

NEURAL NETWORKS

OVERVIEW



What is Deep Learning?

ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



MACHINE LEARNING

Ability to learn without explicitly being programmed

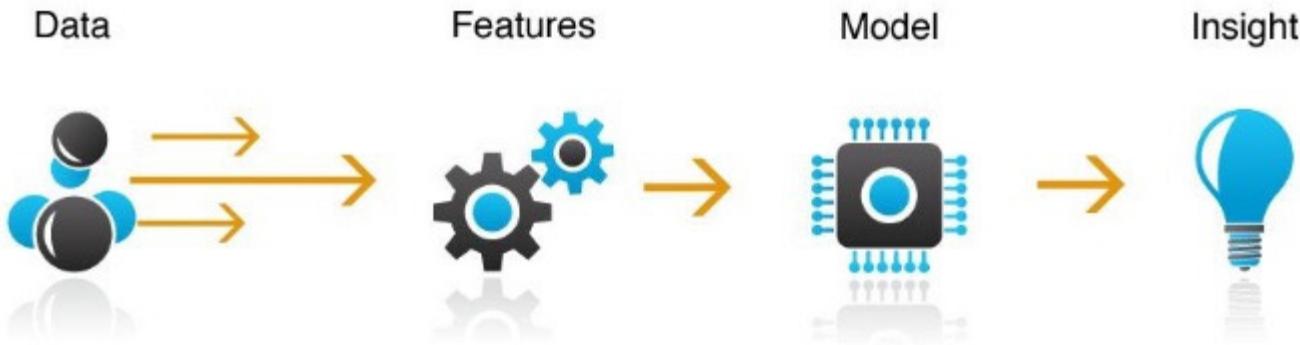


DEEP LEARNING

Extract patterns from data using neural networks



SHALLOW LEARNING



Transforms the input data into **only one** or **two** successive **representation layers**

e.g. SVM: high-dimensional non-linear projections

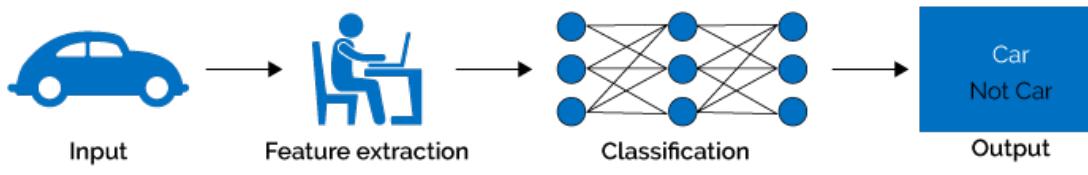
But, complex problems require more refined representations

Feature Engineering = humans manually engineer layers of representations of the data

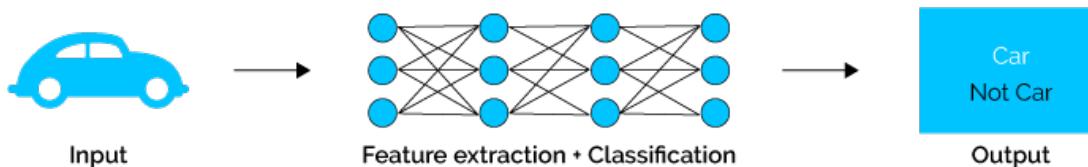
Good features -> solve the problem more elegantly, using fewer resources, less data

DEEP LEARNING

Machine Learning



Deep Learning



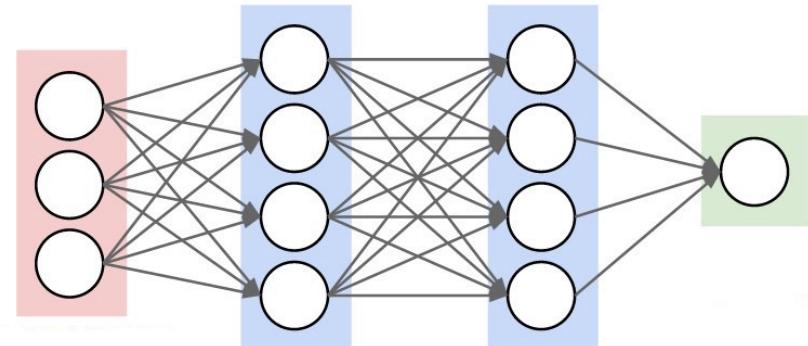
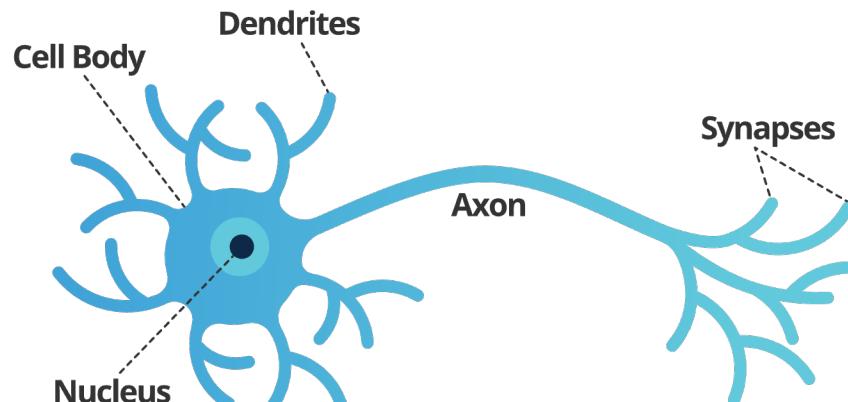
Deep Learning often involves **tens or hundreds** of successive **layers** of representations. Hierarchical representation learning - increasingly complex representations are developed.

These intermediate incremental representations are learned jointly.

Deep Learning completely **automates feature engineering**

Artificial Neural Networks

Brain Neuron Structure



Neural Networks

Good news:

The **universal approximation theorem** states that any continuous function $f : [0, 1]^n \rightarrow [0, 1]$ can be approximated arbitrarily well by a neural network with at least 1 hidden layer with a finite number of weights.

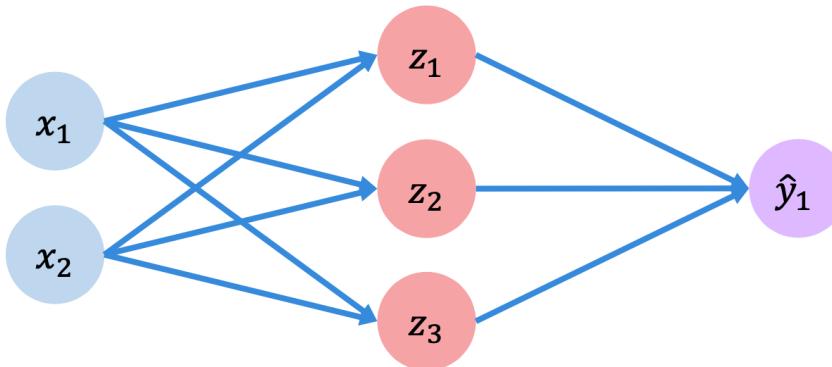
Bad news:

A feedforward network with a single layer is sufficient to represent any function, but **the layer may be infeasibly large and may fail to learn and generalize** correctly.

