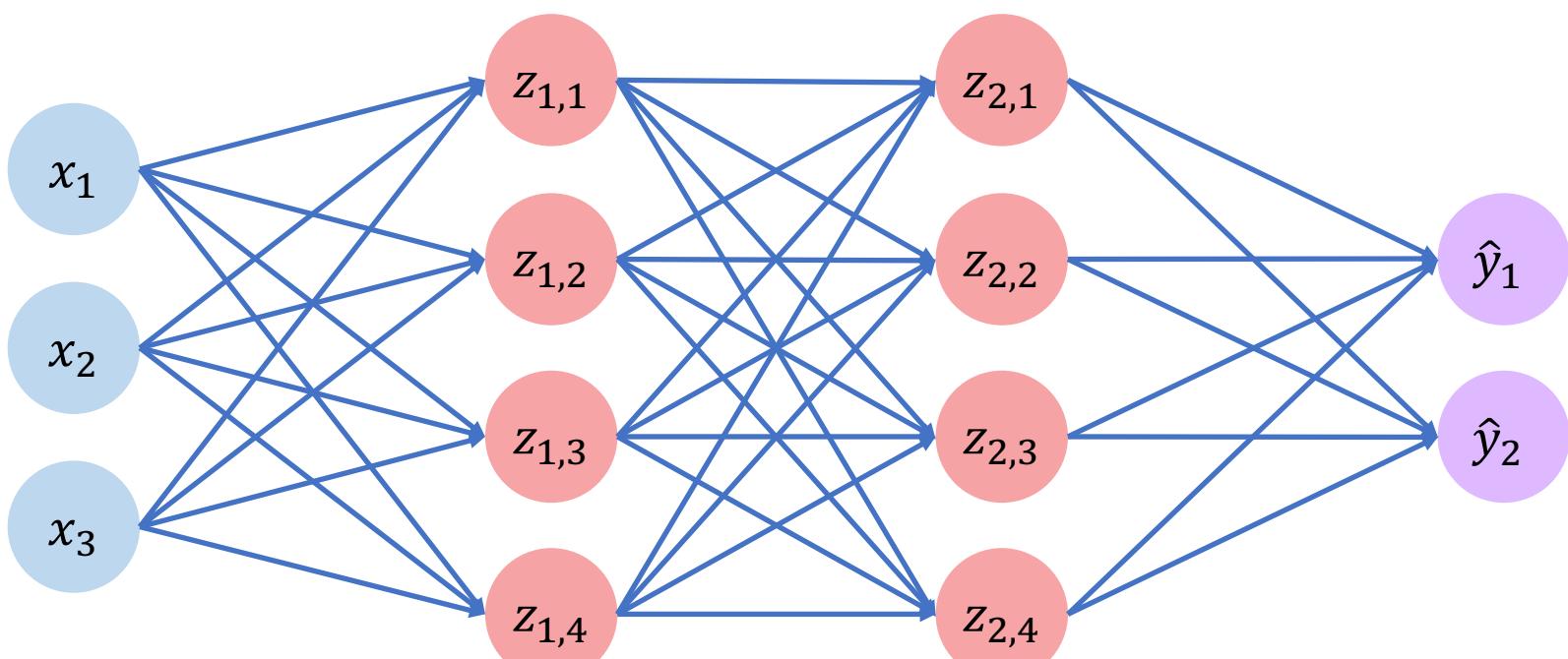
← **Ideal fit** →



Overfitting
Too complex, extra parameters,
does not generalize well

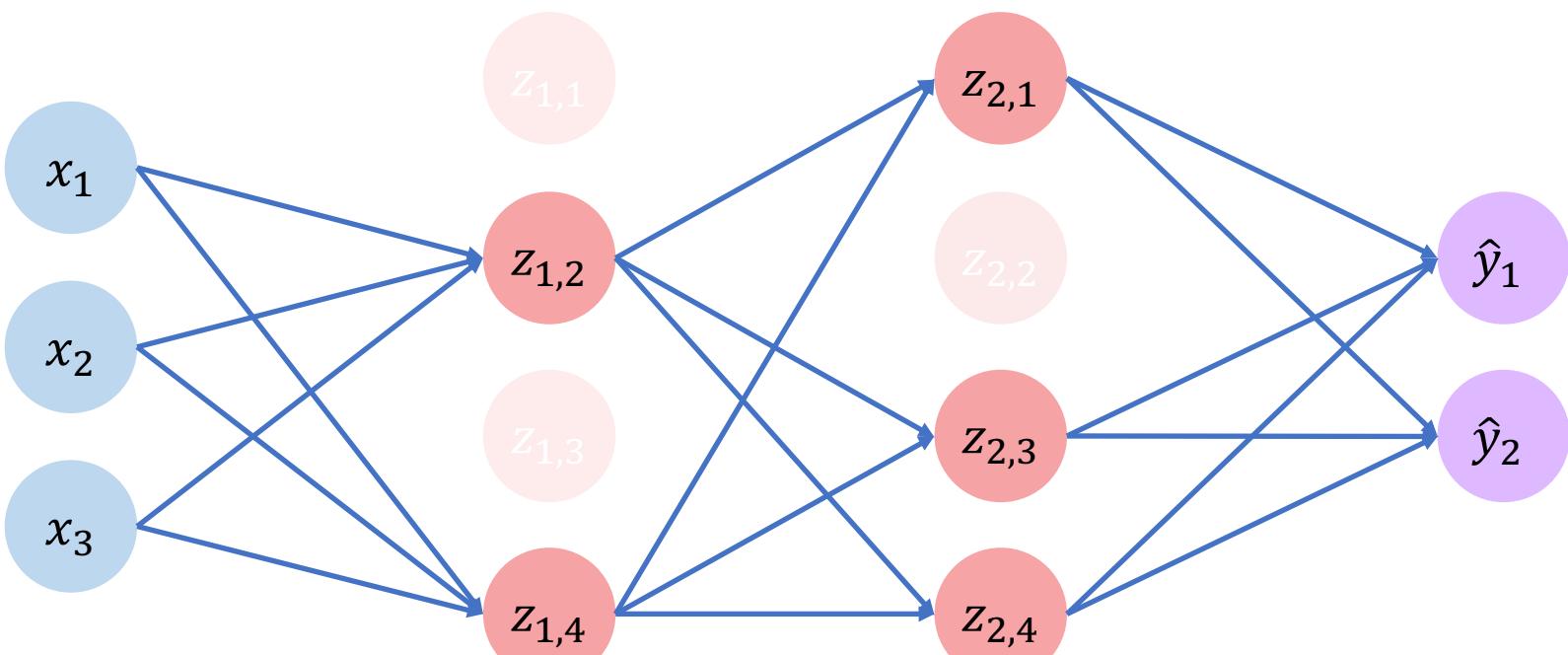
Regularization: dropout

