

Neurons and the brain

Tuesday, 19 September 2023 20:49

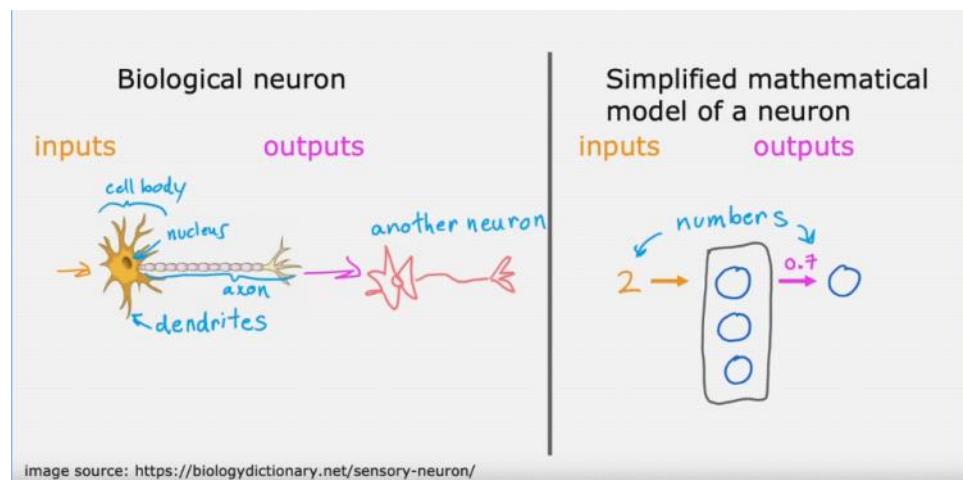
The original goal of neural networks was to mimic the brain

Neural networks were used originally for speech they have advanced from there to things like language processing (NLP)

Today's neural networks have almost nothing to do with how the brain works, just the motivation was from the brain.

Mathematical model of a neuron:

An input is given to a group of neurons to undertake a task, then their output is given to more neurons



Neural networks are now possible thanks to the amount of data now available

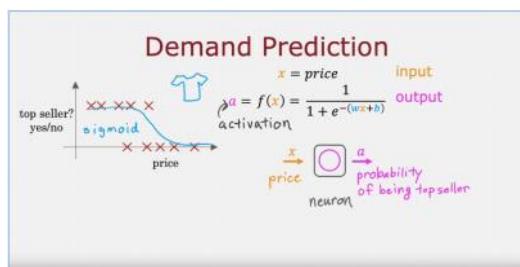
These Neural networks are great if you have a lot of data to use

The uptick in GPU power allowed Neural Networks to be made

Demand Prediction

Wednesday, 20 September 2023 07:02

Using Logistic Regression as a neural network:

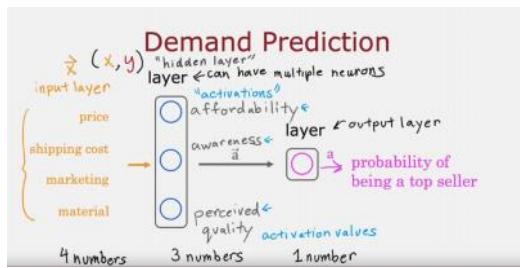


$A = \text{activation}$ is a neurons output



Neuron contains our logistic regression algorithm, so building a neural network means wiring loads of these together

Neural network demand prediction:



Each of these neurons is logistic regression (for this example)

A group of neurons is called a layer

The Input vector is given to first layer of neurons (models)

All layers passed to are "hidden layers" as a "a" vector (activation values)

We cannot decide what is produced by hidden layers nor what weighting different neurons it gives it features

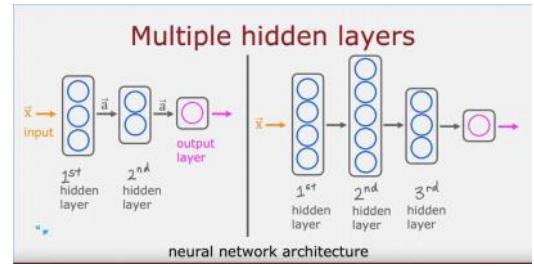
They produce activation values (a vector)

Which is sent from layer to layer (remember we don't know what the different features / inputs are)

We can only decide the amount of layers and neurons

Each neuron learns via gradient descent to fulfill its objectives

Although what features are in hidden layers can't be decided, the amount is.