Common Loss Functions

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

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



$f(x)$	y
30	90
80	20
85	95
\vdots	\vdots

Final Grades
(percentage)

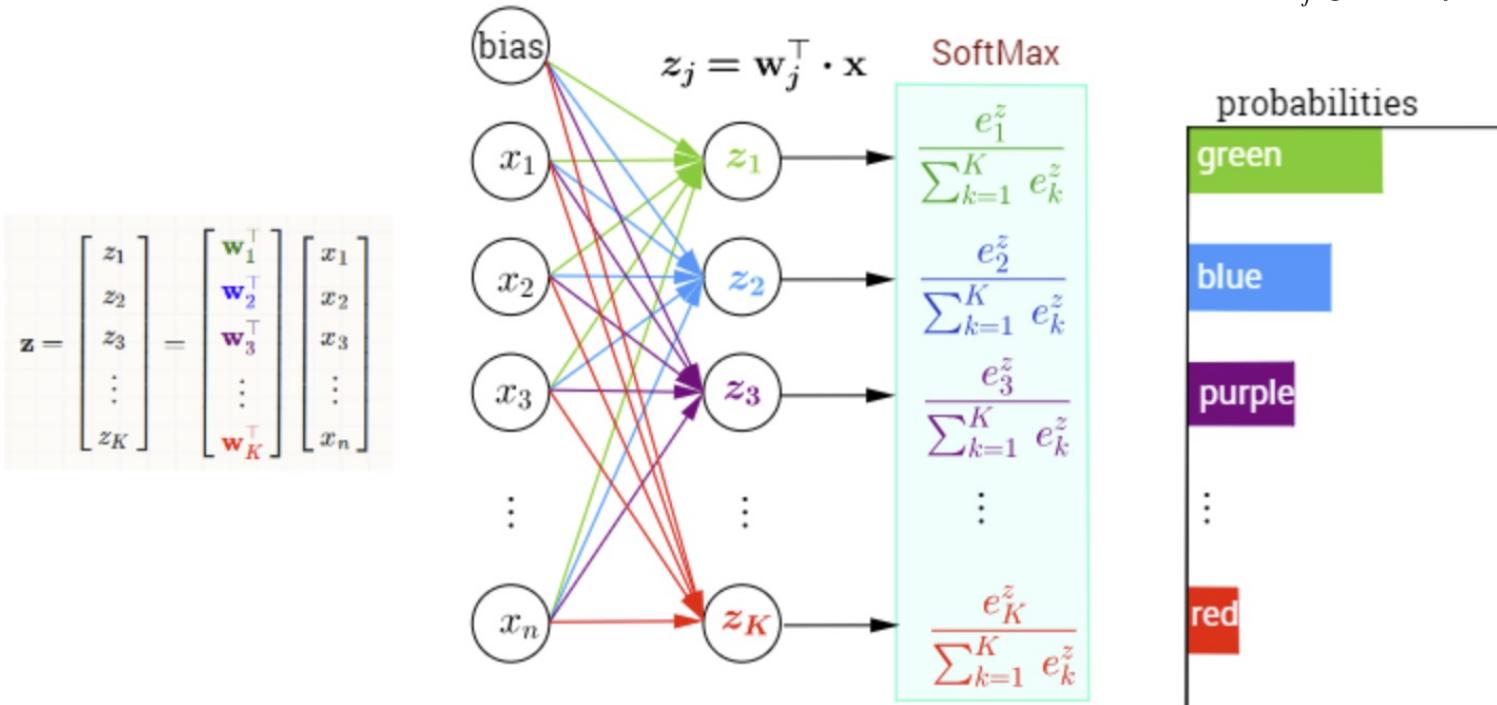
$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \underbrace{(y^{(i)} - f(x^{(i)}; \mathbf{W}))^2}_{\text{Actual Predicted}}$$



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

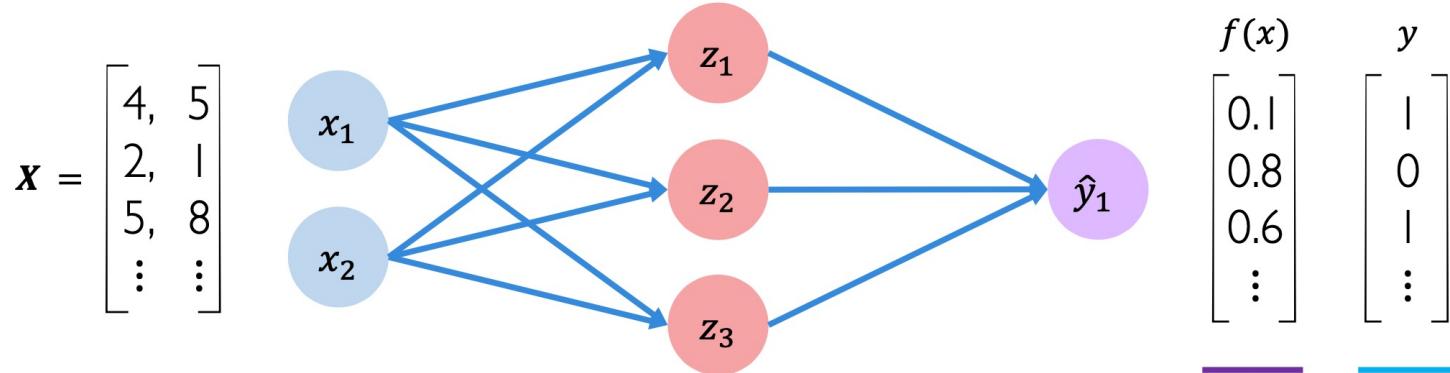
Common Loss Functions

$$\Phi(\vec{v})_i = \frac{\exp(v_i)}{\sum_{j=1}^d \exp(v_j)} \quad \forall i \in \{1 \dots d\}$$