- During training, randomly set some activations to 0



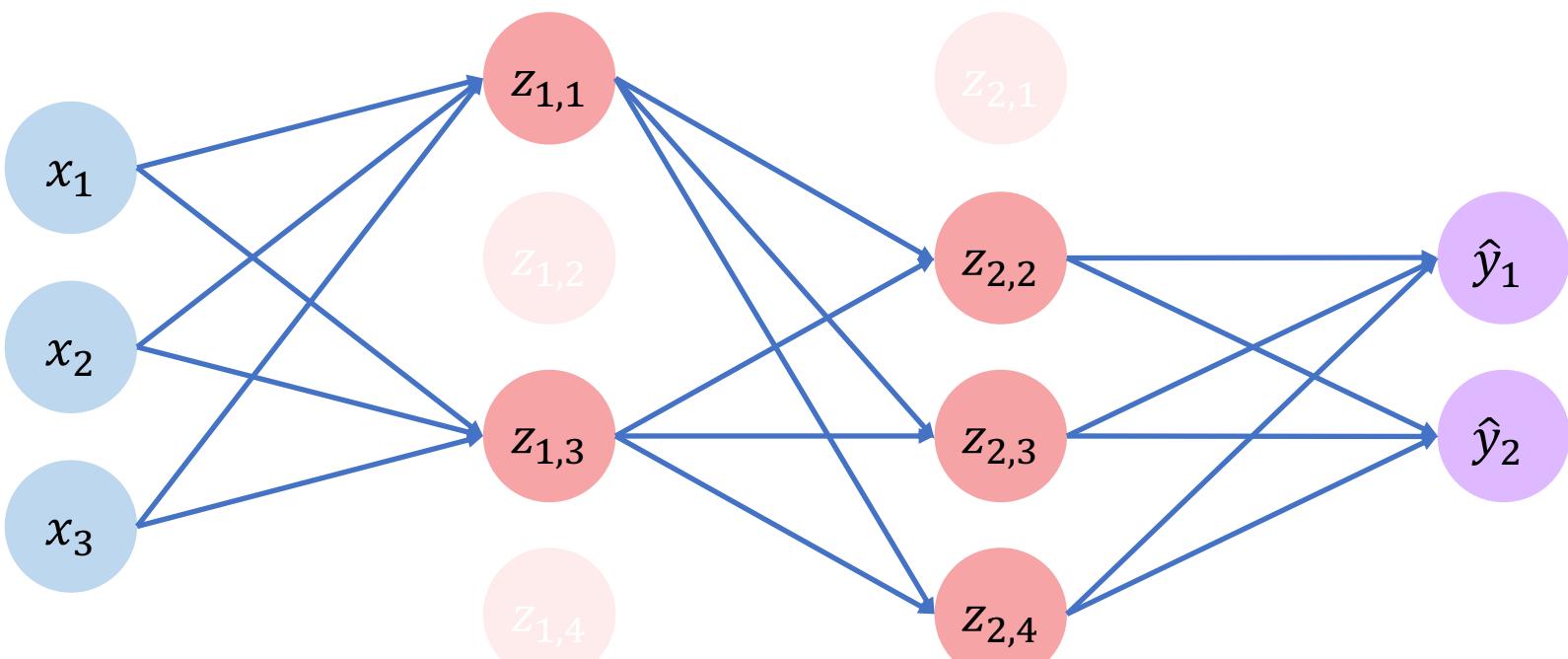
Regularization: dropout

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