The amount of layers can also be decided

Left 3 activation values in vector-> 2 Vector activation values in vector -> 1 output

So the amount of data used as features for the next layer can be decided

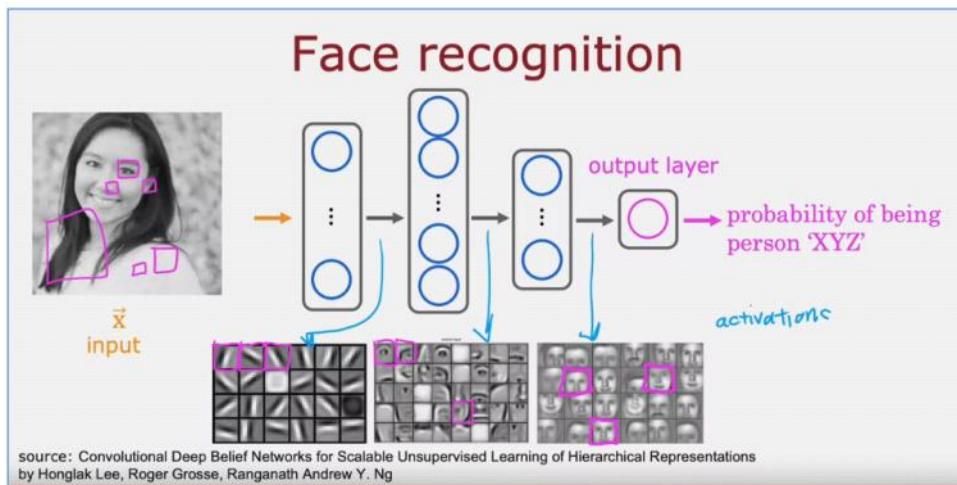
We cannot decide what features will be produced or what weights it gives its inputs by each hidden layer

Example: Recognizing Images

Wednesday, 20 September 2023 08:01

An image is stored in the computer as a matrix of intensity values

There are a lot of these values depending on image size



Each layer does a different job to advance the next layer

First layer looks for lines, second for nose, third for face shape

The neural network will learn what's best for each layer

Each neuron learns its objective to help fulfill the layers objective

Then because it's logistic regression, you get probability at the end

We can't decide what each hidden layer weights features with nor can we decide what they produce

Because we don't decide and the input or output of each layer, it can find the best way to detect a face

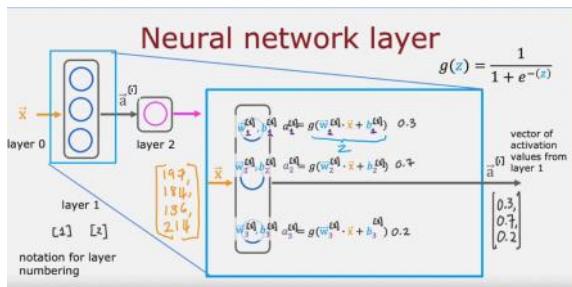
Neural Network Layer (Thresholding)

Wednesday, 20 September 2023 08:44

Square brackets means a layer in a neural network

Threshold are applied after your output layer

Layer 1:



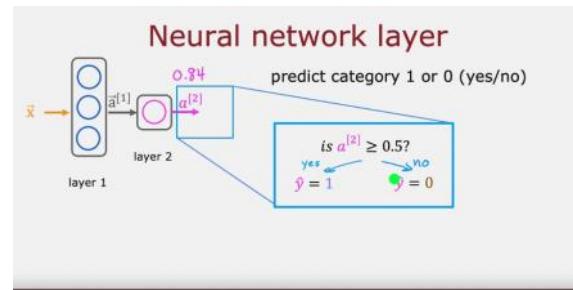
Each Neuron within each layer fulfils an objective when taught from gradient descent

Remember, each layer contains a model (logistic regression etc.)