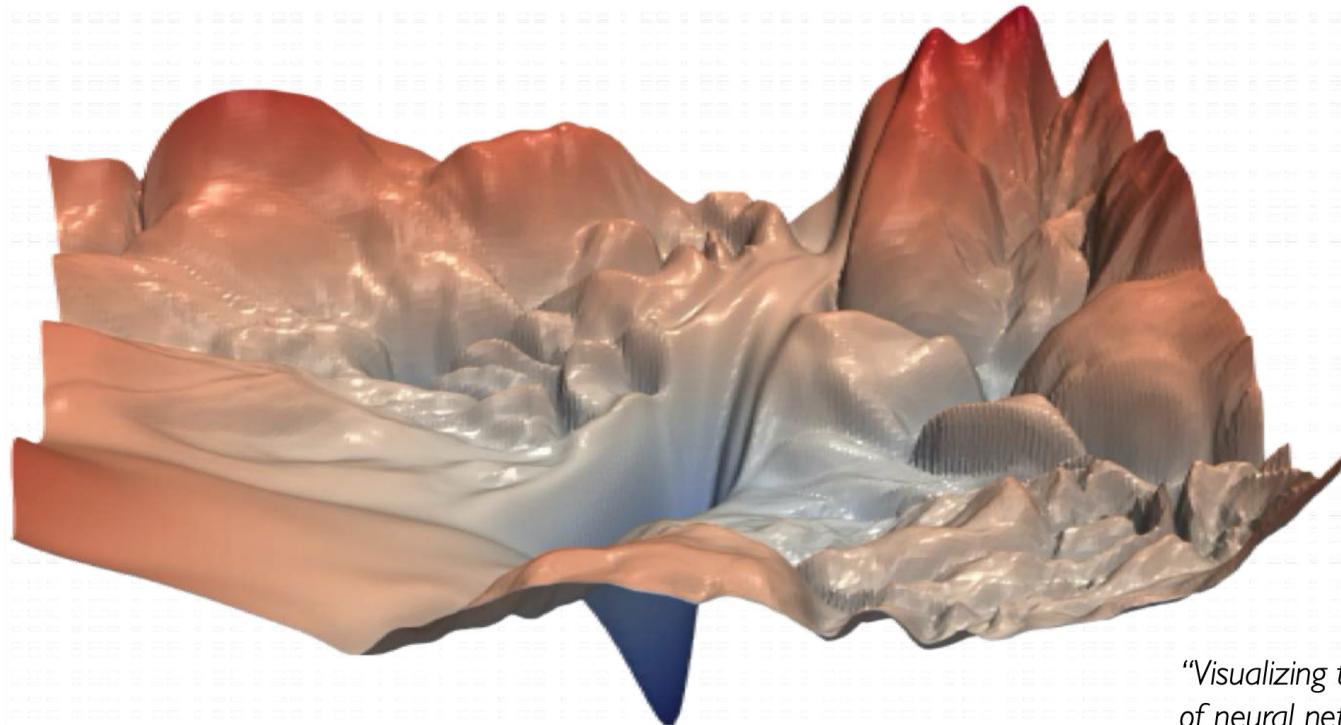
Common Loss Functions

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



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

Loss Functions



*"Visualizing the loss landscape
of neural nets". Dec 2017.*

Training Neural Networks with Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

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

2. Loop until convergence:

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

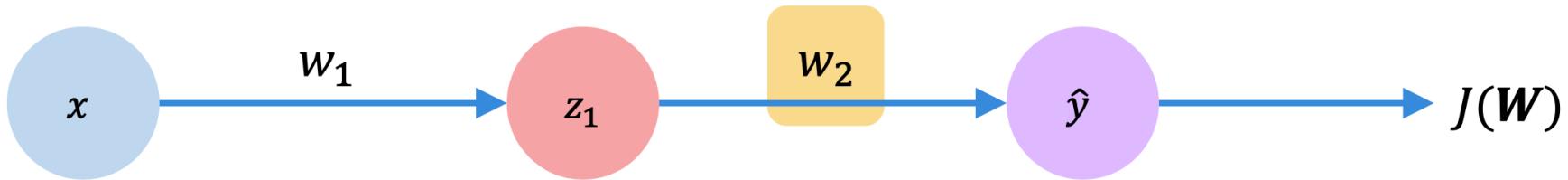
```
 grads = tf.gradients(ys=loss, xs=weights)
```

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

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

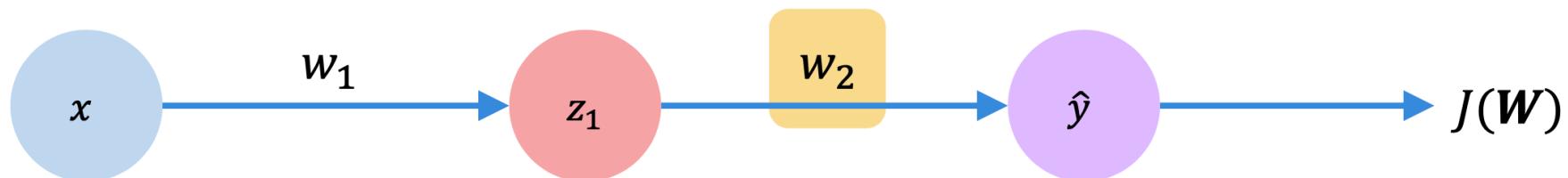
5. Return weights

Training Neural Networks: Backpropagation



How does a small change in one weight (ex. w_2) affect the final loss $J(\mathbf{W})$?

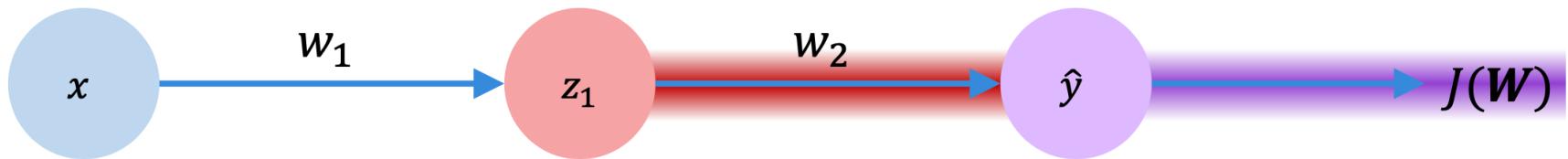
Training Neural Networks: Backpropagation



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

Let's use the chain rule!

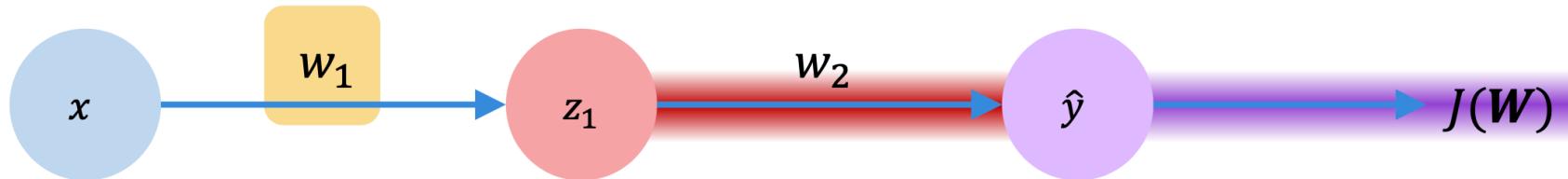
Training Neural Networks: Backpropagation



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



Training Neural Networks: Backpropagation

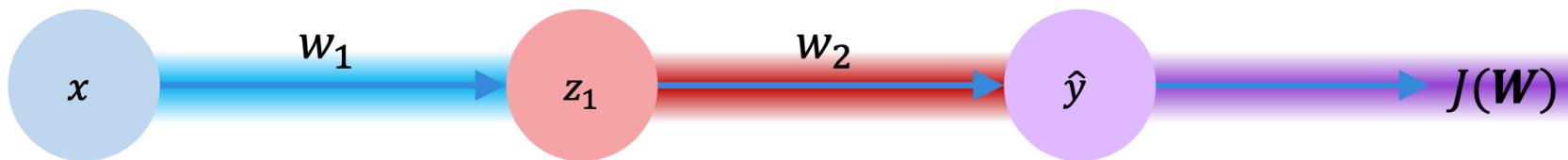


$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

Apply chain rule!

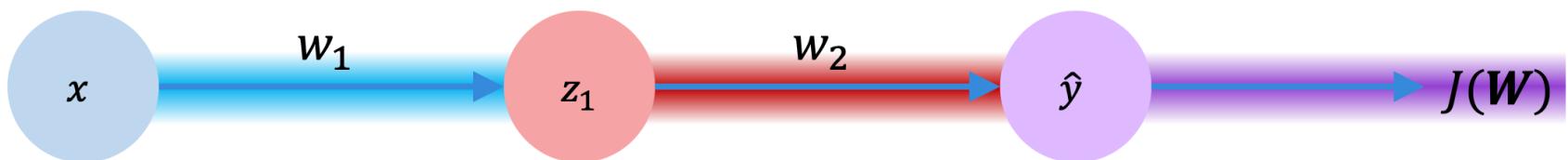
Apply chain rule!