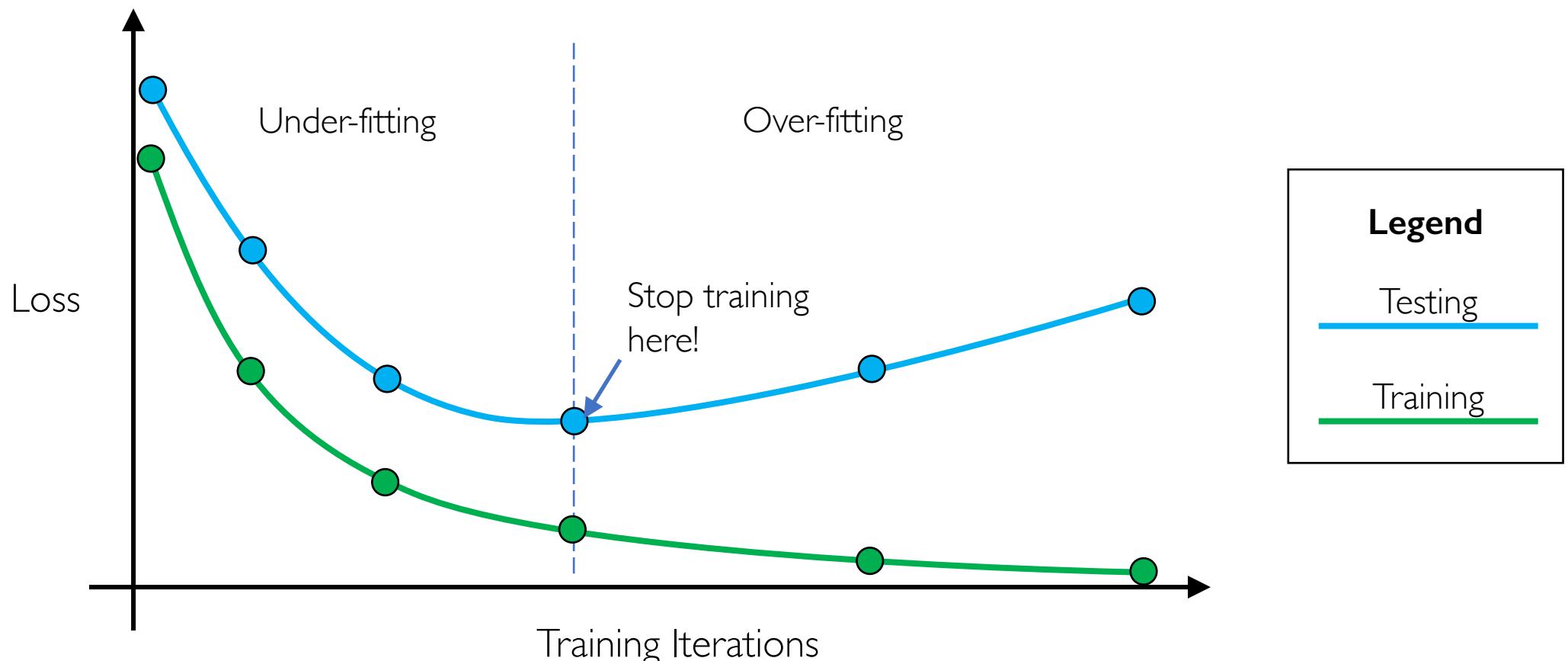
Regularization: dropout

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



Regularization: early stopping

- Stop training before we have a chance to overfit

