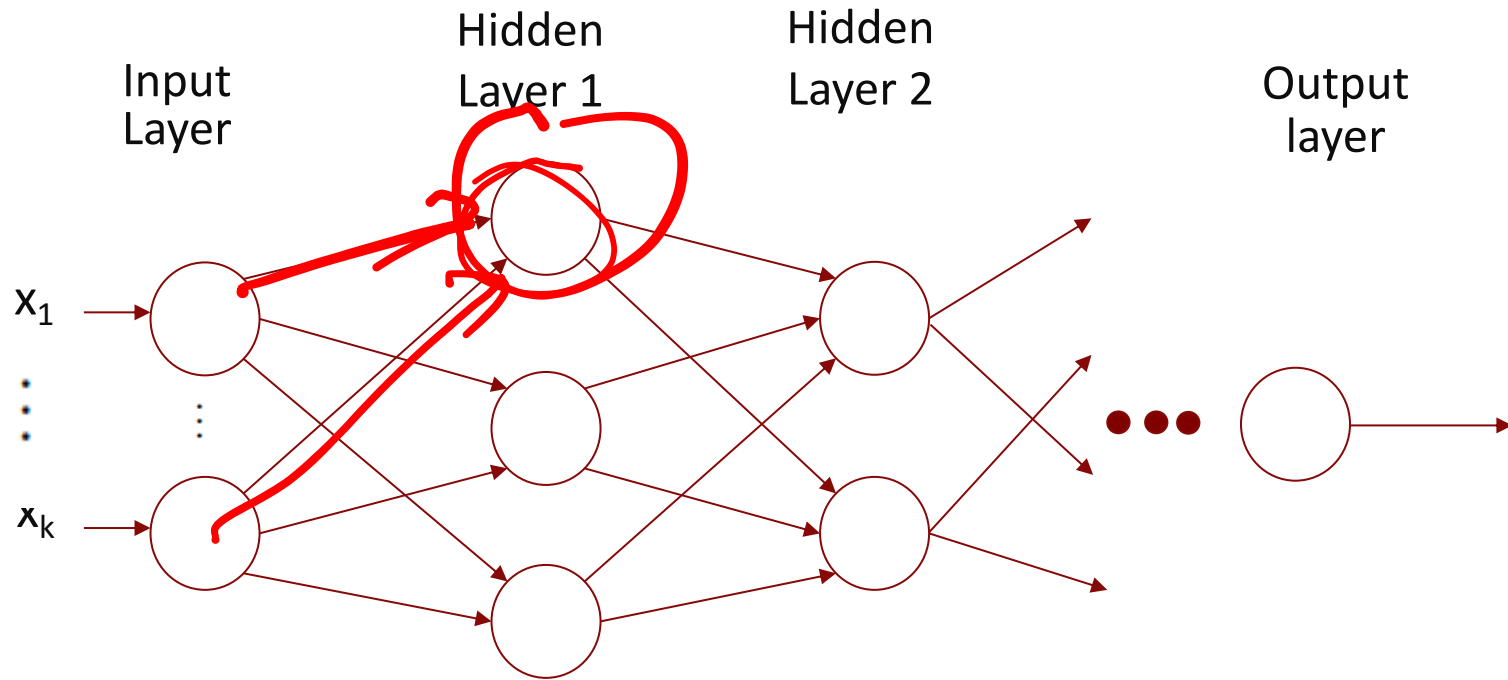
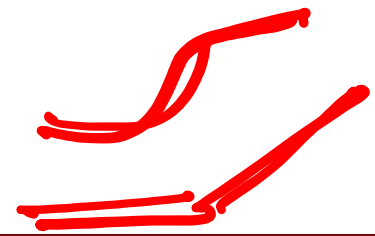


Lecture 2:

Training Deep Neural Networks

Very fractured

Recap: Designing a DNN



User chooses the **# of hidden layers**, **# units in each layer**, the **activation function(s)** for the hidden layers and for the output layer

Application: Predicting heart disease

Predicting Heart Disease

Using a dataset of patients made available by the Cleveland Clinic, we will build our first DL model to predict if a patient has been diagnosed with heart disease from demographics and bio-markers

x_0, x_1, x_2, x_3, x_4 ←

Column	Description	Feature Type
Age	Age in years	Numerical
Sex	(1 = male; 0 = female)	Categorical
CP	Chest pain type (0, 1, 2, 3, 4)	Categorical
Trestbpd	Resting blood pressure (in mm Hg on admission)	Numerical
Chol	Serum cholesterol in mg/dl	Numerical
FBS	fasting blood sugar in 120 mg/dl (1 = true; 0 = false)	Categorical
RestECG	Resting electrocardiogram results (0, 1, 2)	Categorical
Thalach	Maximum heart rate achieved	Numerical
Exang	Exercise induced angina (1 = yes; 0 = no)	Categorical
Oldpeak	ST depression induced by exercise relative to rest	Numerical
Slope	Slope of the peak exercise ST segment	Numerical
CA	Number of major vessels (0-3) colored by fluoroscopy	Both numerical & categorical
Thal	3 = normal; 6 = fixed defect; 7 = reversible defect	Categorical
Target	Diagnosis of heart disease (1 = true; 0 = false)	Target

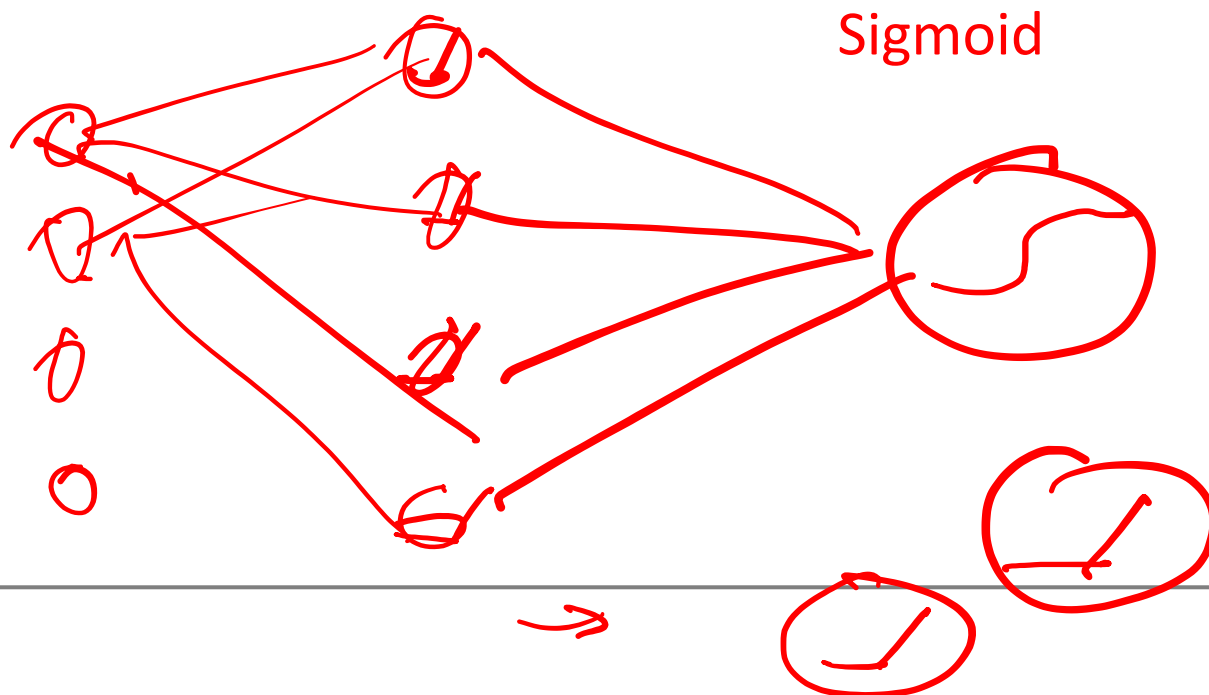
What we want to predict

Let's design our NN

- We design i.e., “lay out” the network
 - Choose *the number of hidden layers* and *the number of ‘neurons’ in each layer*
 - Pick the *right output layer* based on the type of the output

Let's design our NN

- We design i.e., “lay out” the network
 - Choose **the number of hidden layers** and **the number of ‘neurons’ in each layer** 1 hidden layer with 16 ReLU neurons
 - Pick the **right output layer** based on the type of the output



Let's visualize this NN

Input
Layer

X_1



X_{29}

Column	Description	Feature Type
Age	Age in years	Numerical
Sex	(1 = male; 0 = female)	Categorical
CP	Chest pain type (0, 1, 2, 3, 4)	Categorical
Trestbpd	Resting blood pressure (in mm Hg on admission)	Numerical
Chol	Serum cholesterol in mg/dl	Numerical
FBS	fasting blood sugar in 120 mg/dl (1 = true; 0 = false)	Categorical
RestECG	Resting electrocardiogram results (0, 1, 2)	Categorical
Thalach	Maximum heart rate achieved	Numerical
Exang	Exercise induced angina (1 = yes; 0 = no)	Categorical
Oldpeak	ST depression induced by exercise relative to rest	Numerical
Slope	Slope of the peak exercise ST segment	Numerical
CA	Number of major vessels (0-3) colored by fluoroscopy	Both numerical & categorical
Thal	3 = normal; 6 = fixed defect; 7 = reversible defect	Categorical
Target	Diagnosis of heart disease (1 = true; 0 = false)	Target

There are only 13 input variables but some of them are categorical so we one-hot-encode them, resulting in 29 inputs (details in colab).

Let's visualize this NN

Input
Layer

x_1

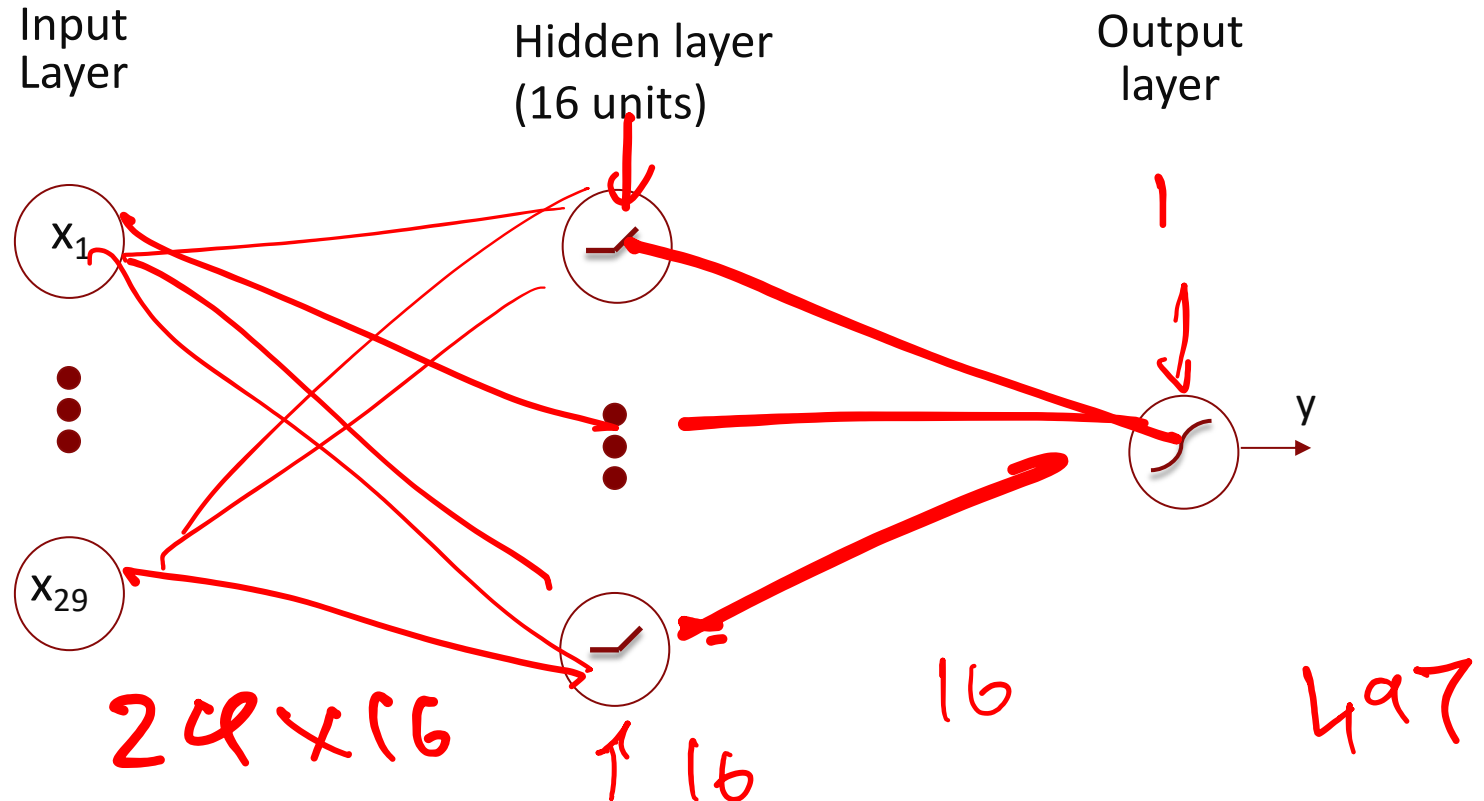
⋮

x_{29}

Hidden layer
(16 units)