Training Neural Networks: Backpropagation



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

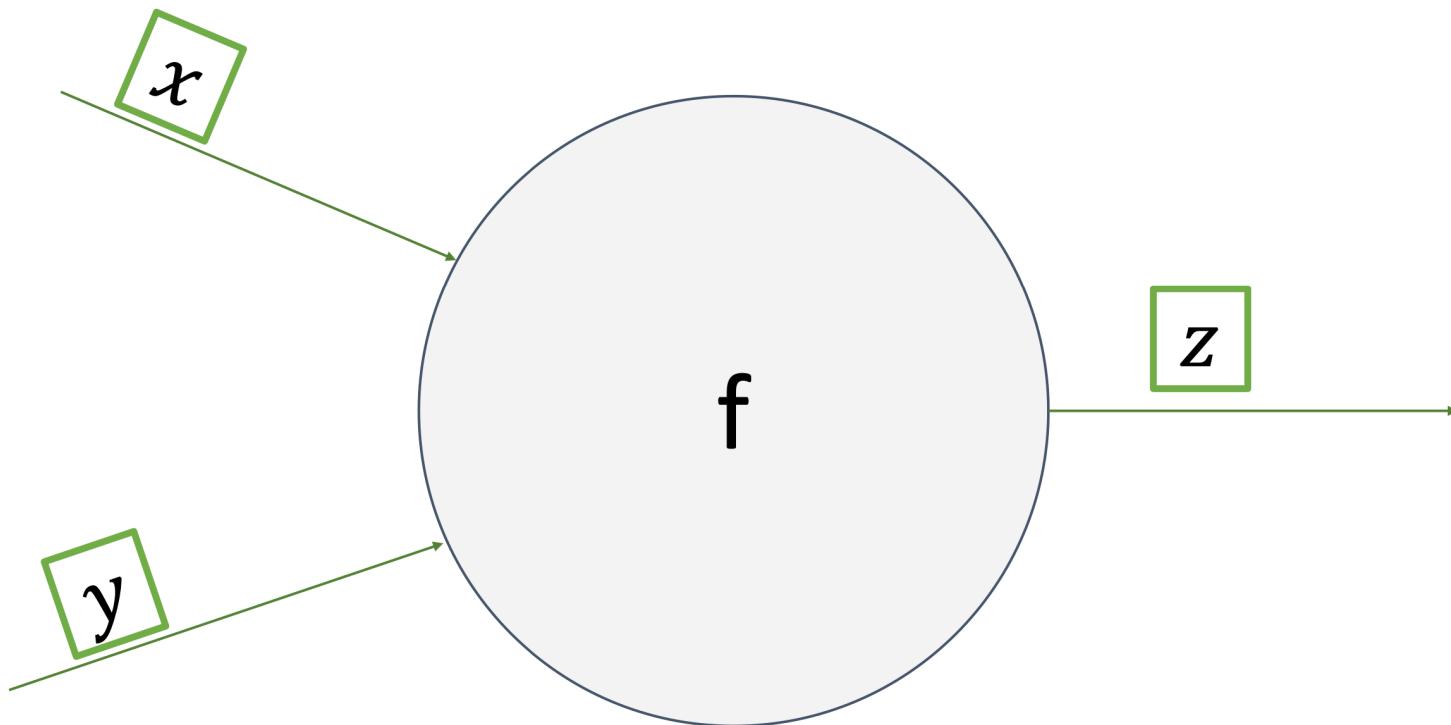
Training Neural Networks: Backpropagation



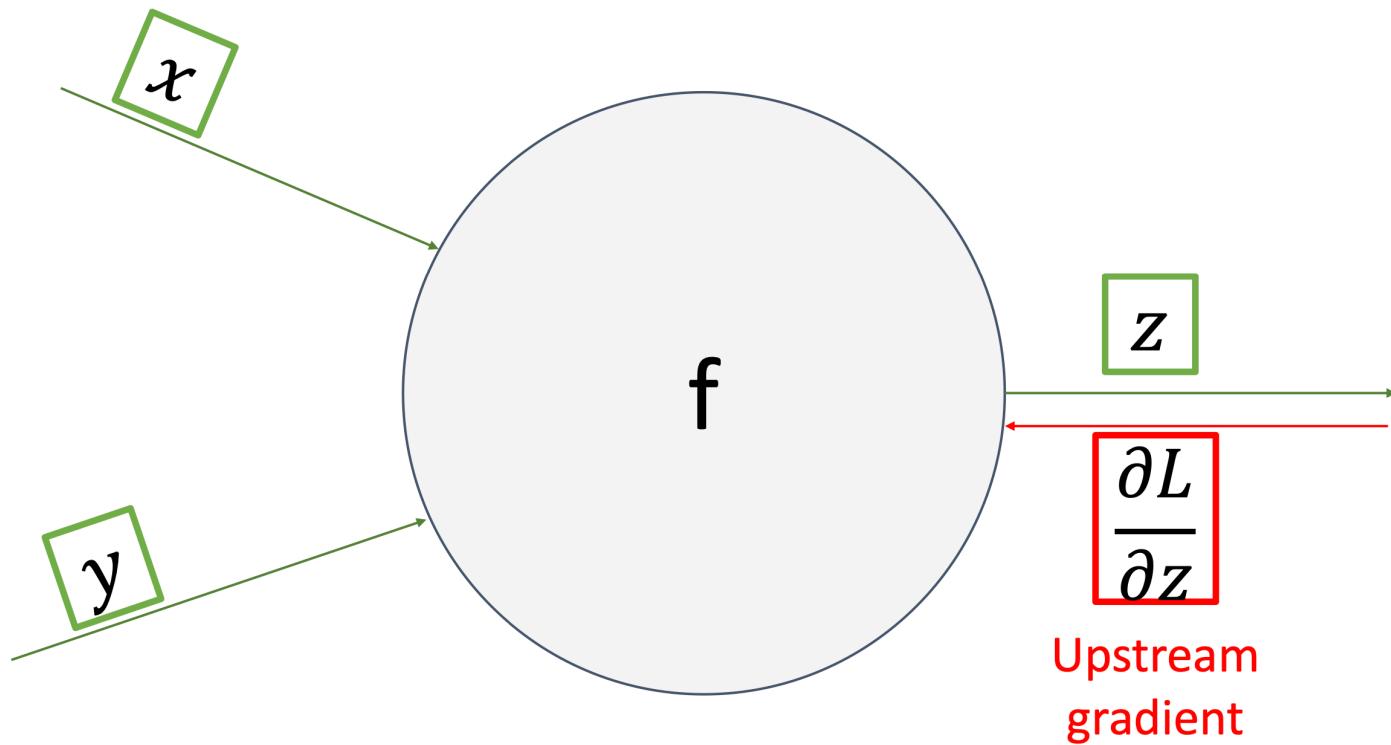
$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underbrace{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue}}$$