The output of "a"[1] (activation values from layer one) becomes the input for layer 2

Output layer:

As part of your output layer you can install thresholding



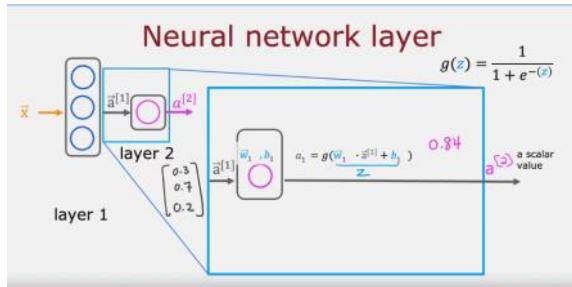
Summary:

Every layer inputs a vector of numbers and applies a model to it.
and passes its results from layer to layer to make predictions

Thresholds are applied after output layer

Layer 2:

Output from layer 1 becomes input for layer 2



This takes in "a" 1 which the activation values from layer 1



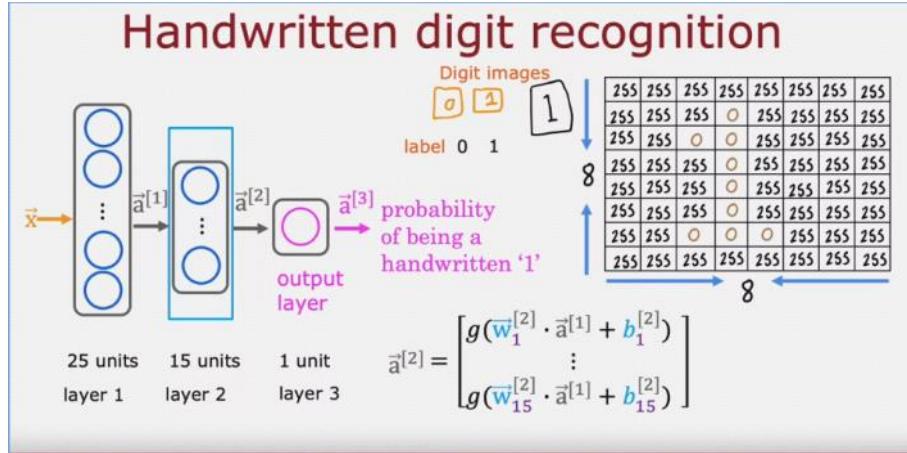
It produces "a" 2 activation values

Inference: making predictions (forward propagation) (Units)

Tuesday, 26 September 2023 06:20

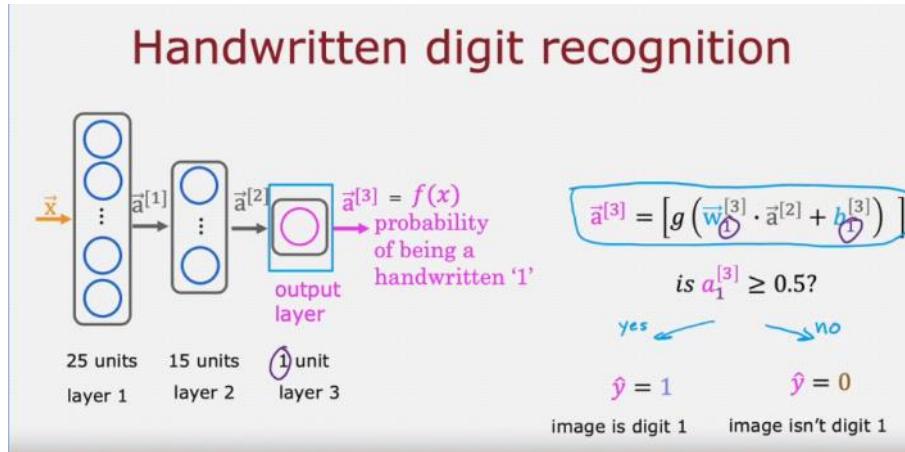
Inference is when you use a trained model

A unit is a neuron



Layer two has 15 units so its computes go to w15 & b15

This is called forward propagation because you compute from left to right through the network using the previous layers activation values



$F(x)$ is your neural networks output

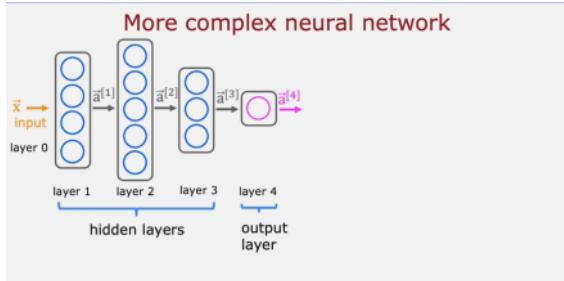
With forward propagation, the number of neurons should decrease as you get closer to output layer (as above)

More Complex Neural Networks (L is for)