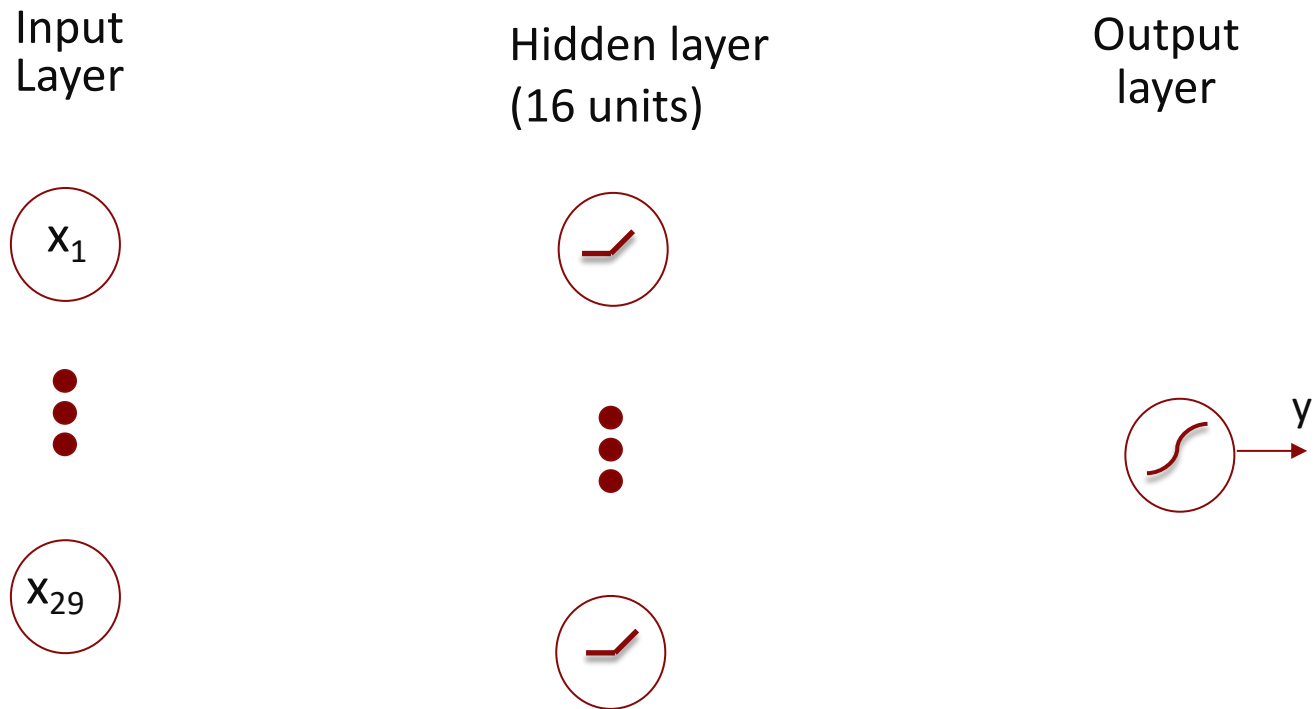
Output
layer

Let's visualize this NN




How many parameters (i.e., weights and biases) does this network have?

Let's visualize this NN



How many parameters (i.e., weights and biases) does this network have?

497

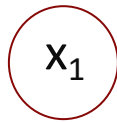


We will now “translate” this network into Keras code to demonstrate how easy it is.

We will give a fuller intro to Keras/Tensorflow and train this model in Colab a bit later.

We first define the input layer in Keras

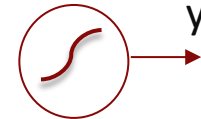
Input
Layer



Hidden layer
(16 units)



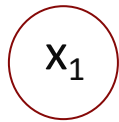
Output
layer



```
input = keras.Input(shape=29)
```

Next, we define the hidden layer

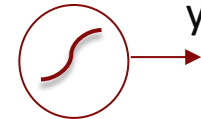
Input
Layer



Hidden layer
(16 units)



Output
layer



```
input = keras.Input(shape=29)
```

```
keras.layers.Dense(16, activation="relu")
```

Since this layer is fully connected to the previous and later layers, we use 'Dense'

Input
Layer

x_1

⋮

x_{29}

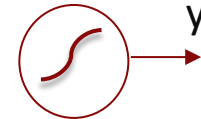
Hidden layer
(16 units)



⋮



Output
layer

 y

```
input = keras.Input(shape=29)
```

```
keras.layers.Dense(16, activation="relu")
```

We specify the number of neurons we want in this layer ...

Input
Layer

x_1

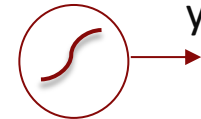


x_{29}

Hidden layer
(16 units)



Output
layer



```
input = keras.Input(shape=29)
```

```
keras.layers.Dense(16, activation="relu")
```

... and the activation function

Input
Layer

x_1

⋮

x_{29}

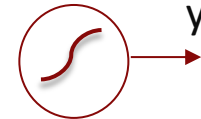
Hidden layer
(16 units)



⋮



Output
layer



```
input = keras.Input(shape=29)
```

```
keras.layers.Dense(16, activation="relu")
```


Next, we “feed” the input to this layer ...

Input
Layer

x_1

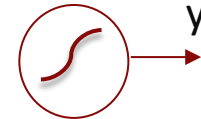


x_{29}

Hidden layer
(16 units)



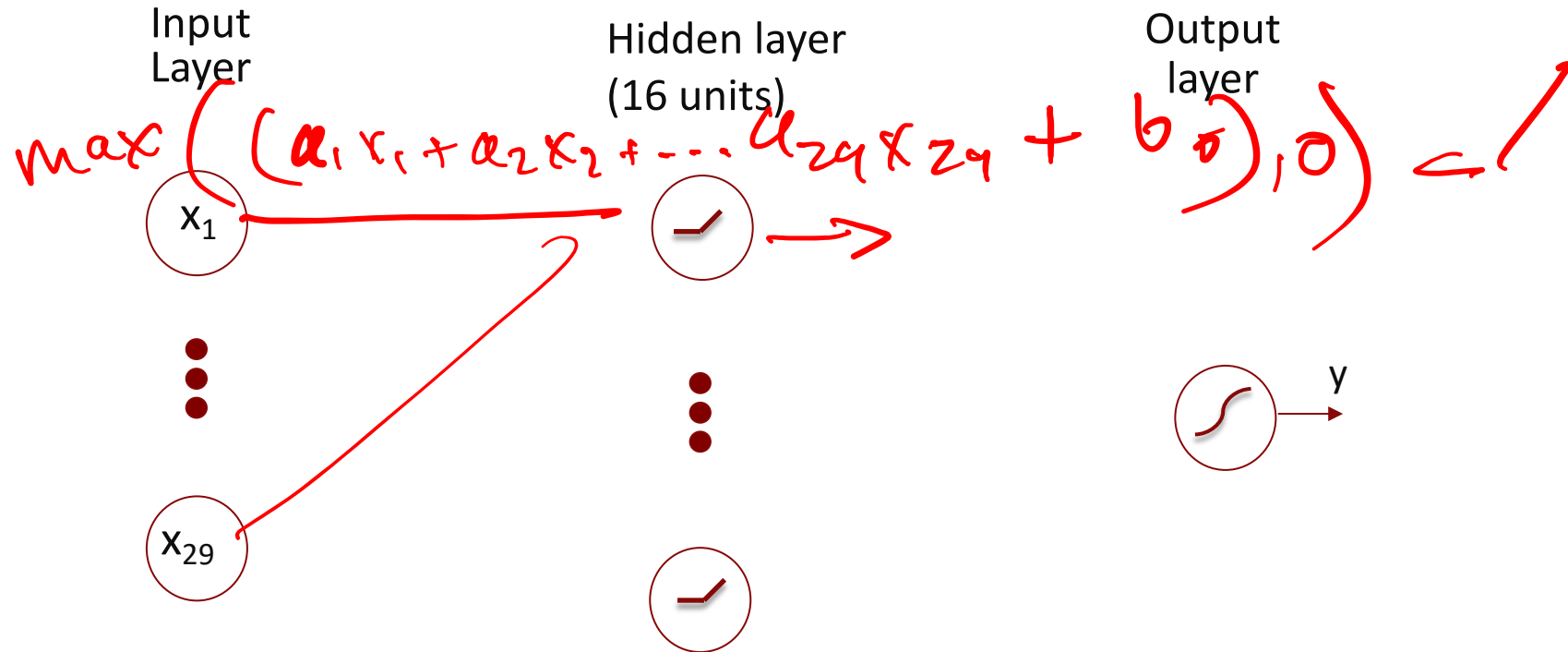
Output
layer



```
input = keras.Input(shape=29)
```

```
keras.layers.Dense(16, activation="relu")(input)
```

... and give a name to the output of this layer



```
input = keras.Input(shape=29)
```

```
h = keras.layers.Dense(16, activation="relu")(input)
```

Finally, we come to the output layer!

Input
Layer

x_1

⋮

x_{29}


Hidden layer
(16 units)



⋮



Output
layer

 y

```
input = keras.Input(shape=29)
```

```
h = keras.layers.Dense(16, activation="relu")(input)
```

```
keras.layers.Dense(1, activation="sigmoid")
```

We have just one unit in this layer ...

Input
Layer

x_1

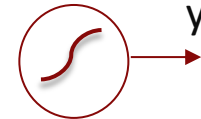


x_{29}

Hidden layer
(16 units)



Output
layer



```
input = keras.Input(shape=29)
```

```
h = keras.layers.Dense(16, activation="relu")(input)
```

```
keras.layers.Dense(1, activation="sigmoid")
```

... and indicate that we need a sigmoid activation function

Input
Layer

x_1

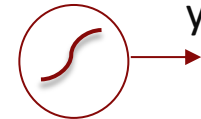


x_{29}

Hidden layer
(16 units)



Output
layer



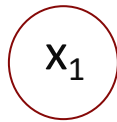
```
input = keras.Input(shape=29)
```

```
h = keras.layers.Dense(16, activation="relu")(input)
```

```
keras.layers.Dense(1, activation="sigmoid")
```

As we did before, we “feed” the output of the hidden layer to this layer ...

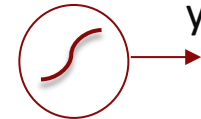
Input
Layer



Hidden layer
(16 units)



Output
layer



```
input = keras.Input(shape=29)
```

```
h = keras.layers.Dense(16, activation="relu")(input)
```

```
keras.layers.Dense(1, activation="sigmoid")(h)
```

... and give the output of this layer a name.

Input
Layer

x_1

⋮

x_{29}

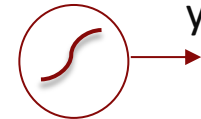
Hidden layer
(16 units)



⋮



Output
layer

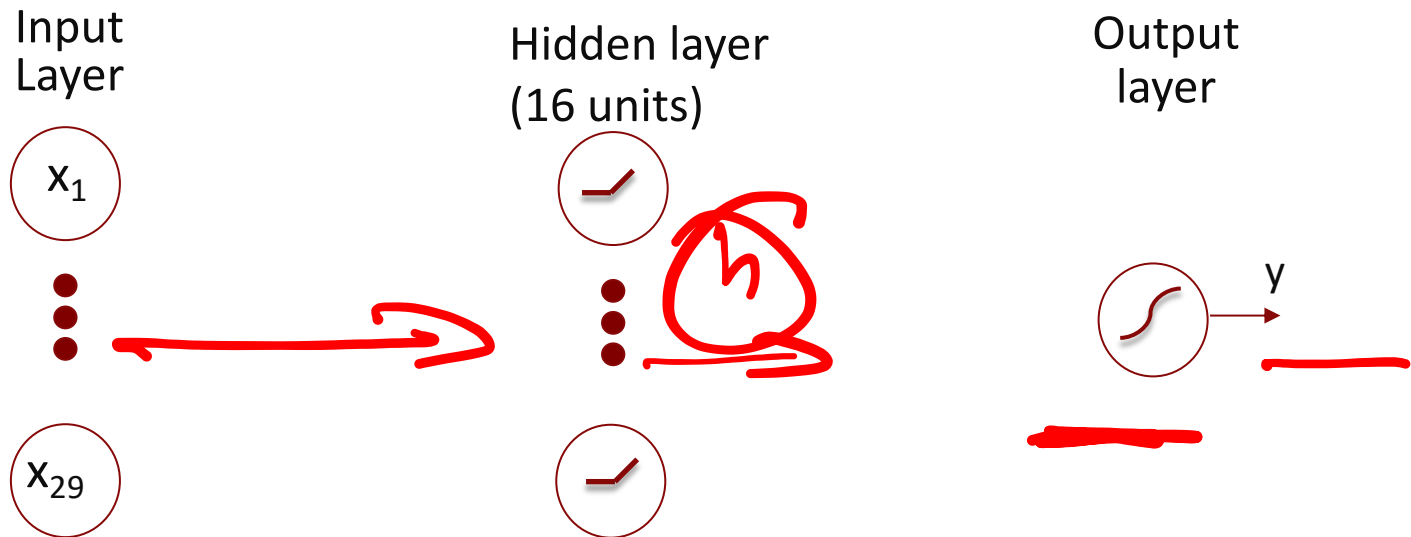


```
input = keras.Input(shape=29)
```

```
h = keras.layers.Dense(16, activation="relu")(input)
```

```
output = keras.layers.Dense(1, activation="sigmoid")(h)
```