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

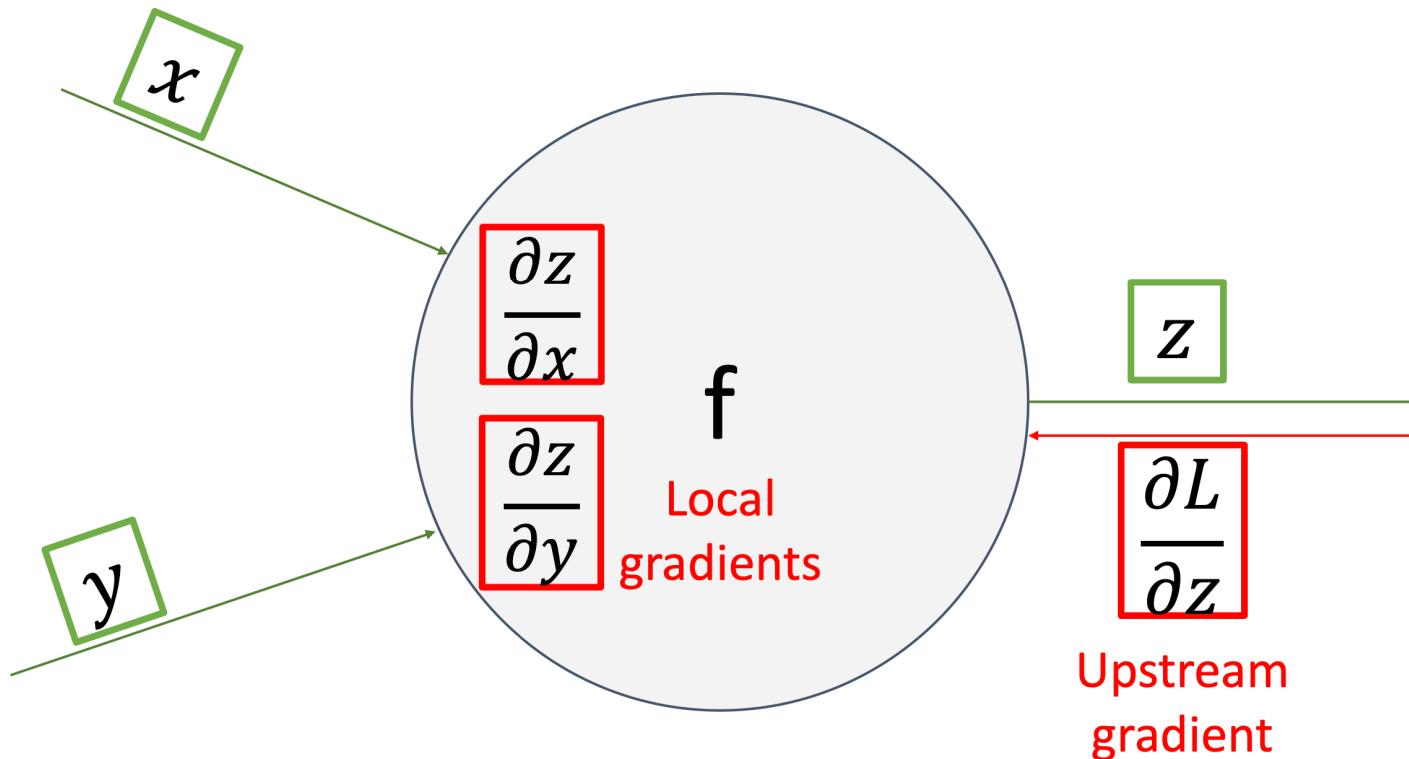
Backpropagation: Upstream Gradient



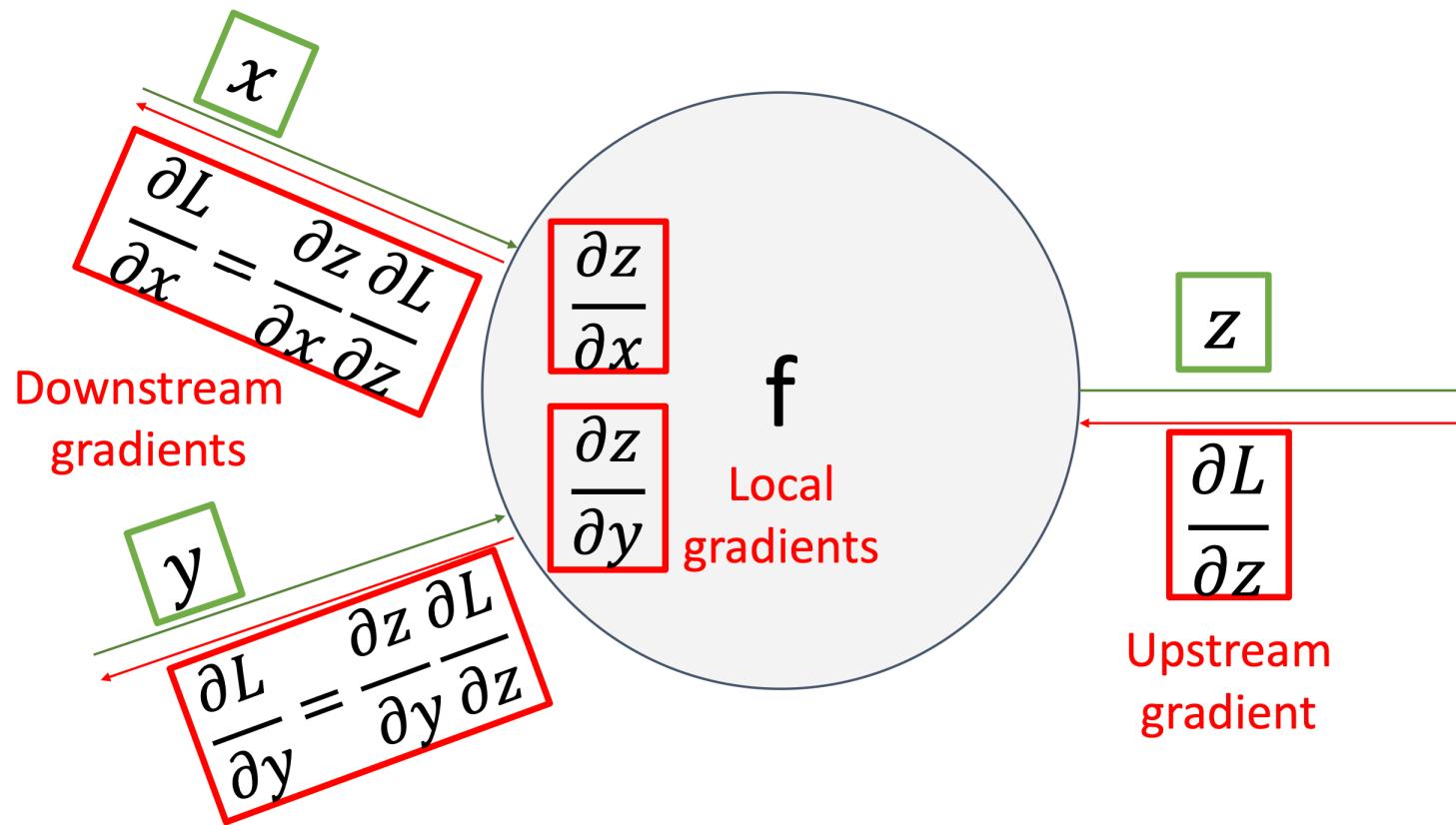
Backpropagation: Upstream Gradient



Backpropagation: Downstream/Upstream Gradient



Backpropagation: Upstream Gradient



Training Neural Networks with Adaptive Learning Rates

Learning rates need not be necessarily fixed

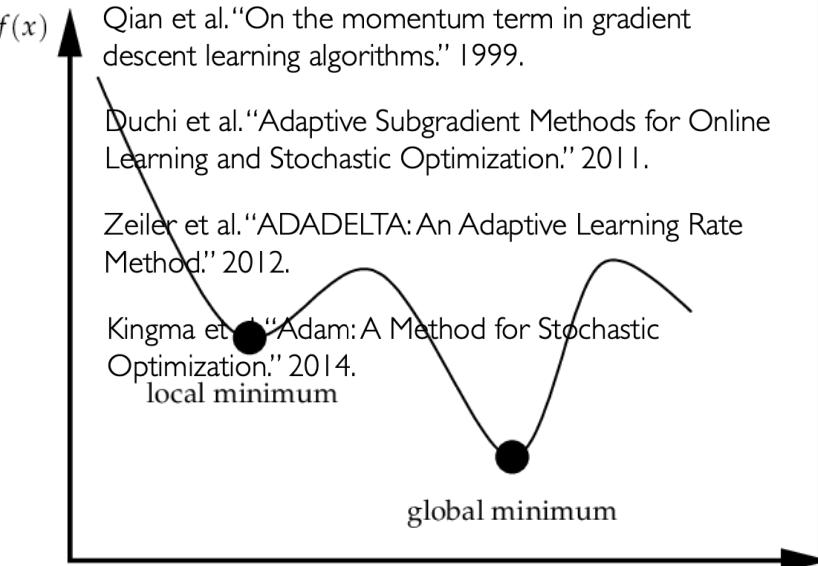
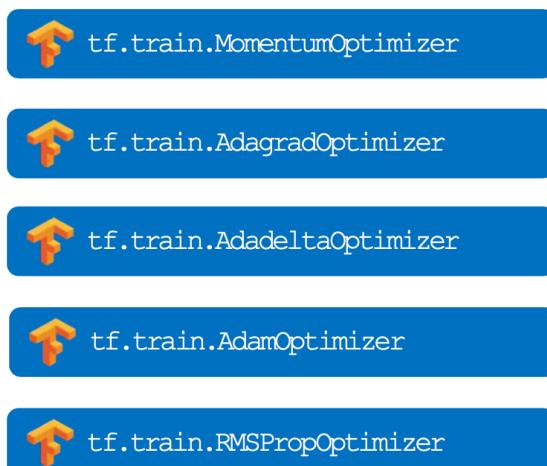
Can be made larger or smaller depending on:

- how large gradient is
- how fast learning is happening
- size of particular weights
- etc...



Training Neural Networks with Adaptive Learning Rates

- Momentum
- Adagrad
- Adadelta
- Adam
- RMSProp



Training Neural Networks with 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

Can be very computational to compute!

Training Neural Networks with Mini-batch/Stochastic Gradient

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
- Mini-batches lead to faster training
 - Can be parallelized for GPUs

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

Sample, Batch, Epoch

- A **sample** is a single row of data. A sample may also be called an **instance**, an **observation**, an **input vector**, or a **feature vector**.
- The **batch size** is a hyperparameter that defines the number of samples to work through before updating the internal model parameters. A training dataset can be divided into one or more batches.
- When *all training samples are used to create one batch*, the learning algorithm is called **batch gradient descent**. When the *batch is the size of one sample*, the learning algorithm is called (true) **stochastic gradient descent**. When the *batch size is more than one sample and less than the size of the training dataset*, the learning algorithm is called **mini-batch gradient descent**. In the case of mini-batch gradient descent, popular batch sizes include 32, 64, 128 and 256 samples.
- The number of **epochs** is a hyperparameter that defines the number times that the learning algorithm will work through the entire training dataset.

Some recent developments..

Gradients without Backpropagation

Atılım Güneş Baydin¹ Barak A. Pearlmutter² Don Syme³ Frank Wood⁴ Philip Torr⁵

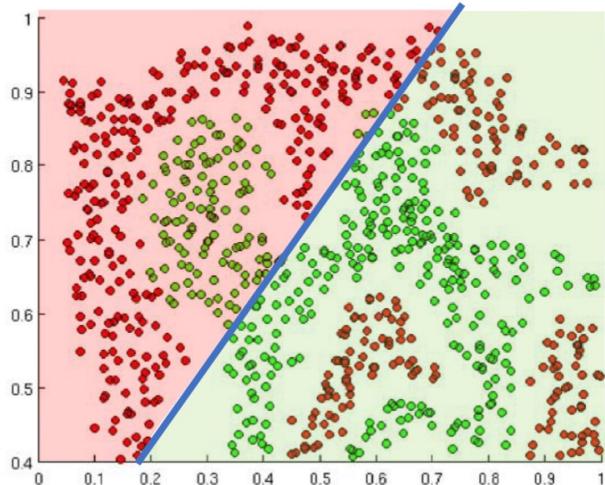
Abstract