Wednesday, 20 September 2023 09:39

Each layer takes in the previous layers activation values

G = activation function (your learning algorithm of choice)



When building neural networks:

L = layer

J = neuron

$$a_j^{[l]} = g(\vec{w}_j^{[l]} \cdot \vec{a}^{[l-1]} + b_j^{[l]})$$

Activation value of layer l , unit(neuron) j

Parameters w & b of layer l , unit j

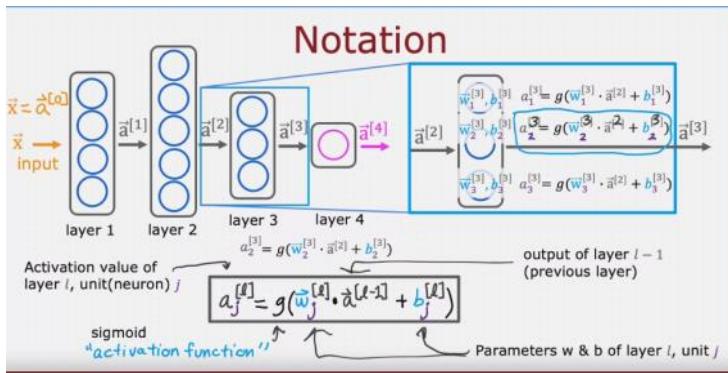
output of layer $l - 1$ (previous layer)

G = activation function (your learning algorithm of choice)

The activation values (output) of each layer of a neural network as shown as "a"[1-4]

Its "a" and then whatever layer 1-4

Notation:



Each layer takes in the previous layers activation values

A[number] = activation values taken in or given out

$$g(\vec{w}_2^{[3]})$$

Bottom right is which neuron in layer (neuron 2)

Top right is which layer overall (layer 3 in this case)

$$\vec{a}^{[2]}$$

This is two because its the activation values from previous layer

Notation Exampled:

$$a_2^{[3]} = g(\vec{w}_2^{[3]} \cdot \vec{a}^{[2]} + b_2^{[3]})$$

G = sigmoid function (logistic regression) or Activation function

$$(\vec{w}_2^{[3]})$$

Layer 3, neuron 2

$$\cdot \vec{a}^{[2]}$$

Activation values from pervious layer (2)

$$b_2^{[3]}$$

Bias for layer 3, neuron 2

Neurons & Layers (Lab)

Tuesday, 26 September 2023 06:53



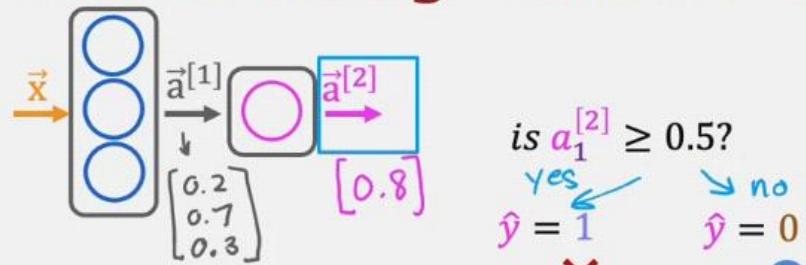
C2_W1_Lab
01_Neuro...

Neurons and Layer Labs

Inference In Code

Tuesday, 26 September 2023 07:12

Build the model using TensorFlow



```
x = np.array([[200.0, 17.0]])
layer_1 = Dense(units=3, activation='sigmoid')
a1 = layer_1(x)

layer_2 = Dense(units=1, activation='sigmoid')
a2 = layer_2(a1)
```

```
if a2 >= 0.5:
    yhat = 1
else:
    yhat = 0
```

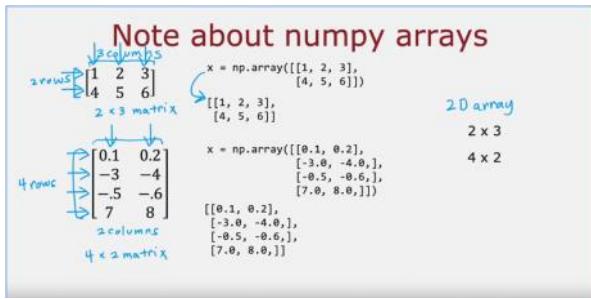
Second layer uses the output of first layer code etc - Forward pagination

Units are the amount of neurons

Then you can threshold after output layer

Data in TensorFlow (how tensorflow stores data)