We have defined and connected the layers

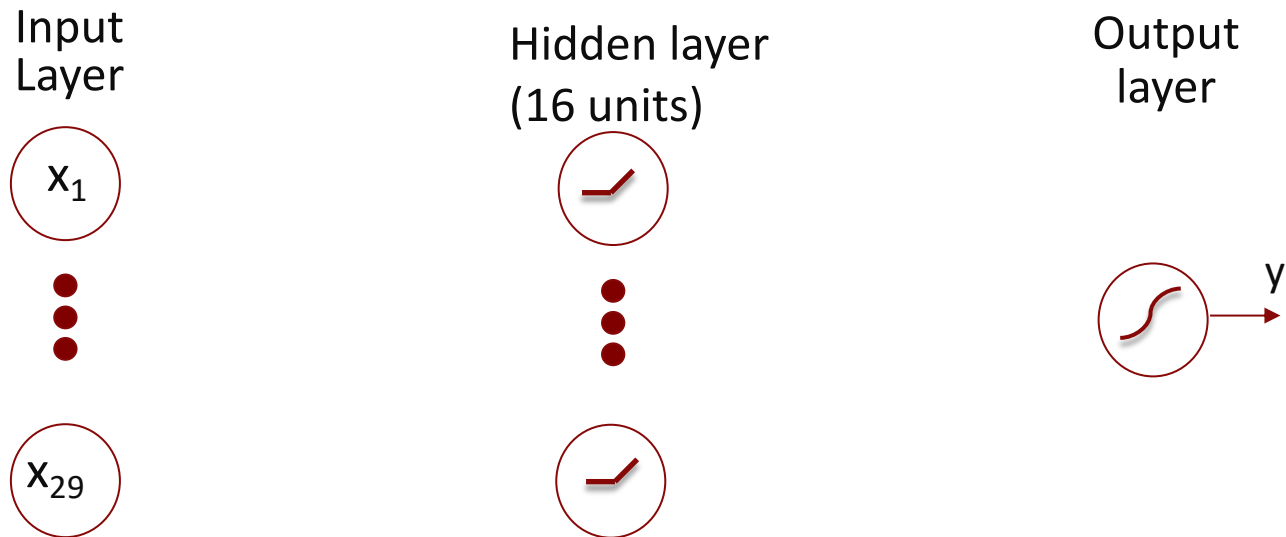


```
input = keras.Input(shape=29)
```

```
h = keras.layers.Dense(16, activation="relu")(input)
```

```
output = keras.layers.Dense(1, activation="sigmoid")(h)
```


We have defined and connected the layers.
The final step is to formally define a model.



```
input = keras.Input(shape=29)
h = keras.layers.Dense(16, activation="relu")(input)
output = keras.layers.Dense(1, activation="sigmoid")(h)
model = keras.Model(input, output)
```

That's it!

A Neural Model for Heart Disease Prediction

```
input = keras.Input(shape=29)
h = keras.layers.Dense(16, activation="relu")(input)
output = keras.layers.Dense(1, activation="sigmoid")(h)
model = keras.Model(input, output)
```

We will show how to train this model with real data and use it for prediction after we cover some *conceptual* building blocks

Training a Deep Neural Network

Recap: Training Linear and Logistic Regression Models

Linear Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

+ Data

lm

$$y = 2.8 + 0.89x_1 - 3.9x_2 + \dots + 1.06x_n$$

Recap: Training Linear and Logistic Regression Models

Linear Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

+ Data

lm

$$y = 2.8 + 0.89x_1 - 3.9x_2 + \dots + 1.06x_n$$

Logistic Regression

$$y = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$

+ Data

glm

$$y = \frac{1}{1 + e^{-(2.8 + 0.89x_1 - 3.9x_2 + \dots + 1.06x_n)}}$$

Recap: Training Linear and Logistic Regression Models

Linear Regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

+ Data

lm

$$y = 2.8 + 0.89x_1 - 3.9x_2 + \dots + 1.06x_n$$

Logistic Regression

$$y = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$

+ Data

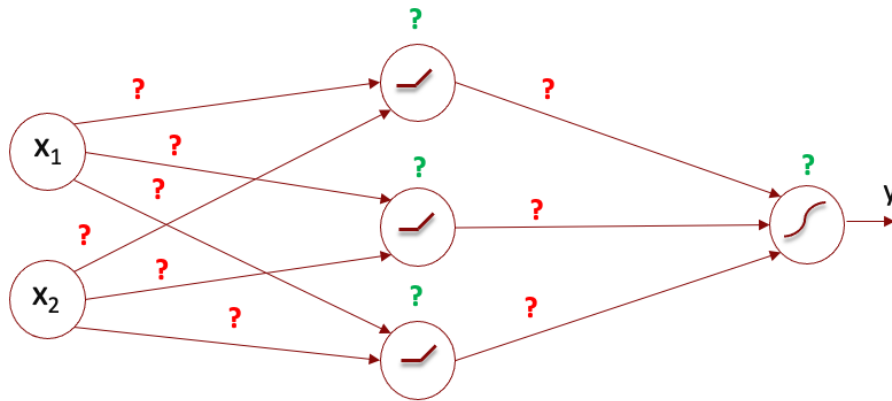
glm

$$y = \frac{1}{1 + e^{-(2.8 + 0.89x_1 - 3.9x_2 + \dots + 1.06x_n)}}$$

Recall

- Training is finding values for the weights/coefficients so that the model's predictions come as close to the actual values as possible
- 'lm' and 'glm' use optimization algorithms under the hood to find these “best” values

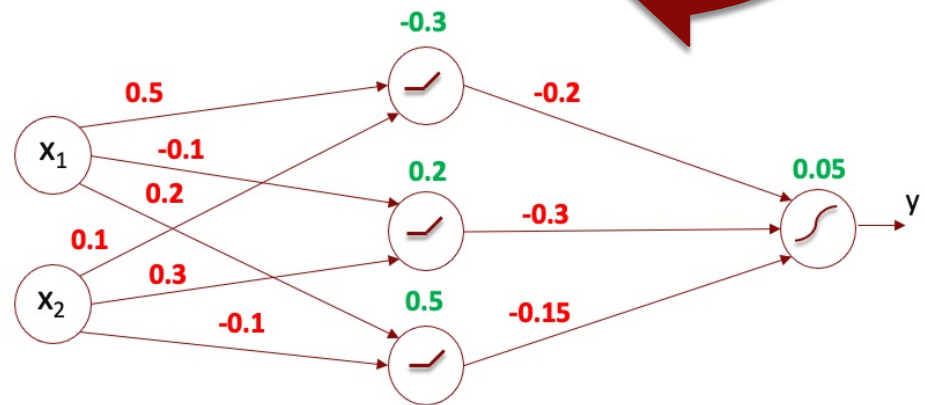
Training a DNN




+ Data

Training

Training a DNN is no different since it is just a (very complex function) with lots of parameters





The essence of training is to find the “best” values for the weights and biases i.e., those that minimize a function that measures the discrepancy between the actual and predicted values

These functions are called loss functions in the DL world

Loss Functions

Loss functions

$$(a_1x_1^1 + a_2x_2^1 + b - y_1)^2 + (a_1x_1^2 + a_2x_2^2 + b - y_2)^2 + \dots$$

- A “loss function” is a function that quantifies the error in a model’s prediction.
 - If the predictions are close to the actual values, the “loss” would be small.
 - A perfect model would have a loss of zero.

x^1, y^1

x^2, y^2

\vdots

x^{10}, y^{10}

model(x^1) = $a_1x_1^1 + a_2x_2^1 + b$

$(\text{model}(x^1) - y^1)^2 +$

$(\text{model}(x^2) - y^2)^2 +$

\vdots

$(\text{model}(x^{10}) - y^{10})^2 +$

Loss functions

- A “loss function” is a function that quantifies the error in a model’s prediction.
 - If the predictions are close to the actual values, the “loss” would be small.
 - A perfect model would have a loss of zero.
- In linear regression, you will recall that we quantify this error using “sum of squared errors”. So, “sum of squared errors” is the loss function used in linear regression

Loss functions


- A “loss function” is a function that quantifies the error in a model’s prediction.
 - If the predictions are close to the actual values, the “loss” would be small.
 - A perfect model would have a loss of zero.
- In linear regression, you will recall that we quantify this error using “sum of squared errors”. So, “sum of squared errors” is the loss function used in linear regression
- The loss function we chose must be matched well with the kind of output that comes out of the model.

Mean Squared Error (MSE) Loss

$$\frac{1}{n} \sum_{i=1}^n \left(y^i - \text{model}(x^i) \right)^2$$

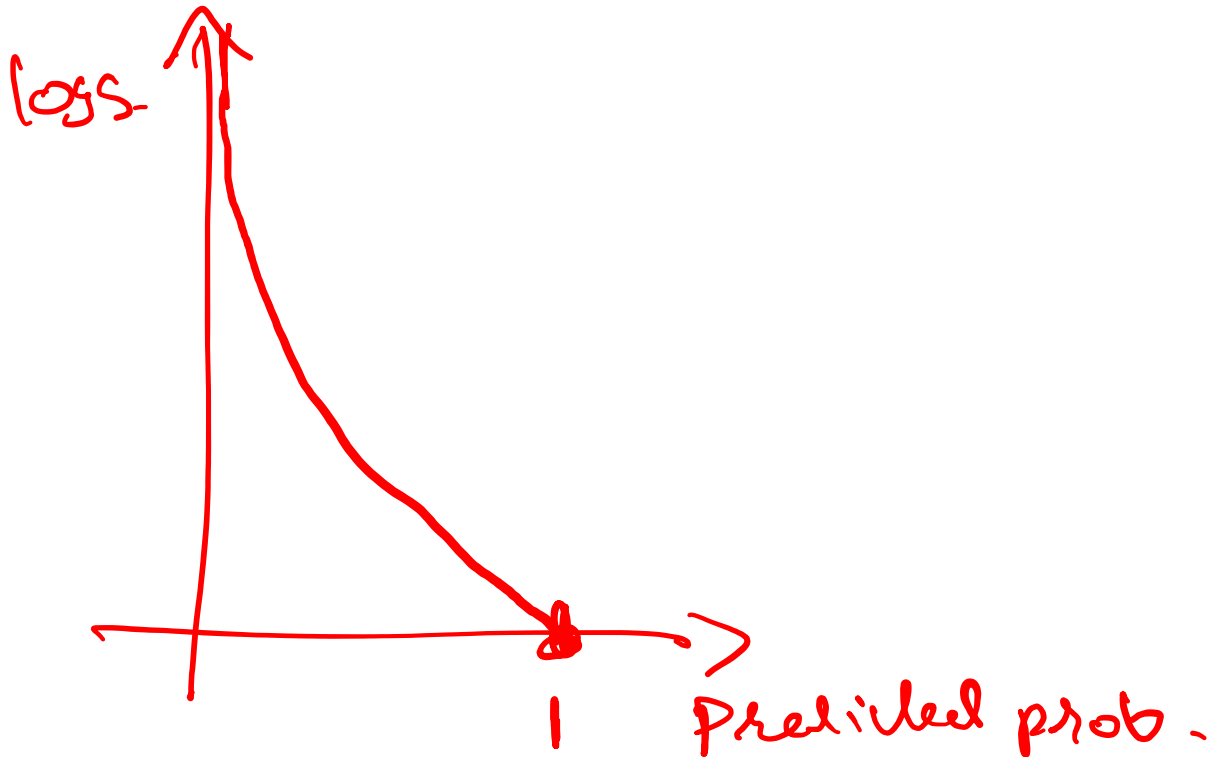
Actual value of
ith data point

Predicted value of
ith data point



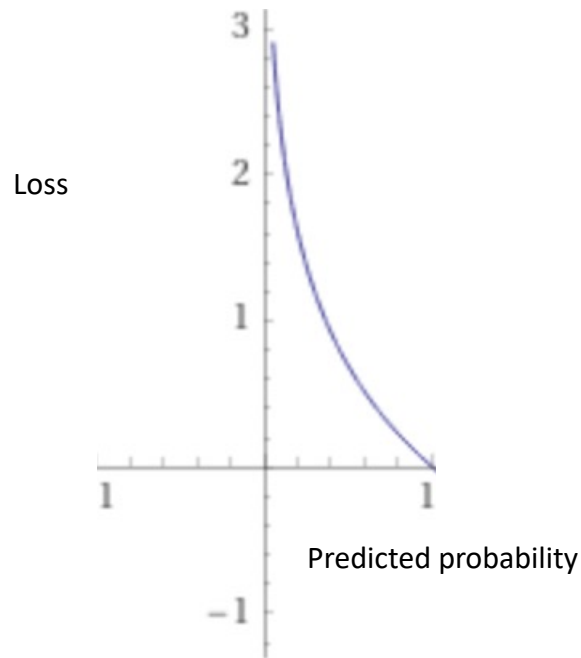
What is a good loss function for the Heart Disease Prediction Model where the prediction is a probability number and the actual output is 0-1?