Using backpropagation to compute gradients of objective functions for optimization has remained a mainstay of machine learning. Backpropagation, or reverse-mode differentiation, is a special case within the general family of automatic differentiation algorithms that also includes the forward mode. We present a method to compute gradients based solely on the directional derivative that one can compute exactly and efficiently via the forward mode. We call this formulation the **forward gradient**, an unbiased estimate of the gradient that can be evaluated in a single forward run of the function, entirely eliminating the need for backpropagation in gradient descent. We demonstrate forward gradient descent in a range of problems, showing substantial savings in computation and enabling training up to twice as fast in some cases.

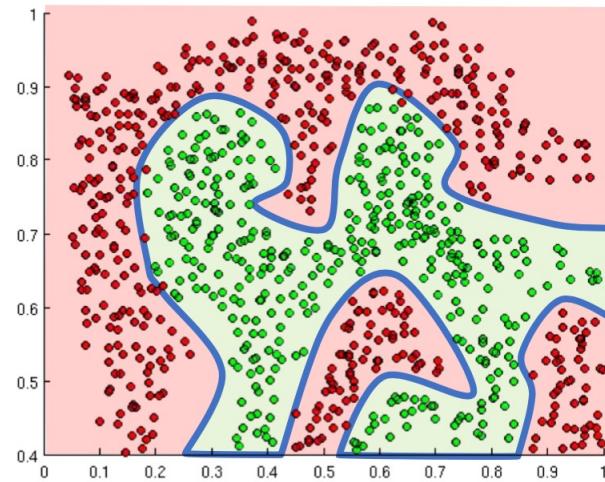
forward–backward algorithm, of which backpropagation is a special case conventionally applied to neural networks. This is mainly due to the central role of scalar-valued objectives in ML, whose gradient with respect to a very large number of inputs can be evaluated exactly and efficiently with a single evaluation of the reverse mode.

Reverse mode is a member of a larger family of AD algorithms that also includes the forward mode (Wengert, 1964), which has the favorable characteristic of requiring only a single forward evaluation of a function (i.e., not involving any backpropagation) at a significantly lower computational cost. Crucially, forward and reverse modes of AD evaluate different quantities. Given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, forward mode evaluates the Jacobian–vector product $J_f v$, where $J_f \in \mathbb{R}^{m \times n}$ and $v \in \mathbb{R}^n$; and reverse mode evaluates the vector–Jacobian product $v^\top J_f$, where $v \in \mathbb{R}^m$. For the case of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ (e.g., an objective function in ML), forward mode gives us $\nabla f \cdot v \in \mathbb{R}$. the directional derivative:

Activation Functions



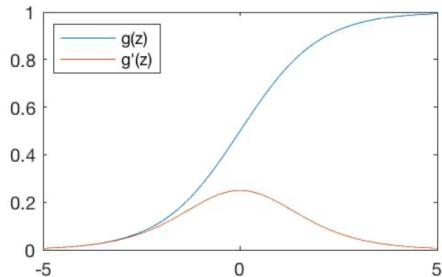
Linear Activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions

Common Activation Functions

Sigmoid Function

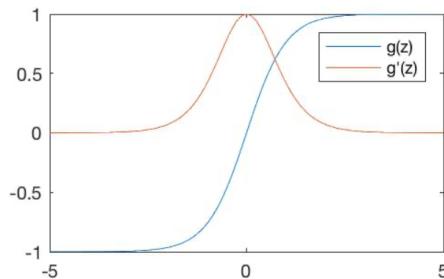


$$g(z) = \frac{1}{1 + e^{-z}}$$

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

 `tf.nn.sigmoid(z)`

Hyperbolic Tangent

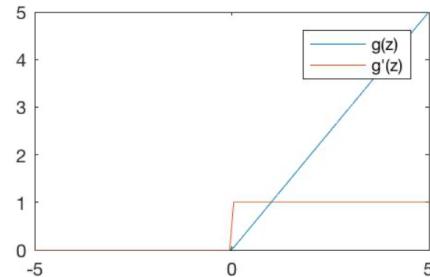


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

 `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

Epoch
000,000Learning rate
0.03Activation
TanhRegularization
NoneRegularization rate
0Problem type
Classification<http://playground.tensorflow.org/>

DATA

Which dataset do you want to use?