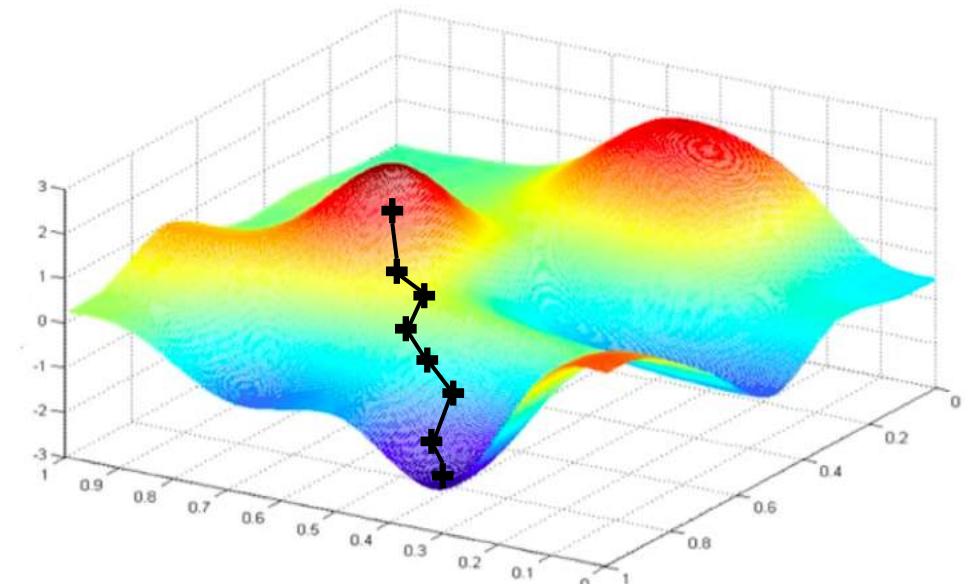


Mini-batches

Algorithm

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

Can be very computational to compute!