For data points with $y = 1$ (i.e., patients with heart disease)..



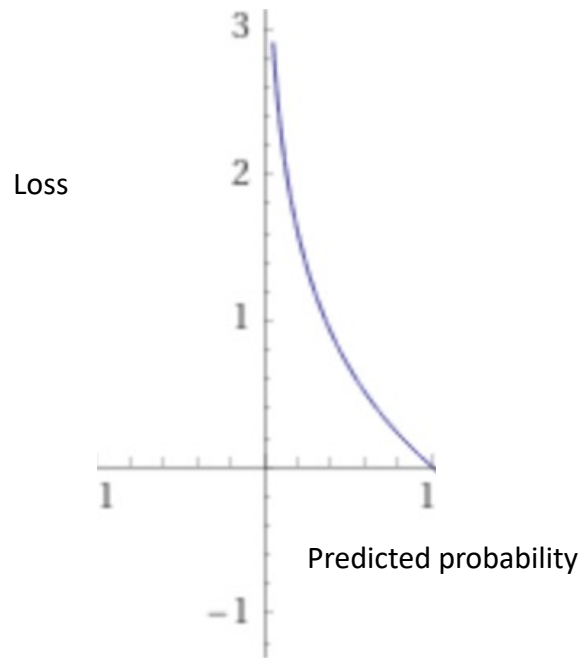
For data points with $y = 1$

For data points with $y = 1$ (i.e., patients with heart disease),
lower predicted probabilities should have **higher** loss



For data points with $y = 1$

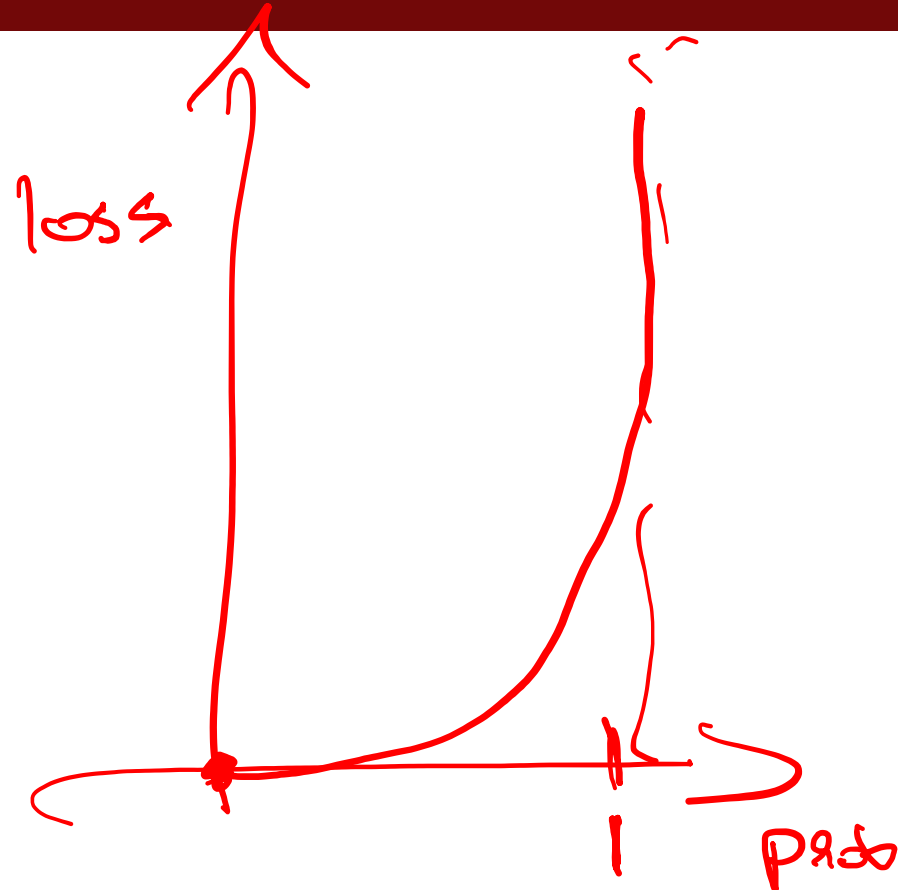
We can capture this requirement using the log function



Predicted probability	<i>$-\log(\text{predicted probability})$</i>
1/1000	9.97
1/10	3.32
1/2	1.0
1	0.0

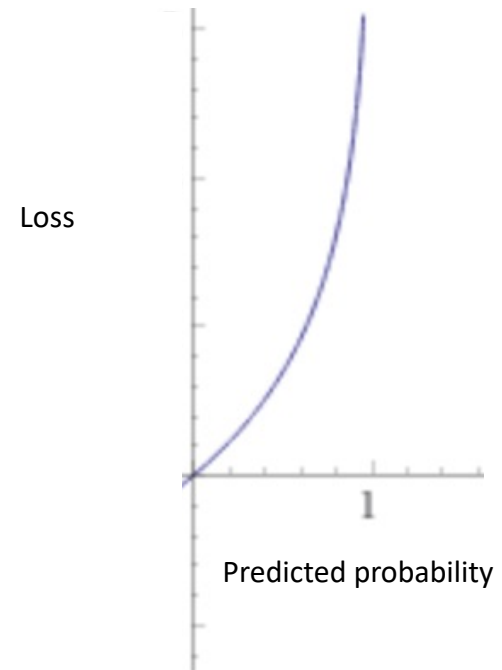
For data points with $y = 1$,
 $\text{loss} = -\log(\text{predicted probability})$

For data points with $y = 0$ (i.e., patients without heart disease) ...



For data points with $y = 0$

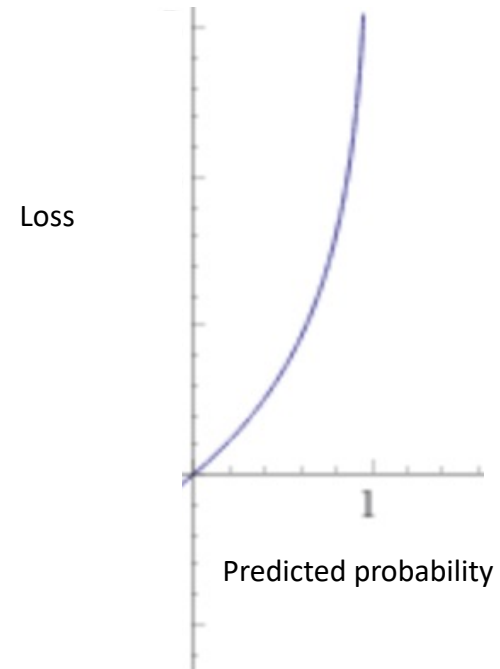
For data points with $y = 0$ (i.e., patients without heart disease), **higher** predicted probabilities should have **higher** loss



For data points with $y = 0$

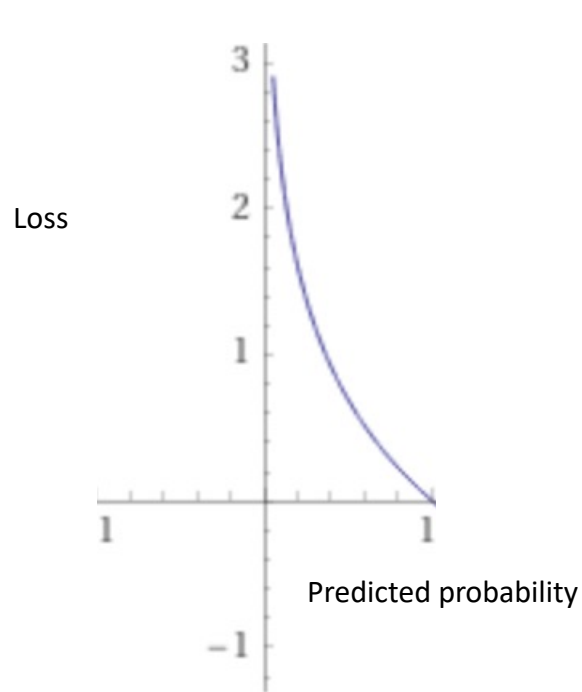
We can capture this requirement as well using the log function

Predicted probability	$-\log(1 - \text{predicted probability})$
1/1000	0.001
1/10	0.15
1/2	1.0
0.999999	19.93

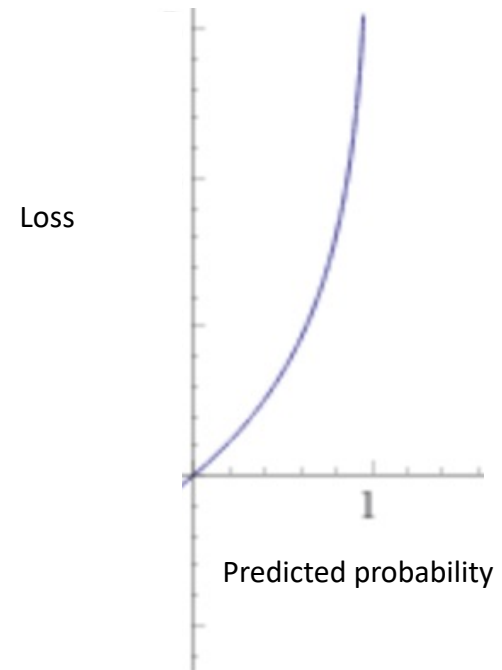


For data points with $y = 0$,
 $\text{loss} = -\log(1 - \text{predicted probability})$

Summary



For data points with $y = 1$,
loss = $-\log(\text{predicted probability})$



For data points with $y = 0$,
loss = $-\log(1 - \text{predicted probability})$

This can be compactly written as a single function

For data points with $y = 1$,
loss = -log(predicted probability)


For data points with $y = 0$,
loss = -log(1 - predicted probability)

$$\frac{1}{n} \sum_{i=1}^n -y^i \log \left(\underbrace{\text{model}(x^i)}_{\substack{\text{Predicted} \\ \text{probability for the} \\ \text{i}^{\text{th}} \text{ data point}}} \right) - (1 - y^i) \log \left(1 - \underbrace{\text{model}(x^i)}_{\substack{\text{Predicted} \\ \text{probability for the} \\ \text{i}^{\text{th}} \text{ data point}}} \right)$$

This is the Binary Cross-Entropy Loss function!

For data points with $y = 1$,
loss = -log(predicted probability)

For data points with $y = 0$,
loss = -log(1 - predicted probability)


$$\frac{1}{n} \sum_{i=1}^n -y^i \log(\text{model}(x^i)) - (1 - y^i) \log(1 - \text{model}(x^i))$$

Minimizing loss functions

Minimizing functions

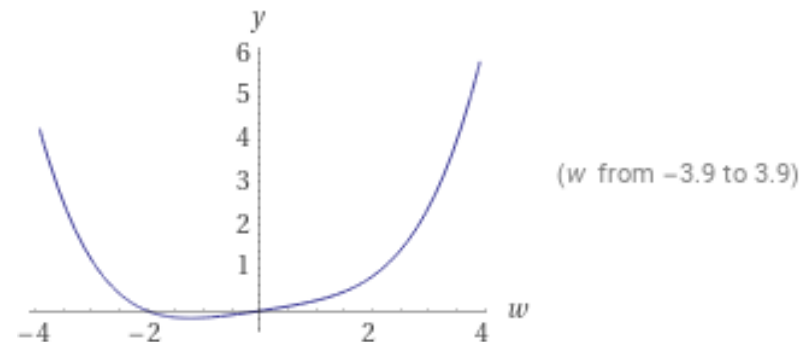


- Loss functions are just a particular kind of function so we will first consider the general problem of minimizing an arbitrary function
- After we develop some intuition about how to do this, we will return to the specific task of minimizing a loss function

Minimizing a single-variable function

Let's say we want to minimize the function:

$$g(w) = \frac{1}{50}(w^4 + w^2 + 10w)$$

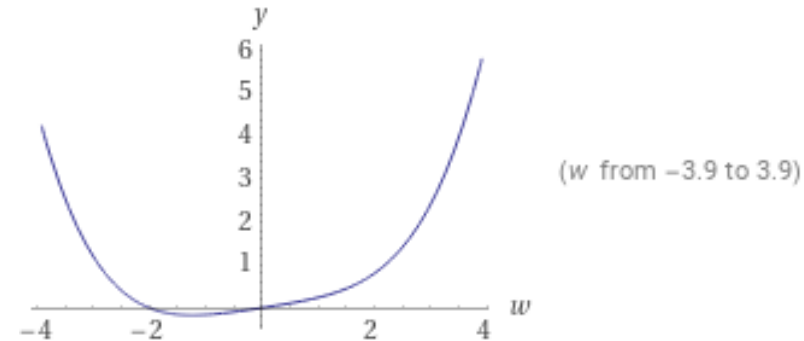


How can we go about this?