Ratio of training to test data: 50%

Noise: 0

Batch size: 10

REGENERATE

FEATURES

Which properties do you want to feed in?



+ - 2 HIDDEN LAYERS

4 neurons

-

2 neurons

+

+

-

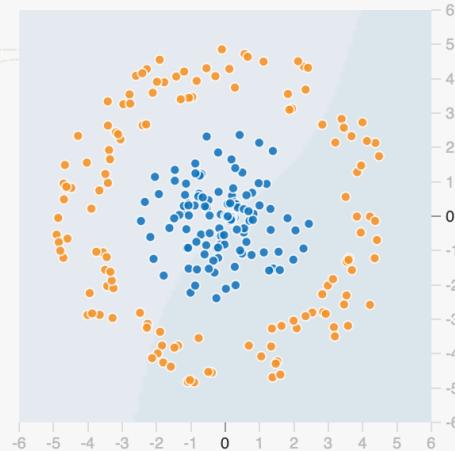
The outputs are mixed with varying **weights**, shown by the thickness of the lines.

This is the output from one **neuron**. Hover to see it larger.

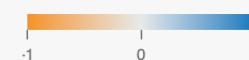
OUTPUT

Test loss 0.495

Training loss 0.510

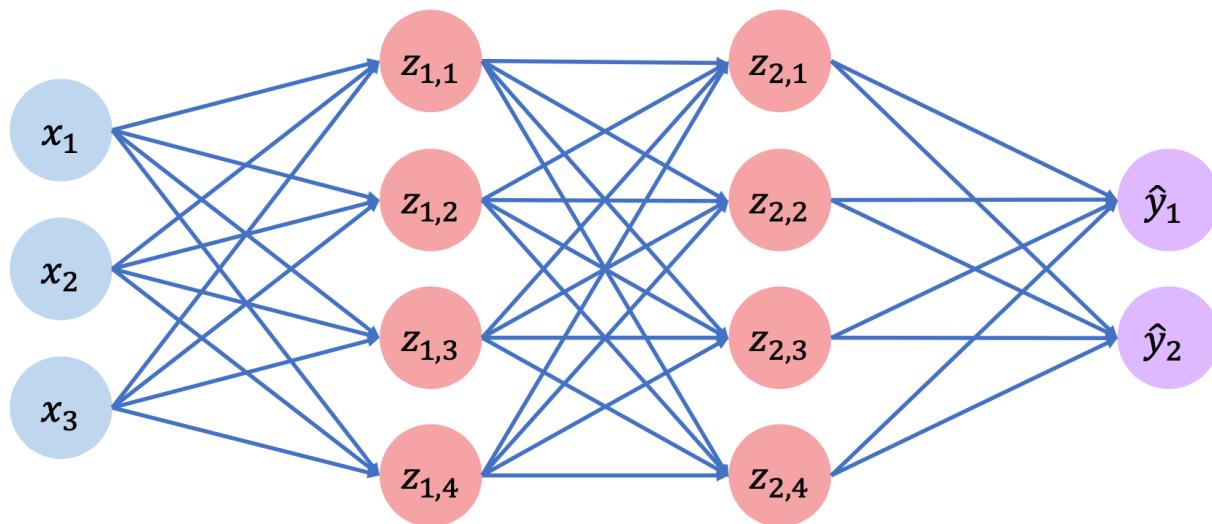


Colors shows data, neuron and weight values.