Tuesday, 26 September 2023 07:38



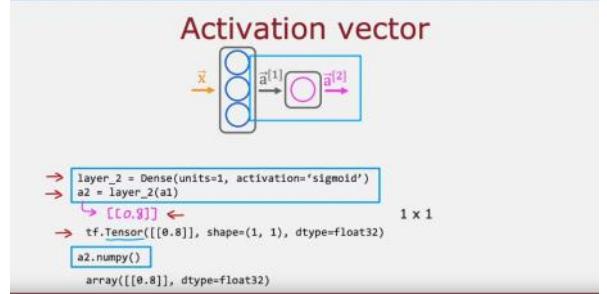
A matrix is a 2d array of data

Matrix can be lots of different sizes (4x2 etc)

Row Vector:



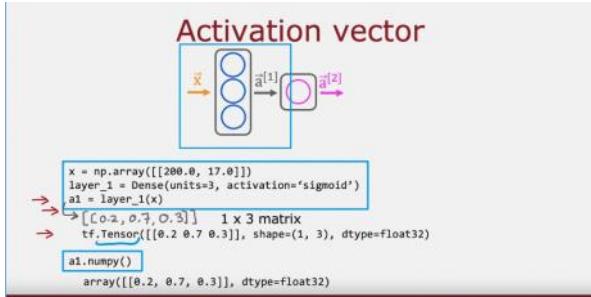
Layer 2 example:



Column Vector:



TesnorFlow wants you to use Matrix's (2d arrays), because its meant to handle large datasets



A1 is activation values for layer 1

2 row 1 column vector is produced by A1

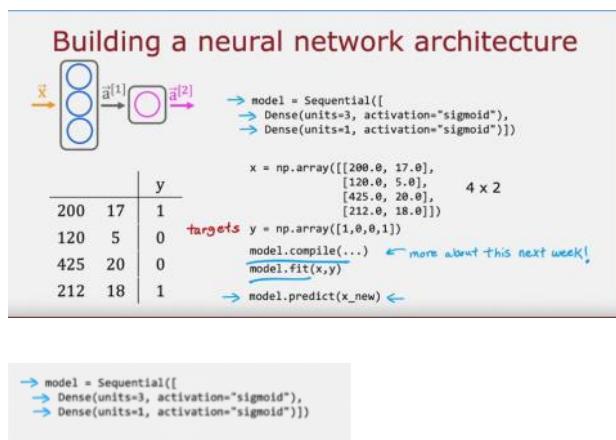
Numpy and tensor arrays are different

A1.numpy() converts tensor matrix to numpy arrys

Building a neural network

Tuesday, 26 September 2023 08:20

Neural network code explained:



Creates a model sequence

```
→ model = Sequential([
    → Dense(units=3, activation="sigmoid"),
    → Dense(units=1, activation="sigmoid")])
```

```
x = np.array([[200.0, 17.0],
    [120.0, 5.0],
    [425.0, 20.0],
    [212.0, 18.0]])
```

targets y = np.array([1,0,0,1])

Your training data in matrix

```
model.compile(...)
```

Compiles your model

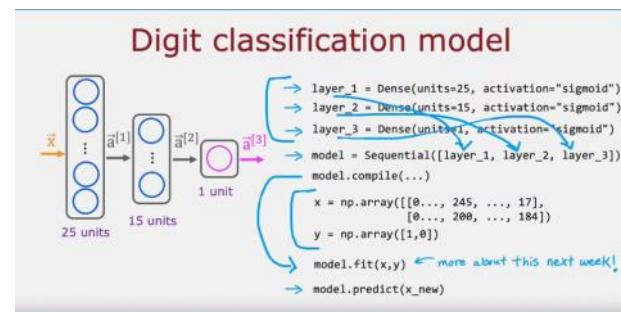
```
model.fit(x,y)
```

Trains with training data

```
→ model.predict(x_new)
```

Does forward propagation, feeds values from previous layer to next (left to right)

Digit Classification model example:



Simple Neural Network (Lab)

Tuesday, 26 September 2023 08:41

Normalization is feature scaling

Epochs controls how many times gradient decent should cover entire database



C2_W1_Lab
02_Coffe...