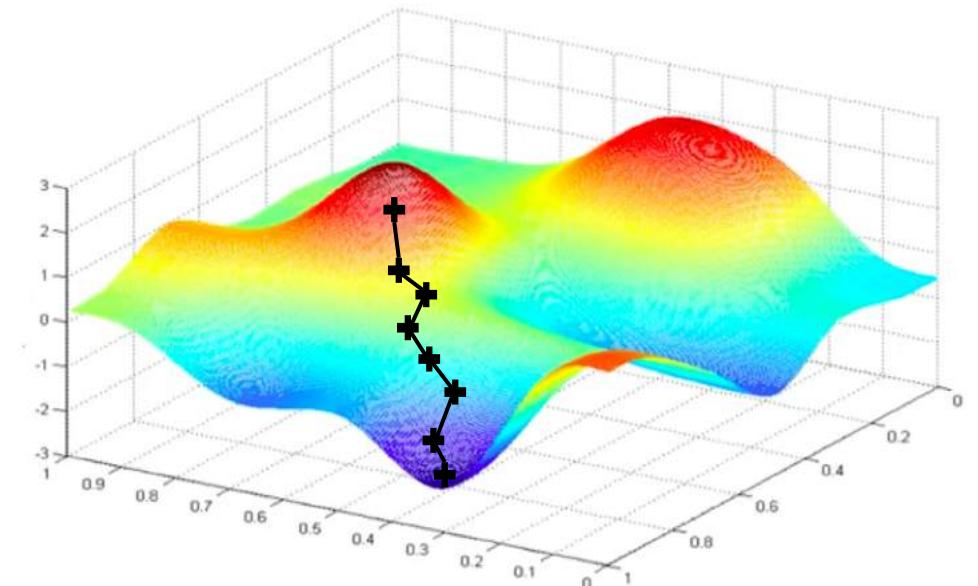


Mini-batches

Algorithm

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

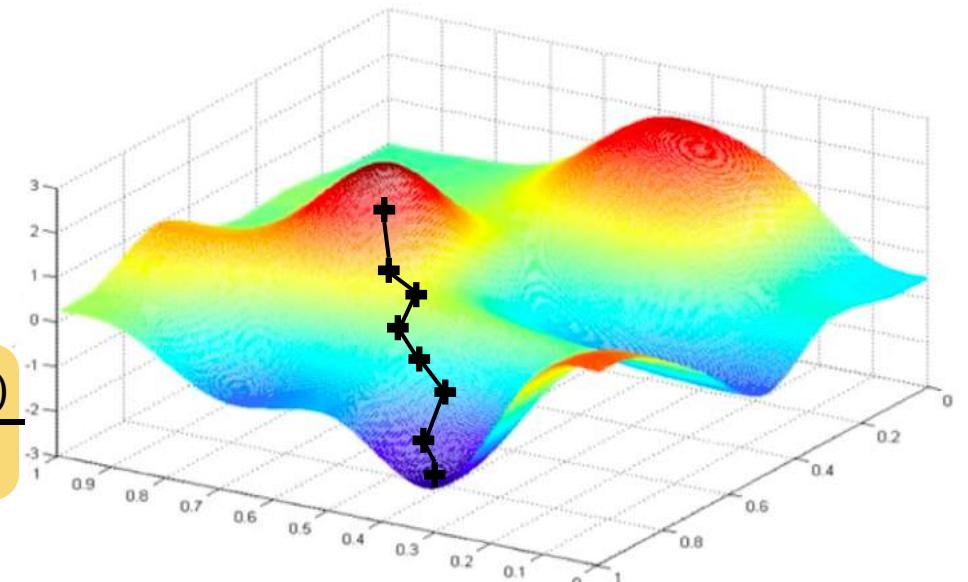
Easy to compute but
very noisy
(stochastic)!