Minimizing a single-variable function

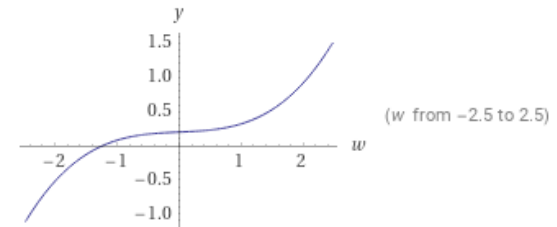
Let's say we want to minimize the function:

$$g(w) = \frac{1}{50}(w^4 + w^2 + 10w)$$



Can we use its derivative?

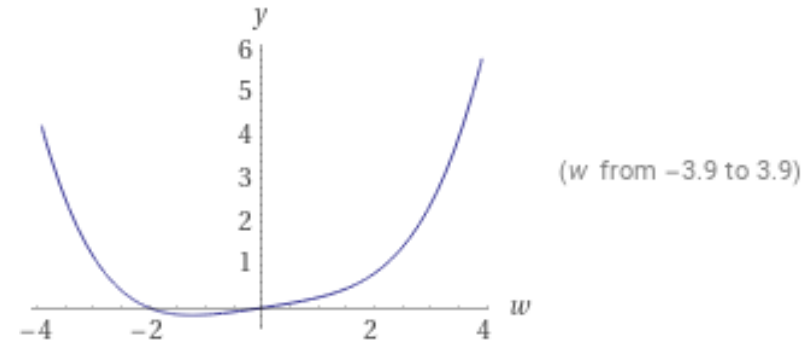
$$\frac{\partial}{\partial w}g(w) = \frac{2}{25}w^3 + \frac{1}{25}w + \frac{1}{5}$$



Minimizing a single-variable function

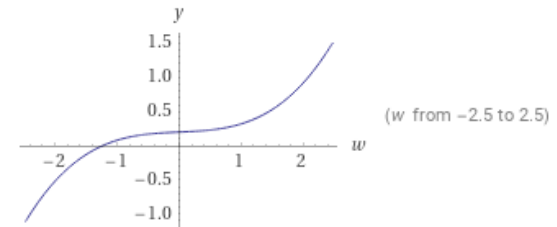
Let's say we want to minimize the function:

$$g(w) = \frac{1}{50}(w^4 + w^2 + 10w)$$



What does the derivative at a point tell us?

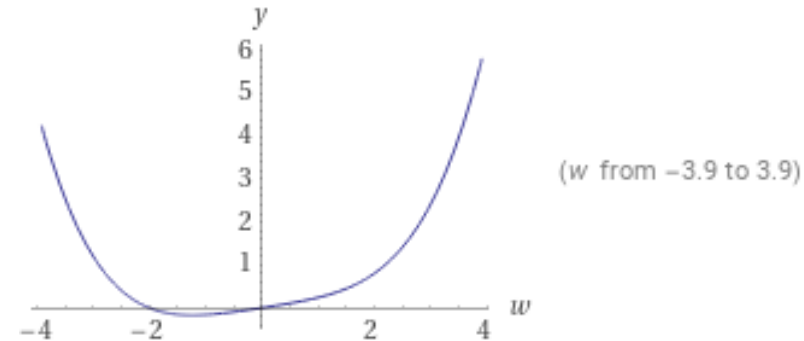
$$\frac{\partial}{\partial w} g(w) = \frac{2}{25}w^3 + \frac{1}{25}w + \frac{1}{5}$$



Minimizing a single-variable function

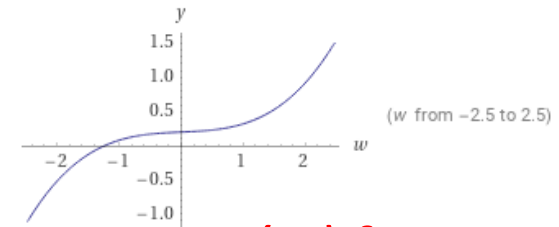
Let's say we want to minimize the function:

$$g(w) = \frac{1}{50}(w^4 + w^2 + 10w)$$



What does the derivative at a point tell us?

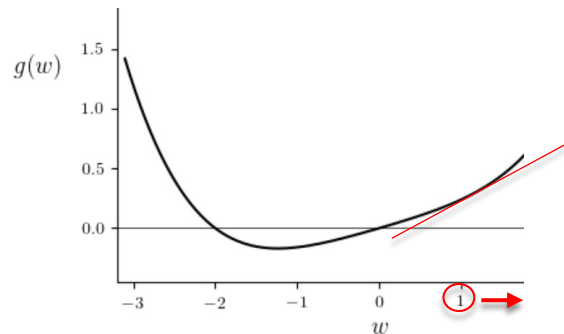
$$\frac{\partial}{\partial w}g(w) = \frac{2}{25}w^3 + \frac{1}{25}w + \frac{1}{5}$$



The derivative (or slope) tells us the change in $g(w)$ for a small increase in w

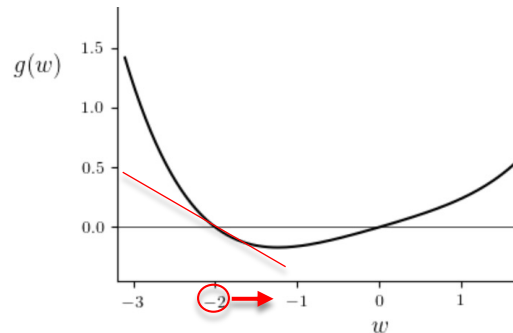
The value of knowing the derivative

If the derivative at a point w is ...	What it means
Positive	Increasing w slightly will increase $g(w)$



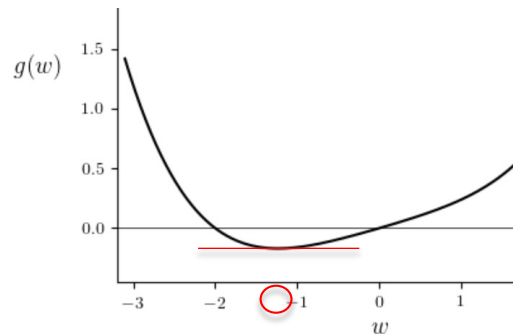
The value of knowing the derivative

If the derivative at a point w is ...	What it means
Positive	Increasing w slightly will increase $g(w)$
Negative	Increasing w slightly will decrease $g(w)$



The value of knowing the derivative

If the derivative at a point w is ...	What it means
Positive	Increasing w slightly will increase $g(w)$
Negative	Increasing w slightly will decrease $g(w)$
~ 0	Changing w slightly won't change $g(w)$



This suggests an algorithm for minimizing $g(w)$

1. Start with some point w
2. Calculate the derivative (i.e., slope) of $g(w)$ at w

If the derivative is ...	What it means	Since we want to minimize loss, do this ...
Positive	Increasing w will increase the loss function	Reduce w slightly
Negative	Increasing w will decrease the loss function	Increase w slightly
~ 0	Changing w won't change the loss function	Stop



3. Go to step 2

This is Gradient Descent!

1. Start with some point w
2. Calculate the derivative (i.e., slope) of $g(w)$ at w

This can
be written
compactly
as

If the derivative is ...	What it means	Since we want to minimize loss, do this ...
Positive	Increasing w will increase the loss function	Reduce w slightly
Negative	Increasing w will decrease the loss function	Increase w slightly
~ 0	Changing w won't change the loss function	Stop


$$w \leftarrow w - \alpha \frac{dg(w)}{dw}$$


3. Go to step 2

Gradient Descent

$$w \leftarrow w - \alpha \frac{dg(w)}{dw}$$

α is called the “**learning rate**” and is our way of ensuring that we increase or decrease w **slightly**

Typically set to small values (e.g., 0.1, 0.001, 0.0001) and determined by trial and error

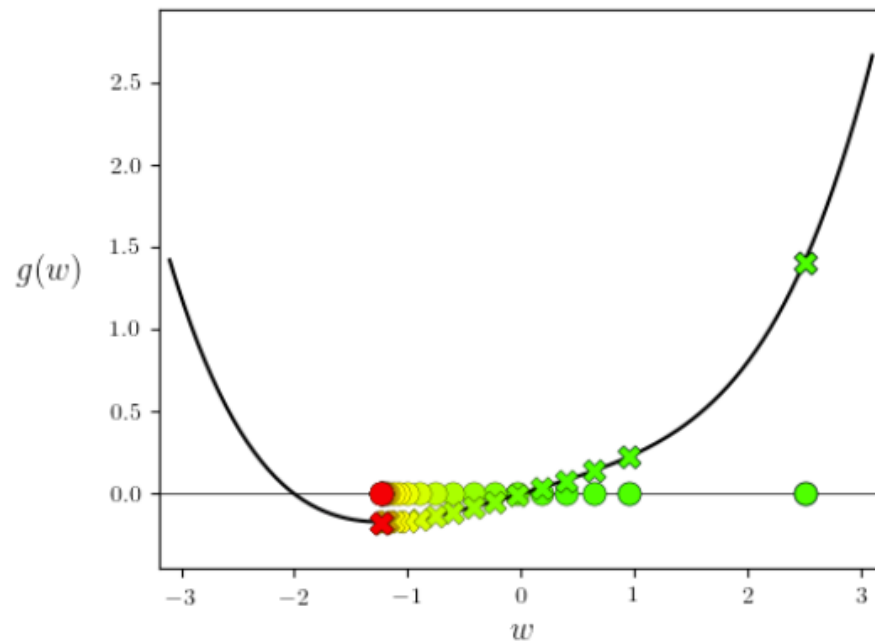
Let's apply this algorithm to $g(w)$

$$\frac{\partial}{\partial w} g(w) = \frac{2}{25}w^3 + \frac{1}{25}w + \frac{1}{5}$$

$$w \leftarrow w - \alpha \frac{dg(w)}{dw}$$

We will start at $w = 2.5$, set $\alpha = 1$ and run the algorithm for a few iterations (switch to animation)

Gradient Descent in action



Minimizing a multi-variable function

$$g(w_1, w_2) = w_1^2 + w_2^2 + 2$$

We can calculate the partial derivative of $g(w_1, w_2)$

$$\left[\frac{\partial g}{\partial w_1}, \frac{\partial g}{\partial w_2} \right] = [2w_1, 2w_2]$$

How should we interpret this?

Minimizing a multi-variable function

$$g(w_1, w_2) = w_1^2 + w_2^2 + 2$$

$$\nabla g = \left[\frac{\partial g}{\partial w_1}, \frac{\partial g}{\partial w_2} \right] = [2w_1, 2w_2]$$

The first number is the change in $g(w)$ for a small increase in w_1 , *with w_2 kept unchanged*. The second number is the change in $g(w)$ for a small increase in w_2 , *with w_1 kept unchanged*

Minimizing a multi-variable function

$$g(w_1, w_2) = w_1^2 + w_2^2 + 2$$

$$\nabla g = \left[\frac{\partial g}{\partial w_1}, \frac{\partial g}{\partial w_2} \right] = [2w_1, 2w_2]$$

The first number is the change in $g(w)$ for a small increase in w_1 , *with w_2 kept unchanged*. The second number is the change in $g(w)$ for a small increase in w_2 , *with w_1 kept unchanged*

This is called the “gradient” of $g(w_1, w_2)$ and written as ∇g

Minimizing a multi-variable function

$$\nabla g = \left[\frac{\partial g}{\partial w_1}, \frac{\partial g}{\partial w_2} \right] = [2w_1, 2w_2]$$

We can do gradient descent on each coordinate by using the corresponding partial derivative.