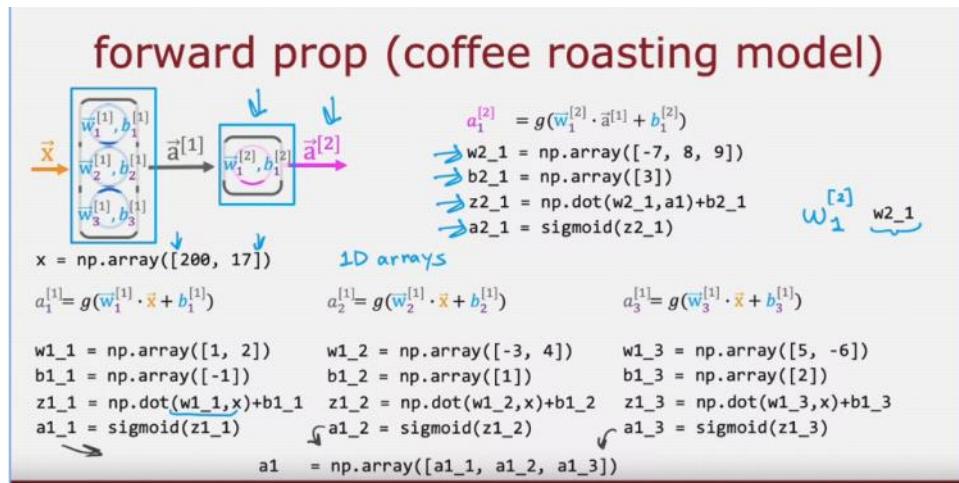


C2_W1_Lab
02_Coffe...

Forward prop in single layer

Tuesday, 26 September 2023 09:52

Non-Tensorflow version:



W2_1 means layer 2, neuron 1's weight etc (underscore means subscript)

Hard coding each neuron for forward prop

General implementation of forward propagation (weight matrix)

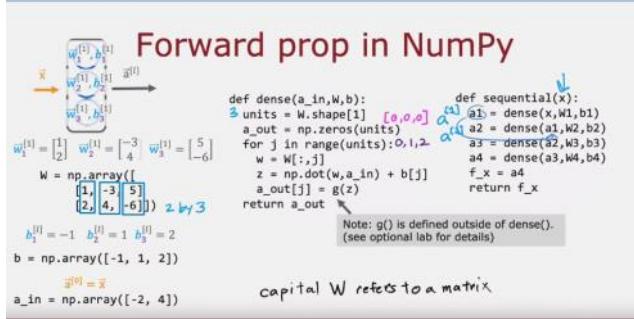
Wednesday, 27 September 2023 14:43

Forward Propagation is to compute your neural network from left to right using the previous layers activation values

Dense function:

Single layer of neural network

Inputs activation values from previous layer then produces next layer activation values



$W = \text{np.array}([\begin{bmatrix} 1, -3, 5 \\ 2, 4, -6 \end{bmatrix }])$ 2 by 3

W should always be capital for matrix

W is loaded with fake data but this is the structure of the matrix

Each column in matrix is the weight vector for a neuron (first column is W1_1 etc)

3 units = W .shape[1]

Equals 3 because of shape of matrix -shape of w just gives us this information

$a_{out} = \text{np.zeros}(units)$

Makes a 3 zero array like units (w.shape)

for j in range(units): 0,1,2

For loop runs for amount of units

$w = W[:, j]$

Gets the weight column for each unit from matrix (first time 1,2, second time -3,4 etc)

$z = \text{np.dot}(w, a_{in}) + b[j]$

Calculates Z for sigmoid function (logistic regression)

$a_{out}[j] = g(z)$

Does activation function

Forward Propagation structure:

```
def sequential(x):
    a1 = dense(x, W1, b1)
    a2 = dense(a1, W2, b2)
    a3 = dense(a2, W3, b3)
    a4 = dense(a3, W4, b4)
    f_x = a4
    return f_x
```

Left to right computing

Taking in previous layers activation values

Simple Neural Network (lab)

Thursday, 28 September 2023 06:28



C2_W1_Lab
03_Coffe...



C2_W1_Lab
03_Coffe...