Mini-batches

Algorithm

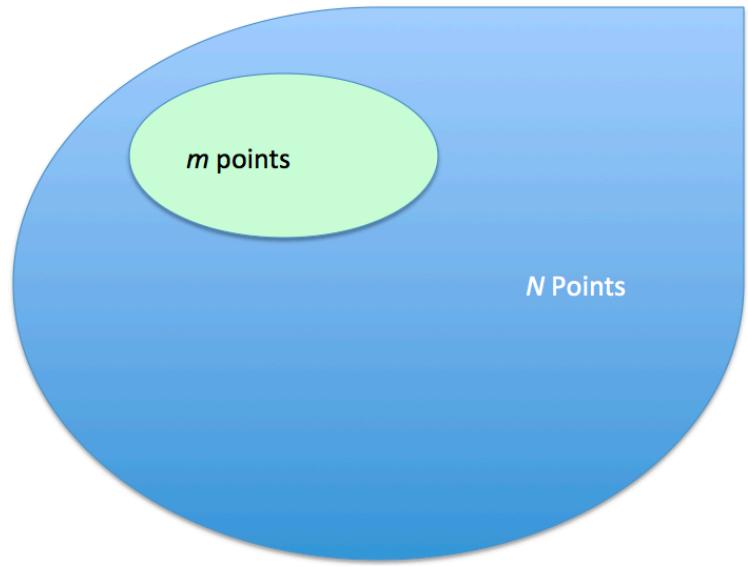
1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient,
$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$$
5. Update weights, $\mathbf{W}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



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

Batch normalization

- Performing normalization for each training mini-batch



Population statistics (with N points)
and batch statistics (with m points)

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.

Optimization

- 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.

Additional details: <https://ruder.io/optimizing-gradient-descent/>

Optimization

$$Vdw=0, Sdw=0, Vdb=0, Sdb=0$$

On iteration t (minibatch)

Compute dw, db using current minibatch.

moment estimate update

$$\begin{cases} Vdw = \beta_1 Vdw + (1-\beta_1) dw, \\ Sdw = \beta_2 Sdw + (1-\beta_2) dw^2 \end{cases}$$

"momentum"

$$\begin{cases} Vdb = \beta_1 Vdb + (1-\beta_1) db, \\ Sdb = \beta_2 Sdb + (1-\beta_2) db^2 \end{cases}$$

"RMSprop"

bias correction

$$\begin{cases} Vdw_{\text{corrected}} = Vdw / (1 - \beta_1^t), \\ Sdw_{\text{corrected}} = Sdw / (1 - \beta_2^t) \end{cases}$$

$$\begin{cases} Vdb_{\text{corrected}} = Vdb / (1 - \beta_1^t), \\ Sdb_{\text{corrected}} = Sdb / (1 - \beta_2^t) \end{cases}$$

Batch gradient descent $W := W - \alpha dw, b := b - \alpha db$

ADAM

$$W := W - \alpha \frac{Vdw_{\text{corrected}}}{\sqrt{Sdw_{\text{corrected}}} + \epsilon}, \quad b := b - \alpha \frac{Vdb_{\text{corrected}}}{\sqrt{Sdb_{\text{corrected}}} + \epsilon}$$

Momentum $W := W - \alpha Vdw_{\text{corrected}}, \quad b := b - \alpha Vdb_{\text{corrected}}$

RMSprop $W := W - \frac{\alpha dw}{\sqrt{Sdw} + \epsilon}, \quad b := b - \frac{\alpha db}{\sqrt{Sdb} + \epsilon}$

α : needs to be tuned

$$\left. \begin{array}{l} \beta_1 = 0.9 \\ \beta_2 = 0.999 \\ \epsilon = 10^{-8} \end{array} \right\} \text{usually fixed}$$

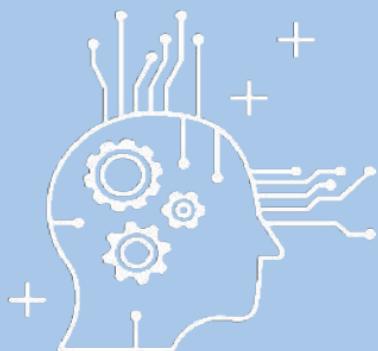
Tutorials

- Google
 - <https://www.tensorflow.org/tutorials/>
- YouTube videos by Siraj Ravel
 - <https://www.youtube.com/watch?v=2FmcHiLCwTU&vl=en>
 - <https://www.youtube.com/watch?v=H1AllrJ- 30>
- Berkeley CS294
 - https://www.youtube.com/watch?v=xZKj7Z1CwHc&list=PLkFD6_40KJlxJMR-j5A1mkxK26gh_qg37&index=23&t=0s
- UNC BIOS735: <https://biodatascience.github.io/statcomp/ml/nl.html>
- ...

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