$$w_1 \leftarrow w_1 - \alpha \left(\frac{\partial g}{\partial w_1} \right)$$

$$w_2 \leftarrow w_2 - \alpha \left(\frac{\partial g}{\partial w_2} \right)$$

Minimizing a multi-variable function

$$\nabla g = [2w_1, 2w_2]$$

$$w_1 \leftarrow w_1 - \alpha \left(\frac{\partial g}{\partial w_1} \right)$$

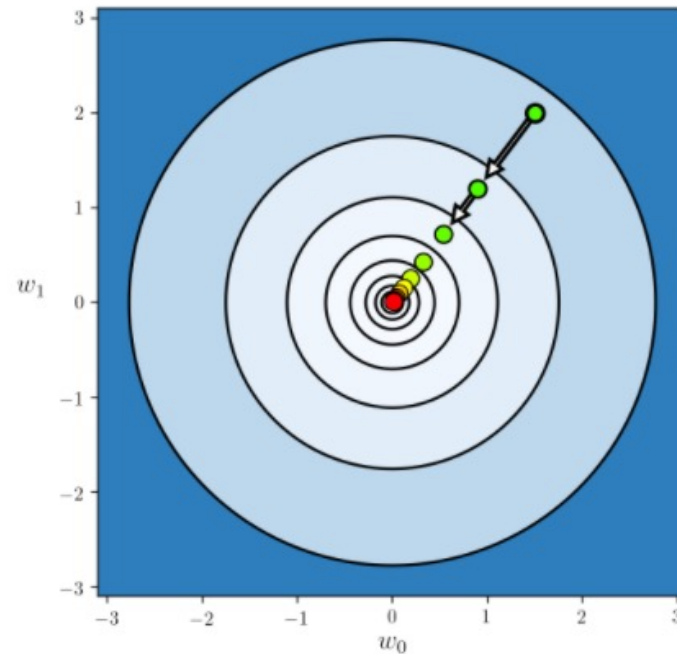
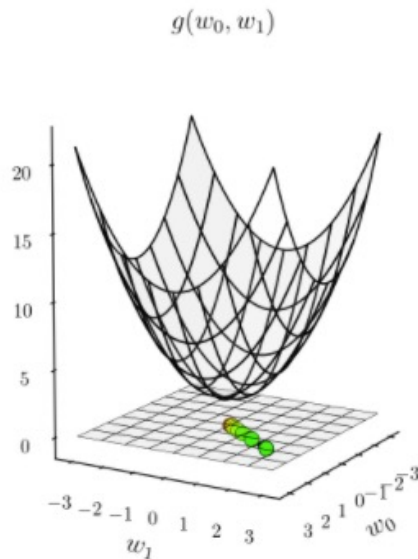
$$w_2 \leftarrow w_2 - \alpha \left(\frac{\partial g}{\partial w_2} \right)$$

As before, this whole thing can be summarized compactly as:


$$w \leftarrow w - \alpha \nabla g(w)$$

Gradient Descent in two dimensions

$$g(w_0, w_1) = w_0^2 + w_1^2 + 2$$



https://kenndanielso.github.io/mlrefined/blog_posts/6_First_order_methods/6_4_Gradient_descent.html



GD may stop near a local minimum (not necessarily a global minimum) or a saddle point but we don't worry about this in practice

Minimizing a **loss function** with gradient descent

$$\textit{Minimize} \quad \frac{1}{n} \sum_{i=1}^n -y^i \log(\textit{model}(x^i)) - (1 - y^i) \log(1 - \textit{model}(x^i))$$

What are the variables we need to change to minimize this function?

Minimizing a **loss function** with gradient descent

$$\textit{Minimize} \quad \frac{1}{n} \sum_{i=1}^n -y^i \log(\textit{model}(x^i)) - (1 - y^i) \log(1 - \textit{model}(x^i))$$

What are the variables we need to change to minimize this function?

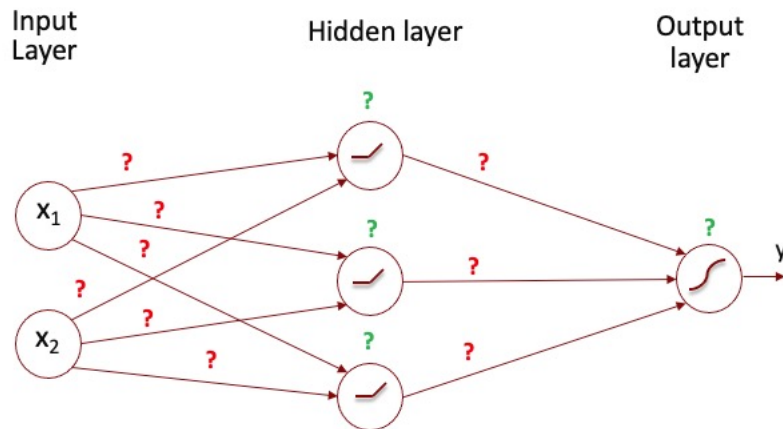
They are the parameters “hiding” inside $\textit{model}(x_i)$

Minimizing a **loss function** with gradient descent

Minimize $\frac{1}{n} \sum_{i=1}^n -y^i \log(\text{model}(x^i)) - (1 - y^i) \log(1 - \text{model}(x^i))$


$$\text{model}(x^i) = \frac{1}{1 + e^{-(w_1 + w_2 \max(0, w_3 + w_4 x_1^i + w_5 x_2^i) + w_6 \max(0, w_7 + w_8 x_1^i + w_9 x_2^i) + w_{10} \max(0, w_{11} + w_{12} x_1^i + w_{13} x_2^i))}}$$

Recall this model
and the NN it
represents



Minimizing a **loss function** with gradient descent

Minimize

$$\frac{1}{n} \sum_{i=1}^n -y^i \log(\text{model}(x^i)) - (1 - y^i) \log(1 - \text{model}(x^i))$$


$$\text{model}(x^i) = \frac{1}{1 + e^{-(w_1 + w_2 \max(0, w_3 + w_4 x_1^i + w_5 x_2^i) + w_6 \max(0, w_7 + w_8 x_1^i + w_9 x_2^i) + w_{10} \max(0, w_{11} + w_{12} x_1^i + w_{13} x_2^i))}}$$

w_1, w_2, \dots, w_{13} are the variables we can change to minimize the loss function

The values of x_1, x_2 and y , on the other hand, are just data

Minimizing a **loss function** with gradient descent

Minimize

$$\frac{1}{n} \sum_{i=1}^n -y^i \log(\text{model}(x^i)) - (1 - y^i) \log(1 - \text{model}(x^i))$$

$$\text{model}(x^i) = \frac{1}{1 + e^{-(w_1 + w_2 \max(0, w_3 + w_4 x_1^i + w_5 x_2^i) + w_6 \max(0, w_7 + w_8 x_1^i + w_9 x_2^i) + w_{10} \max(0, w_{11} + w_{12} x_1^i + w_{13} x_2^i))}}$$

Imagine replacing $\text{model}(x^i)$ with the mathematical expression above wherever $\text{model}(x^i)$ appears in the loss function

Now, your loss function is just a "good old" function of w_1, w_2, \dots, w_{13} and you can apply gradient descent to it as we normally would.

Backpropagation (aka ‘backprop’)

- Backpropagation is an efficient way to compute the gradient of the loss function
- The efficiency stems from exploiting the layer-by-layer architecture of NNs
- By organizing the computation in the form of a “computational graph”, we can incrementally calculate the gradient one layer at a time using **matrix multiplications** (and other simple operations). This approach also eliminates redundant calculations
- It turns out that Graphic Processing Units (GPUs), originally invented to speed up video games, are **perfectly suited for matrix multiplications**!
- Backprop + GPUs → Fast calculation of loss function gradients!

Please see HODL-SP24-Lec-2-Backprop_Example.pdf for a worked-out example

Gradient Descent → Stochastic Gradient Descent

Making Gradient Descent work with large datasets

- Problem: For large datasets (e.g., n in the millions), computing the gradient of the loss function can be very expensive

Making Gradient Descent work with large datasets

- Problem: For large datasets (e.g., n in the millions), computing the gradient of the loss function can be very expensive
- The Solution:
 - At each iteration, instead of using all the n data points in the calculation of the gradient of the loss function, randomly choose just a few of the n observations (called a *minibatch*) and use only these observations to compute the partial derivatives.

Making Gradient Descent work with large datasets

- Problem: For large datasets (e.g., n in the millions), computing the gradient of the loss function can be very expensive
- The Solution:
 - At each iteration, instead of using all the n data points in the calculation of the gradient of the loss function, randomly choose just a few of the n observations (called a *minibatch*) and use only these observations to compute the partial derivatives.
 - This is called Stochastic Gradient Descent (SGD)*

* Strictly speaking, SGD chooses just one observation. What we are describing here is Minibatch Gradient Descent but the term SGD is widely used in the field to describe the latter so we will do the same

Making Gradient Descent work with large datasets

- Problem: For large datasets (e.g., n in the millions), computing the gradient of the loss function can be very expensive
- The Solution:
 - At each iteration, instead of using all the n data points in the calculation of the gradient of the loss function, randomly choose just a few of the n observations (called a *minibatch*) and use only these observations to compute the partial derivatives.
 - This is called Stochastic Gradient Descent (SGD)
 - Because not all n data points are used in the calculation, this only approximates the true gradient but nevertheless works well in practice. In fact, because it is only an approximation of the true gradient, it can escape local minima.
 - SGD comes in many “flavors” and we will see these flavors in the remainder of HODL

Summary of overall training flow

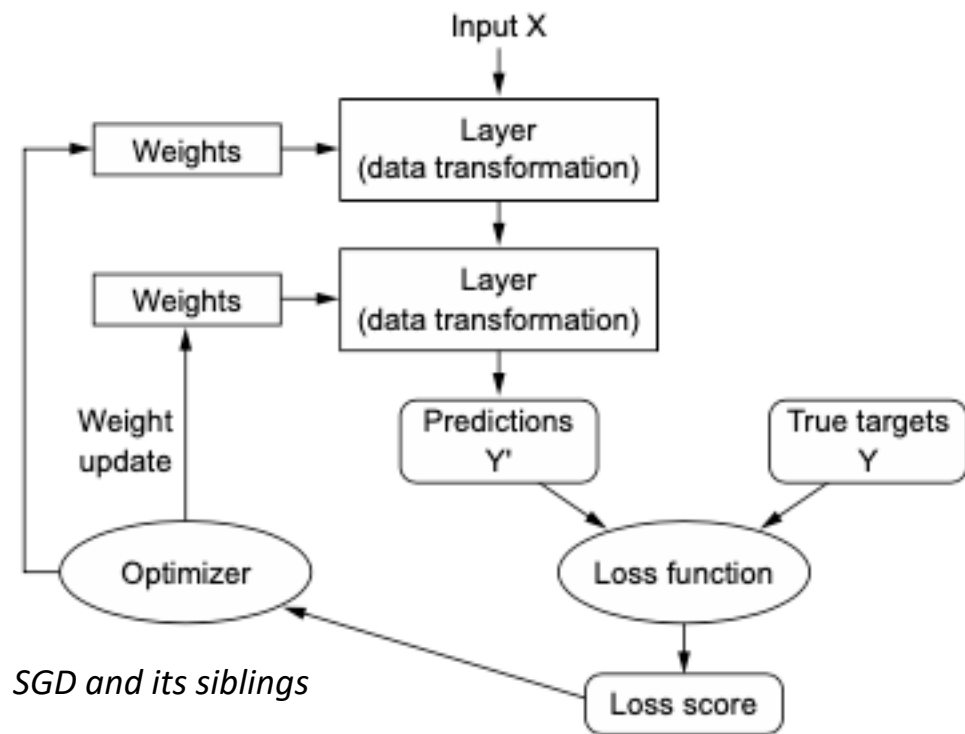


Figure 2.26 Relationship between the network, layers, loss function, and optimizer

Summary: Creating and training a DNN from scratch

- We get the data ready
- We design i.e., “lay out” the network
 - Choose the number of hidden layers and the number of ‘neurons’ in each layer
 - Pick the right output layer based on the type of the output (more on this shortly)
- We pick
 - An appropriate loss function based on the type of the output (more on this shortly)
 - An optimizer from the many SGD flavors that are available and a “good” learning rate
- We set things up in Keras/Tensorflow and start training!

Lightning Intro to Tensorflow/Keras

Tensorflow



Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions

Tensorflow



Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions
- Library of state-of-the-art optimizers

Tensorflow



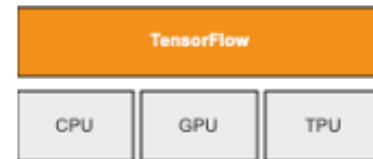
Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions
- Library of state-of-the-art optimizers
- Automatic distribution of computational load across servers

Tensorflow

Tensorflow (TF) is a library that provides

- Automatic calculation of the gradient of (complicated) loss functions
- Library of state-of-the-art optimizers
- Automatic distribution of computational load across servers
- Automatic adaptation of code to work on parallel hardware (GPUs and TPUs)



Keras “sits on top of” Tensorflow ...