 Show test data Discretize output

Regularization: Dropout

- During training, randomly set some activations to 0

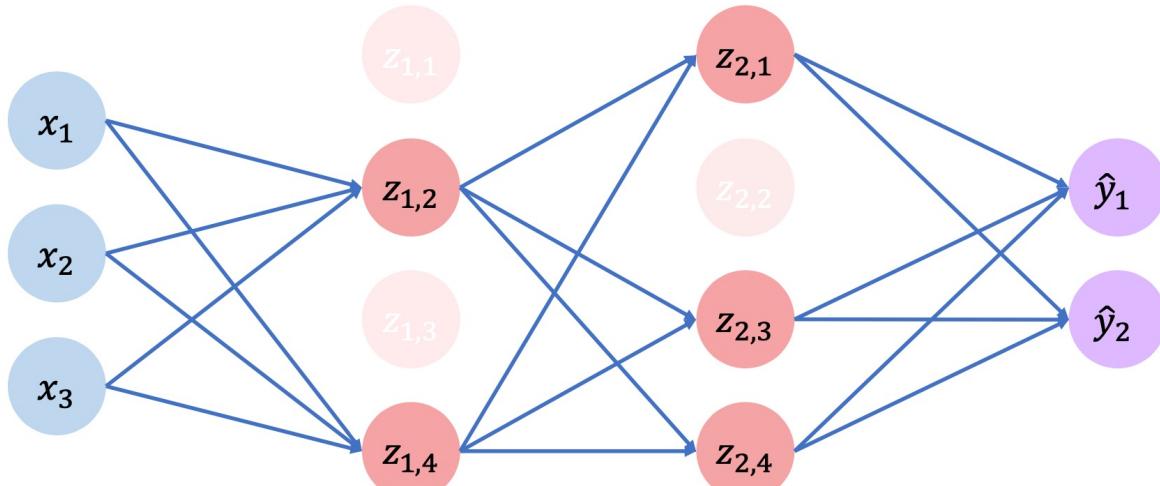


Regularization: Dropout

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

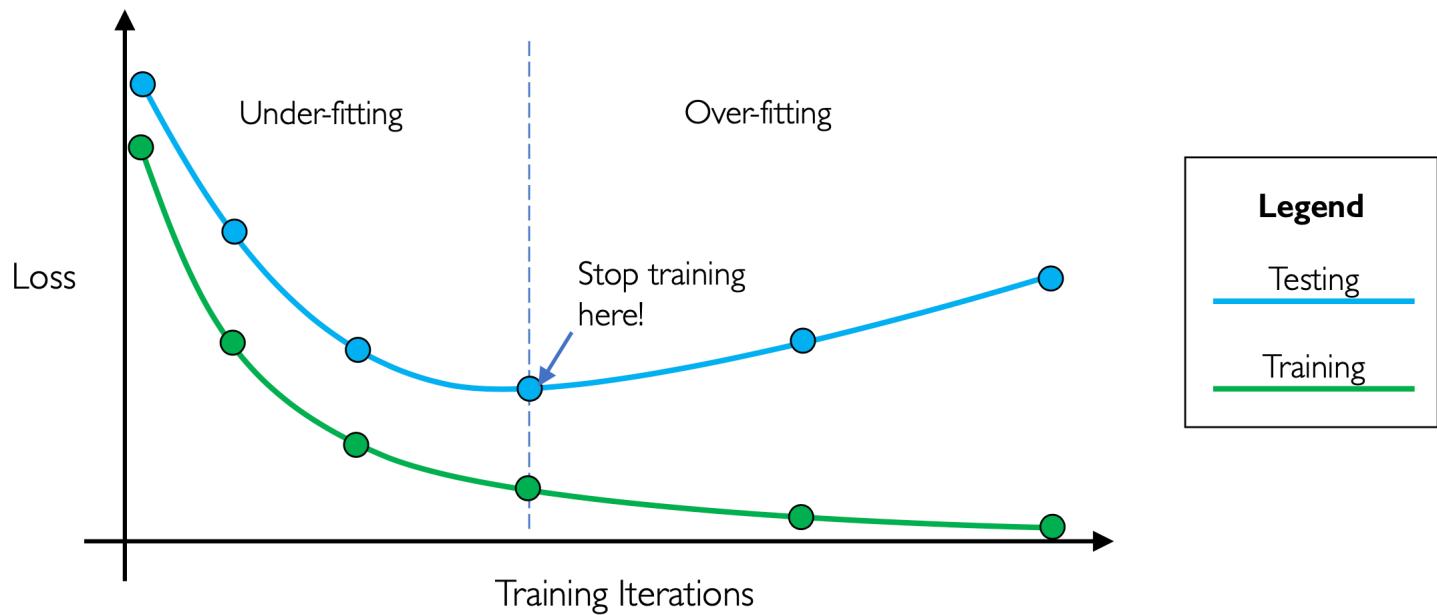


tf.keras.layers.Dropout (p=0.5)



Regularization: Early Stopping

- Stop training before we have a chance to overfit



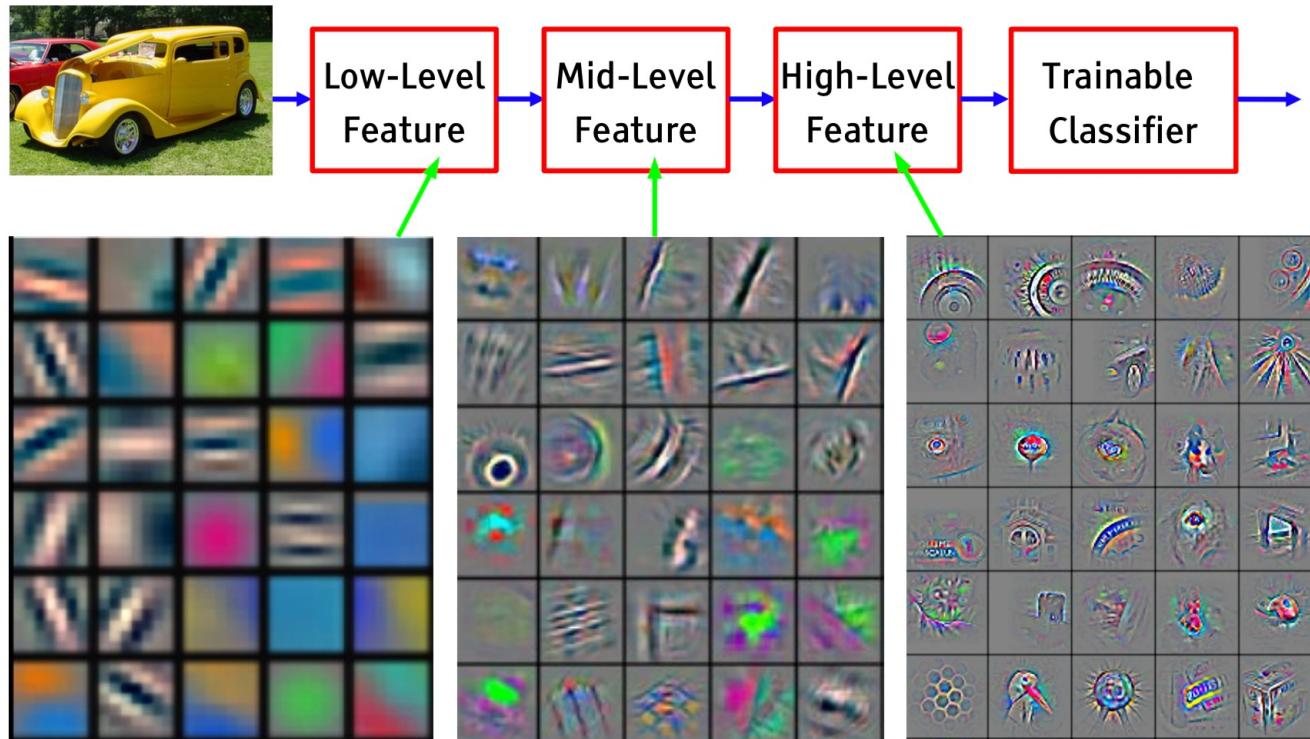
What training typically does not do

Choice of the hyper-parameters has to be done manually:

- Type of activation function
- Choice of architecture (how many hidden layers, their sizes)
- Learning rate, number of training epochs
- What features are presented at the input layer
- How to regularize

It may seem complicated at first - the best way to start is to re-use some existing setup and try your own modifications.

Convolutional Neural Networks (CNNs)



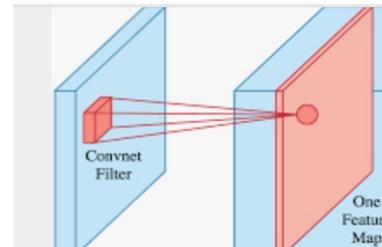
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Convolutional Neural Networks (CNNs)

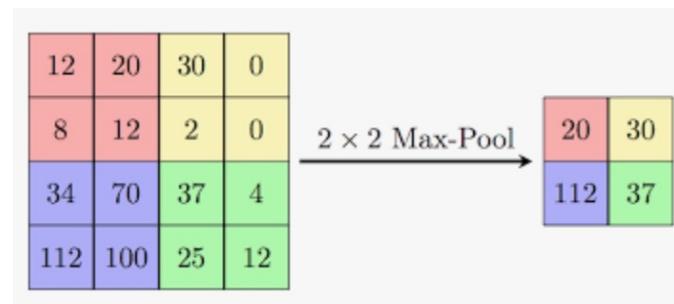
- Practical for computer vision applications
- Hidden layer neurons are not fully connected
- Reduce number of weights (better to avoid overfitting)
- Learn same features irrespective of location in image (edges, eyes, faces...) - concept of weight sharing
- Use same Learning rules (Backpropagation and gradient descent)

Main Layers in a CNN

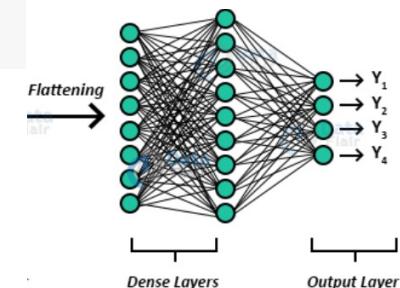
- Convolution Layer



- Pooling Layer



- Dense Layer (Fully-connected layer)



Could be used in classification tasks such as
Sentiment Analysis, Spam Detection or Topic Categorization

Summary

Neural Networks reviewed:

- Loss functions
- Backpropagation
- Activation Functions
- Regularization
- CNNs



**Thanks!
Questions...**