Speculations On Artificial General Intelligence

Thursday, 28 September 2023 06:39

ANI - Artificial Narrow Intelligence models that do specific task, this has the most progress by far

AGI - Artificial General Intelligence, models that do many tasks, little progress here

Practice Lab: Neural Networks for Handwritten Digit Recognition, Binary

Thursday, 28 September 2023 09:57



C2_W1_Ass
ignment



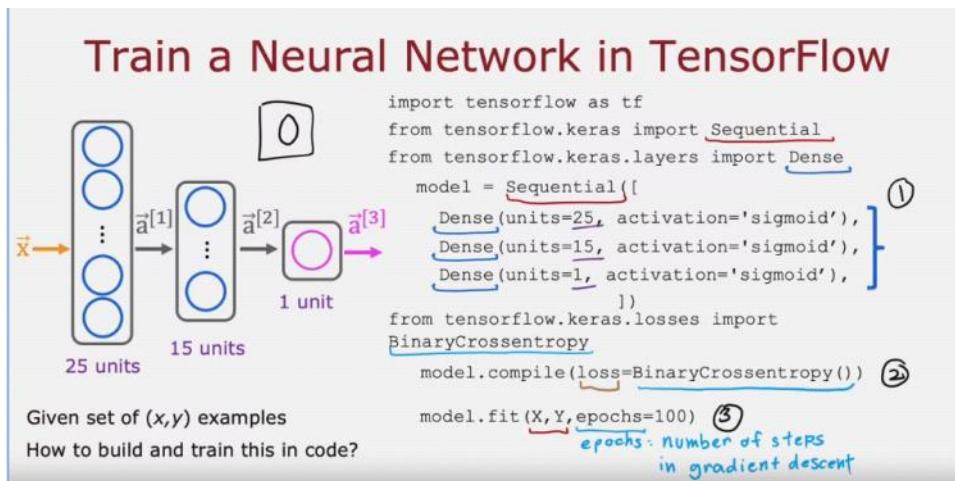
C2_W1_Ass
ignment

Full of really good notes explaining everything

TensorFlow Implementation (Epochs)

Friday, 29 September 2023 06:26

TensorFlow Neural Network + Training:



Epochs is how many steps you take in gradient descent

```
model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='sigmoid'),
])
```

This is your model with forward pagination

```
model.compile(loss=BinaryCrossentropy())
```

Your loss function is your training function

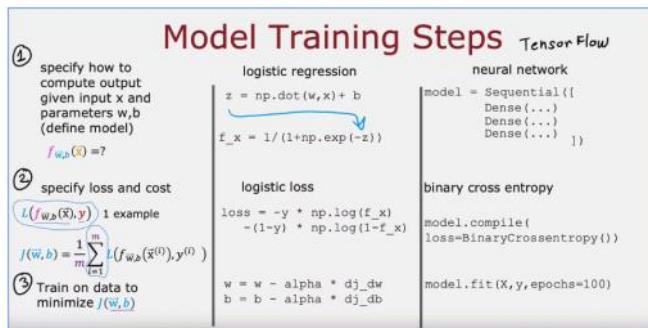
```
model.fit(X, Y, epochs=100) ③  
epochs: number of steps in gradient descent
```

Training the model running 100 steps of gradient descent

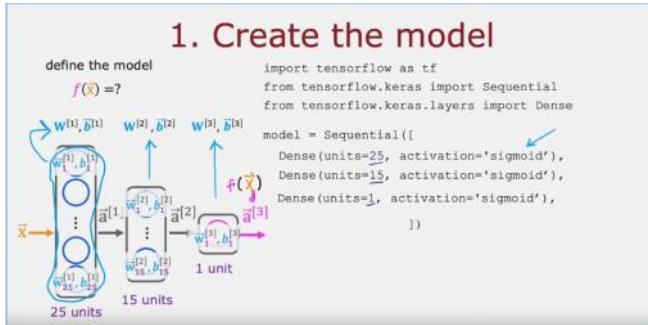
Training Details + Generally making a model

Friday, 29 September 2023 06:38

Neural Network Training:



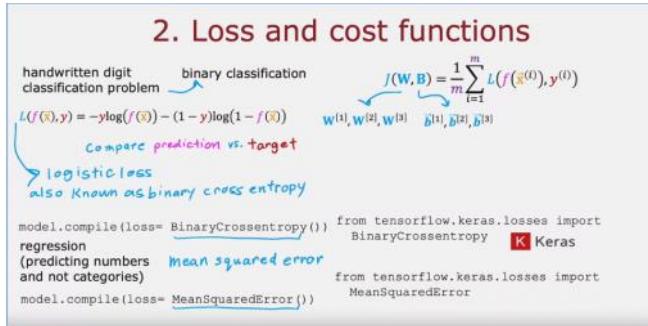
Create Model (1):