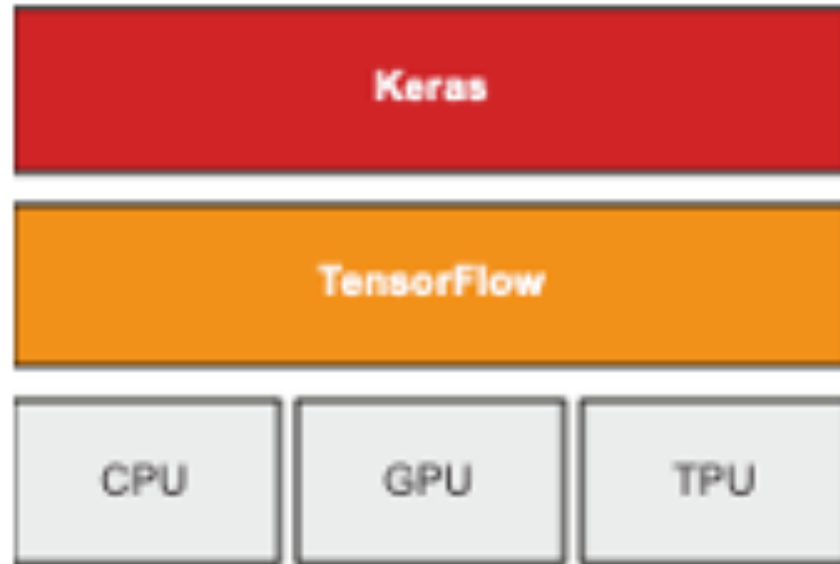



Image: Page 70 of textbook

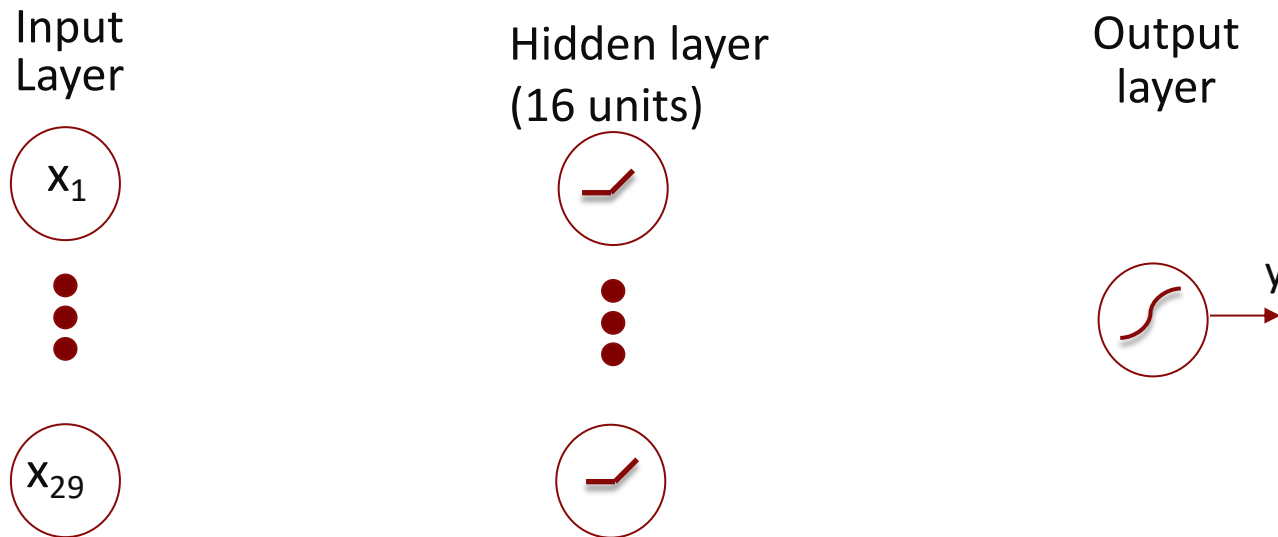
... and provides “convenience” features

- Pre-defined **layers**
- Incredibly flexible ways to specify network **architectures**
- Easy ways to **preprocess** data
- Easy ways to **train** models and **report** metrics
- **Pre-trained models** you can download and customize



Check out the wealth of introductory
and advanced material, with
accompanying colabs, at
[tensorflow.org](https://www.tensorflow.org) and keras.io

Let's revisit the Heart Disease Prediction Model we defined earlier



```
input = keras.Input(shape=29)
h = keras.layers.Dense(16, activation="relu")(input)
output = keras.layers.Dense(1, activation="sigmoid")(h)
model = keras.Model(input, output)
```

Training Checklist

- We get the data ready (will cover in the colab)
- We design i.e., “lay out” the network **1 hidden layer with 16 ReLU neurons**
 - Choose *the number of hidden layers* and *the number of ‘neurons’ in each layer*
 - Pick the *right output layer* based on the type of the output **Sigmoid**
- We pick
 - An appropriate *loss function* based on the type of the output **binary crossentropy**
 - An *optimizer from the many SGD flavors* that are available **“adam”**
- We set things up in Keras/Tensorflow and start training!

Colab

Predicting Heart Disease

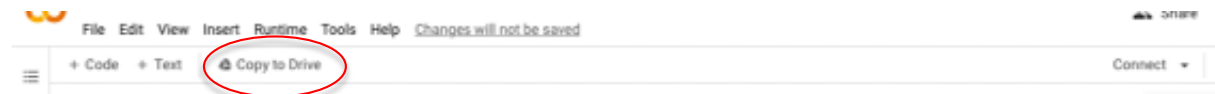
Before we start coding ...

- Don't worry if you don't understand every detail of what we will do in class.
- But go through the Colab notebooks carefully later, play around with the code and make sure you understand every line

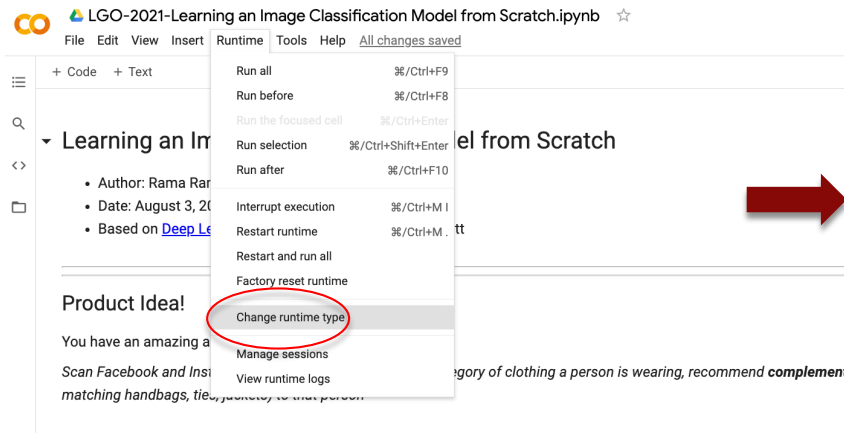
Colab Instructions

<https://colab.research.google.com/drive/1bR9Tx87L4HxB94rlp-mhy4Aj4XJtpcGS?usp=sharing>

Step 1 Make your own copy of the notebook



Step 2 Request a GPU for your notebook



Notebook settings

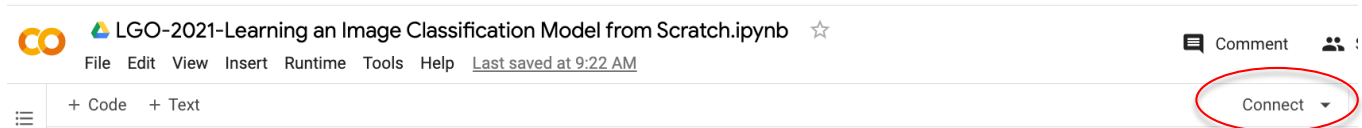
Hardware accelerator
GPU

To get the most out of Colab, avoid using a GPU unless you need one. [Learn more](#)

☐ Omit code cell output when saving this notebook

Cancel Save

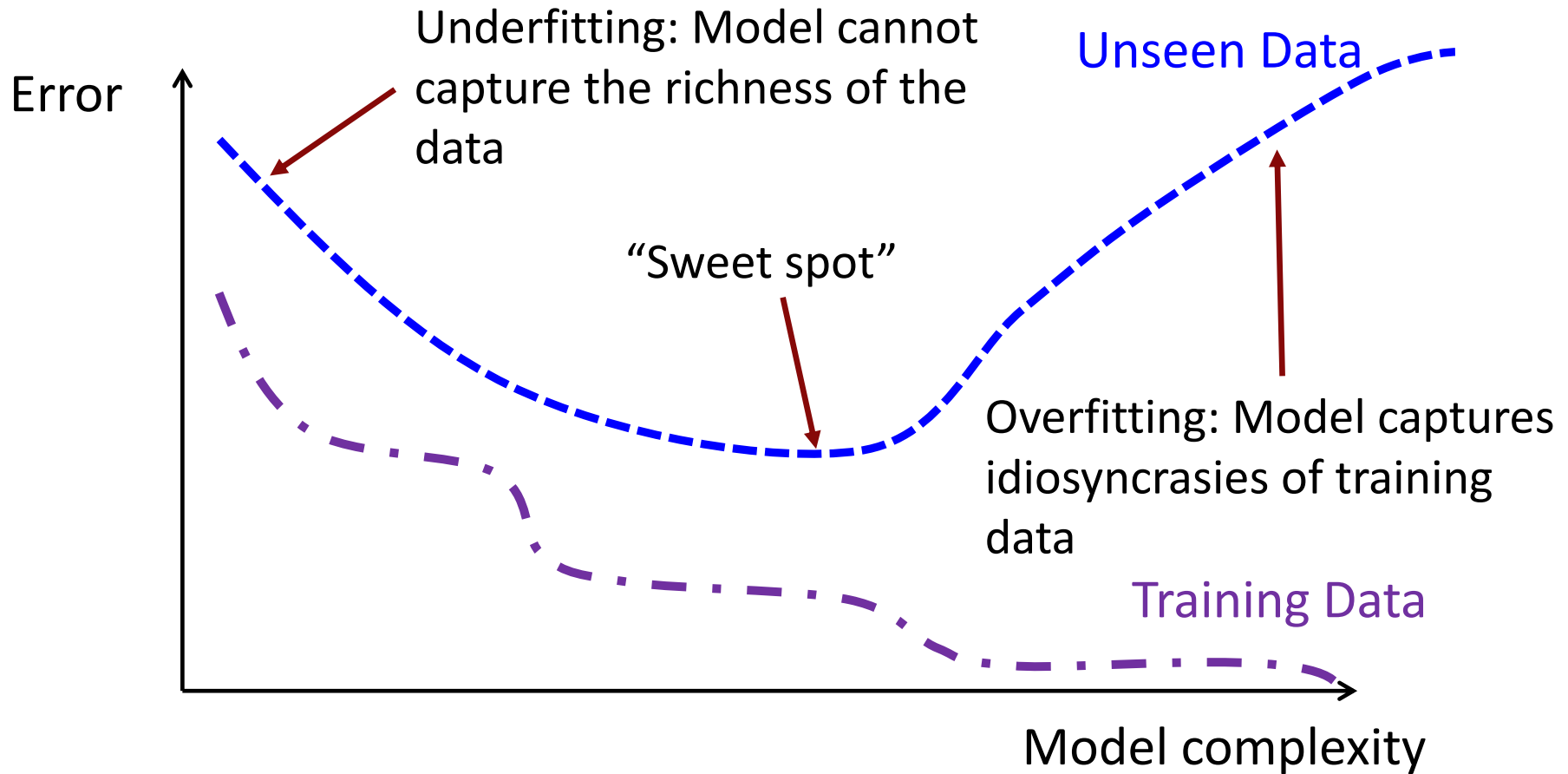
Step 3 Start your notebook



You need to do steps 1 and 2 just the first time you use a notebook. From the second time onwards, jump to Step 3.

Overfitting and Regularization

Recall Underfitting vs. Overfitting

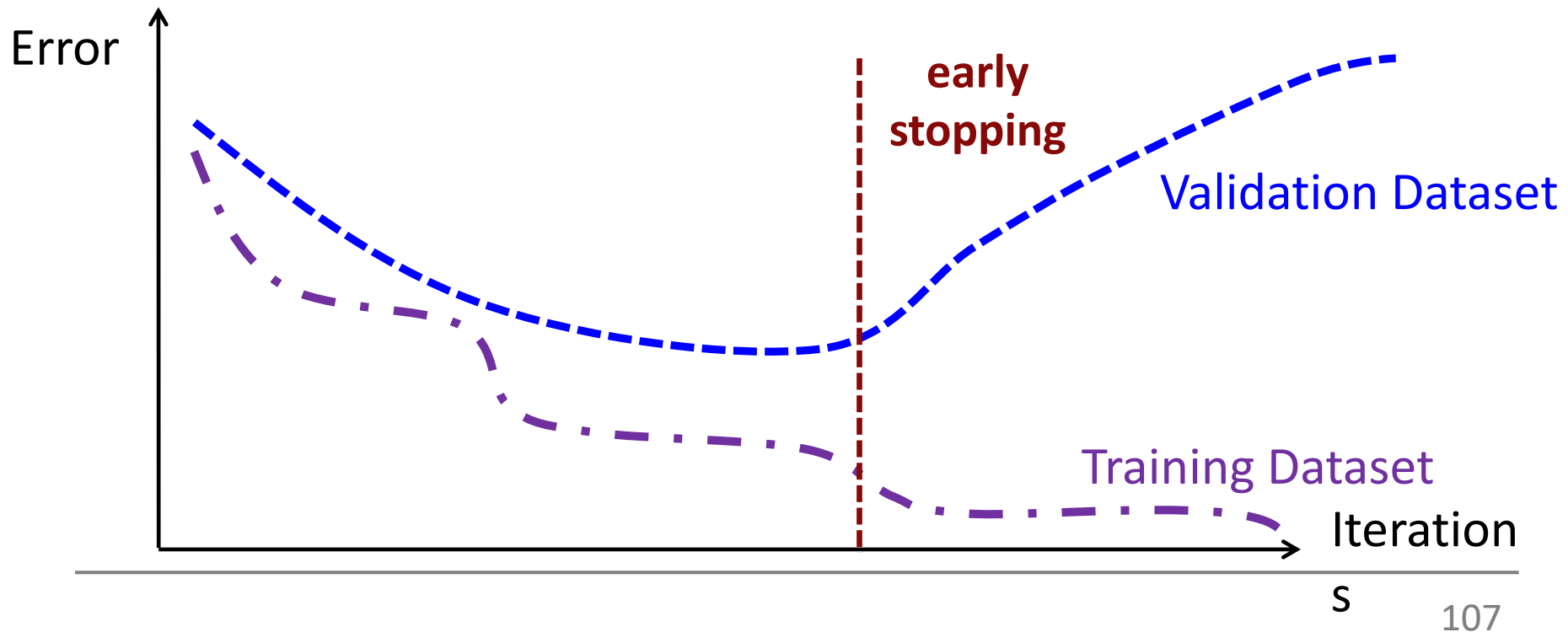


Overfitting in Neural Networks

- To learn smart representations of complex, unstructured data, the NN needs to have large “capacity” i.e., many layers and many neurons in each layer
- But this raises the likelihood of overfitting so we need to add *regularization*
- Several regularization methods have been developed to address this problem

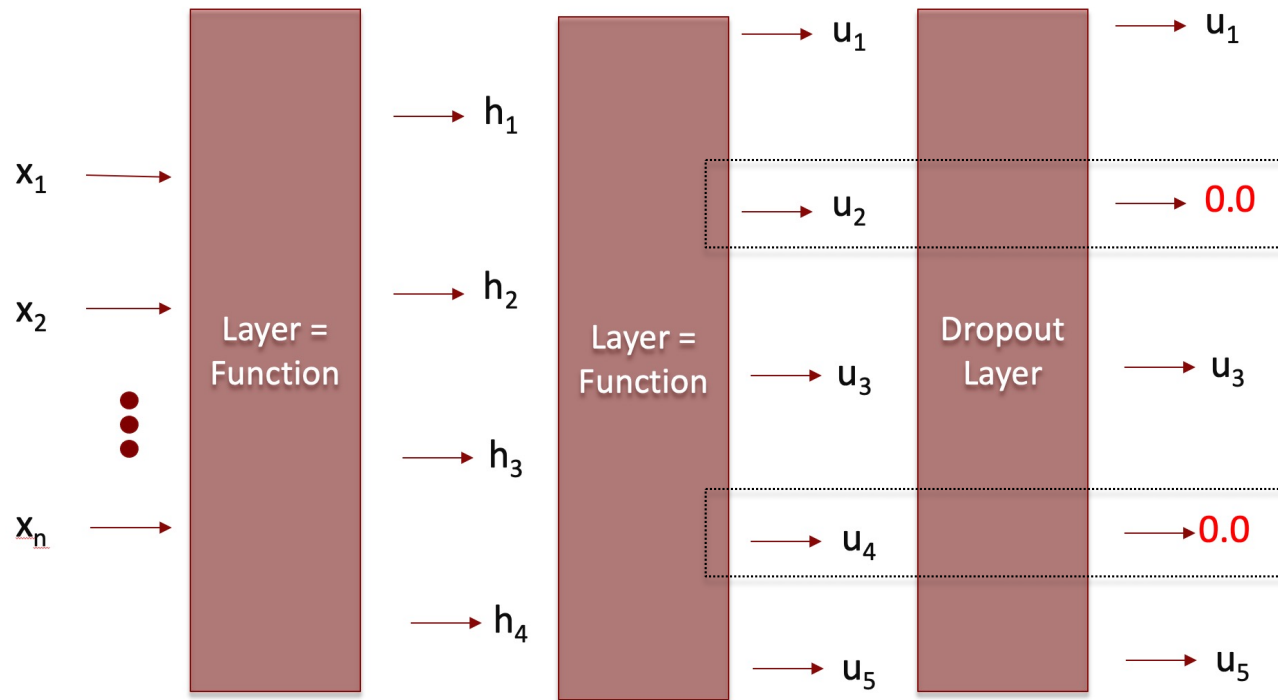
Regularization strategy: *Early Stopping*

Stop the training early before the training loss is minimized by monitoring the loss on a validation dataset.



Regularization strategy: *Dropout*

Randomly zero out the output from some of the nodes (typically 50% of the nodes) in a hidden layer (implemented as a “dropout layer” in Keras)



Training Checklist

- We get the data ready (will cover in the colab)
- We design i.e., “lay out” the network **1 hidden layer with 16 ReLU neurons**
 - Choose *the number of hidden layers* and *the number of ‘neurons’ in each layer*
 - Pick the *right output layer* based on the type of the output **Sigmoid**
- We pick
 - An appropriate *loss function* based on the type of the output **binary crossentropy**
 - An *optimizer from the many SGD flavors* that are available **“adam”**
- We decide on a *regularization strategy* **Early stopping**
- We set things up in Keras/Tensorflow and start training!

Appendix



Why Accuracy and MSE aren't good loss functions when the output is a probability

Why Accuracy and MSE aren't ideal loss functions for a probability output

First approach (accuracy):

Set threshold θ (say $\theta = 0.5$) and predict 1 if $\text{model}(x^i) \geq \theta$.

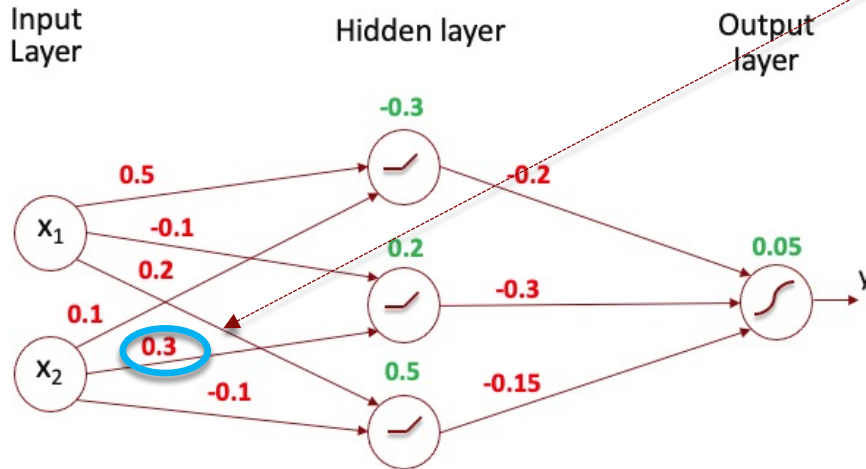
If predicted output equals the true output, the accuracy loss is 0.0

If they differ, the accuracy loss is 1.0

Problem: The gradient of the accuracy loss will be zero most of the time. As a result, the weights won't change in the gradient descent step (i.e., *learning won't happen*)

How the gradient of “Accuracy loss” can be zero

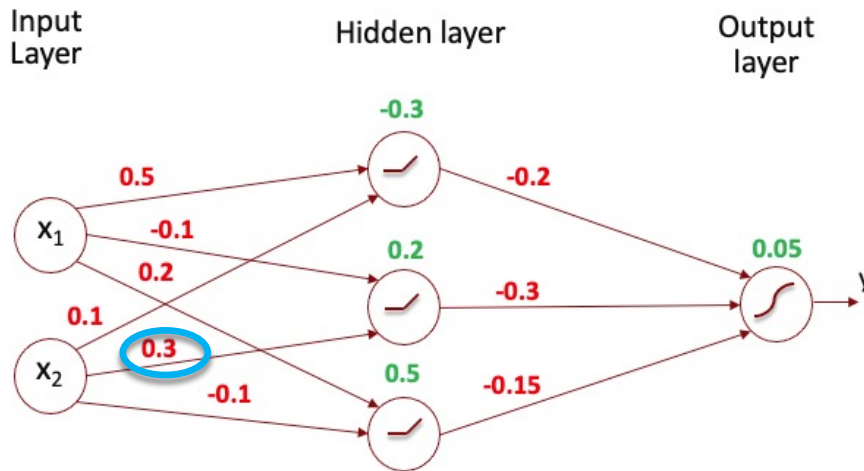
Imagine increasing this weight from 0.3 to 0.301, keeping the other weights unchanged.



How the gradient of “Accuracy loss” can be zero

Imagine increasing this weight from 0.3 to 0.301, keeping the other weights unchanged.

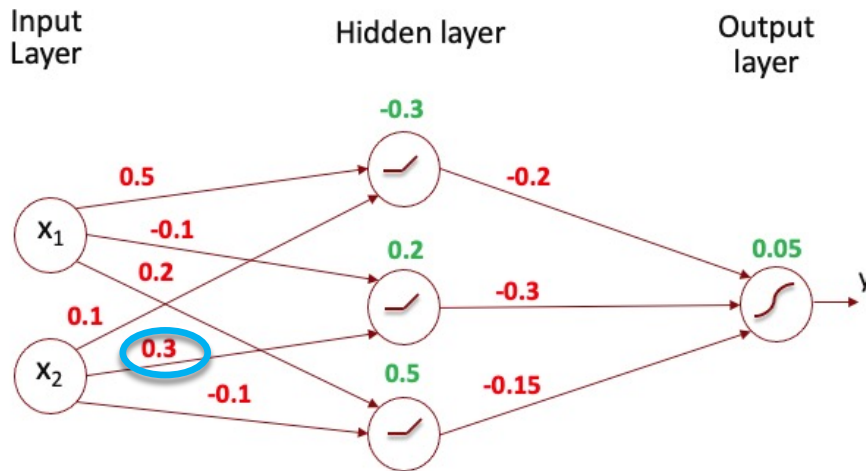
For every data point x , the predicted value $\text{model}(x)$ will also change but very slightly. **But this change will change the classification of x only if the predicted value went from below the threshold of (say) 0.5 to above 0.5 (or vice-versa).**



How the gradient of “Accuracy loss” can be zero

Imagine increasing this weight from 0.3 to 0.301, keeping the other weights unchanged.

For every data point x , the predicted value $\text{model}(x)$ will also change but very slightly. But this change will change the classification of x only if the predicted value went from below the threshold of (say) 0.5 to above 0.5 (or vice-versa).



This “0.5 crossing” will happen only if the predicted value for an x is very close to 0.5 before the change. **This is very unlikely to be the case for most/all the points.**

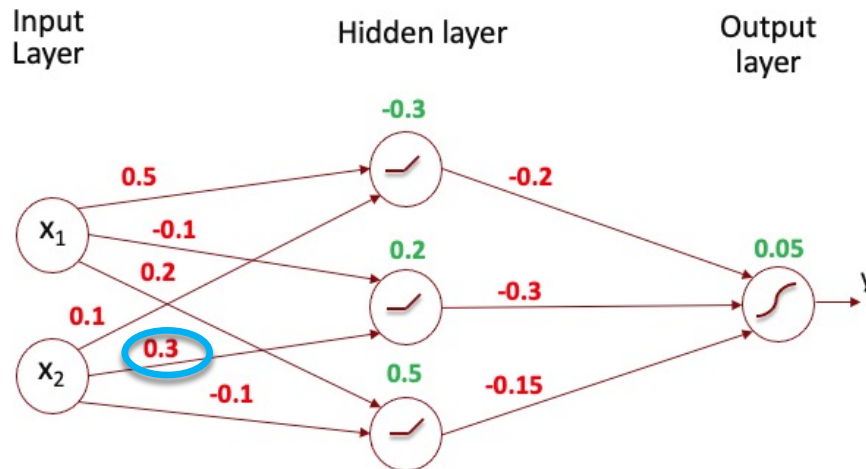
How the gradient of “Accuracy loss” can be zero

Imagine increasing this weight from 0.3 to 0.301, keeping the other weights unchanged.

For every data point x , the predicted value $\text{model}(x)$ will also change but very slightly. But this change will change the classification of x only if the predicted value went from below 0.5 to above 0.5 (or vice-versa).

This “0.5 crossing” will happen only if the predicted value for an x is very close to 0.5 before the change. This is very unlikely to be the case for most/all the points.

As a result, the predicted classifications won't change for most/all points! And therefore, the “accuracy loss” won't change either ==> **the partial derivative of the loss with respect to the weight we changed will be 0.0 => gradient descent will stop!**



Why Accuracy and MSE aren't ideal loss functions for a probability output

First approach (accuracy):

Set threshold θ (say $\theta = 0.5$) and predict 1 if $\text{model}(x^i) \geq \theta$.

Accuracy loss is 1 if predicted and true outputs differ

Problem: The gradient of the accuracy loss will be zero most of the time. As a result, the weights don't change in the gradient descent step (i.e., *learning won't happen*)

Second approach (MSE):

Average “sum of squared losses” (as in linear regression)

Why Accuracy and MSE aren't ideal loss functions for a probability output

First approach (accuracy):

Set threshold θ (say $\theta = 0.5$) and predict 1 if $\text{model}(x^i) \geq \theta$.

Accuracy loss is 1 if predicted and true outputs differ

Problem: The gradient of the accuracy loss will be zero most of the time. As a result, the weights don't change in the gradient descent step (i.e., *learning won't happen*)

Second approach (MSE):

Average “sum of squared losses” (as in linear regression)

Problem: *The size of the gradient does not reflect how “bad” the predictions are (i.e., very wrong predictions are not penalized as aggressively). As a result, learning happens slowly*