Sequential is forward prop

Dense defines each layer along with activation function

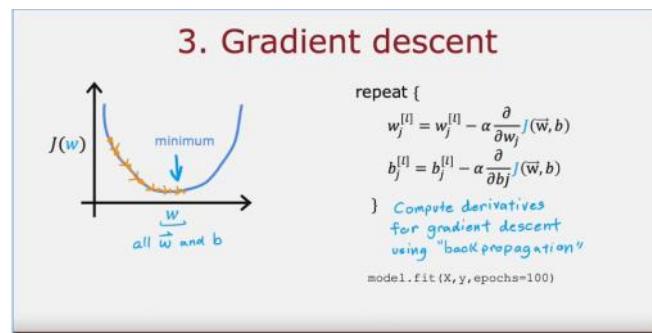
Loss and cost functions (2):



Loss defines your loss function and TensorFlow turns it into cost function

Google for more cost functions named weirdly

Gradient Descent(3):



Model.fit() does gradient descent and back propagation

How Neural Networks are Implemented Efficiently (matmul)

Thursday, 28 September 2023 06:51

Forward Propagation With Vectors:

For loops vs. vectorization

```
x = np.array([200, 17])
W = np.array([[1, -3, 5],
              [-2, 4, -6]])
b = np.array([-1, 1, 2])

def dense(a_in,W,b):
    units = W.shape[1]
    a_out = np.zeros(units)
    for j in range(units):
        w = W[:,j]
        z = np.dot(w, a_in) + b[j]
        a_out[j] = g(z)
    return a_out

[1,0,1]
```

*X = np.array([[200, 17]]) 2Darray
W = np.array([[1, -3, 5], [-2, 4, -6]]) same
B = np.array([-1, 1, 2]) 1x3 2Darray*

vectorized *Z = np.matmul(A_in,W) + B
A_out = g(Z) matrix multiplication
return A_out*

[[1,0,1]]

```
def dense(A_in,W,B):
vectorized Z = np.matmul(A_in,W) + B
A_out = g(Z) matrix multiplication
return A_out
```

matmul - Times Activation values and Weights together because they are both matrix (vectorization)

Matrix multiplication

Thursday, 28 September 2023 07:07

Matrix multiplication works by transposing data and multiplying it by columns

Transposing a matrix:

Lying a matrix out flat into rows from columns

You must do transposing to do matrix multiplication

Matrix Matrix Multiplication:

matrix matrix multiplication

$$A = \begin{bmatrix} 1 & -1 \\ 2 & -2 \\ -1 & -2 \end{bmatrix}$$

rows

$$A^T = \begin{bmatrix} 1 & 2 & -1 \\ 3 & 4 & 5 \\ 6 & 0 & 0 \end{bmatrix}$$

columns

$$W = \begin{bmatrix} 3 \\ 4 \\ 6 \end{bmatrix}$$
$$Z = A^T W = \begin{bmatrix} \leftarrow & \vec{a}_1^T & \rightarrow \\ \leftarrow & \vec{a}_2^T & \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow \\ \vec{w}_1 & \vec{w}_2 \\ \downarrow & \downarrow \end{bmatrix}$$
$$\begin{aligned} \text{row1 col1} &= \begin{bmatrix} \vec{a}_1^T \vec{w}_1 & \vec{a}_1^T \vec{w}_2 \\ \vec{a}_2^T \vec{w}_1 & \vec{a}_2^T \vec{w}_2 \end{bmatrix} \text{row1 col2} \\ \text{row2 col1} &= \begin{bmatrix} (-1 \times 3) + (-2 \times 4) \\ -3 + -8 \end{bmatrix} \text{row2 col2} \\ &= \begin{bmatrix} 11 \\ -11 \end{bmatrix} \end{aligned}$$
$$= \begin{bmatrix} 11 & 17 \\ -11 & -17 \end{bmatrix}$$

m

$$A = \begin{bmatrix} 1 & -1 \\ 2 & -2 \\ -1 & -2 \end{bmatrix}$$
$$A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \end{bmatrix}$$

v

Transposed A (flatten out into rows

matrix matrix multiplication

$$\mathbf{A} = \begin{bmatrix} 1 & -1 \\ 2 & -2 \end{bmatrix}$$

rows

$$\mathbf{A}^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \end{bmatrix}$$

columns

$$\mathbf{W} = \begin{bmatrix} 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$
$$\mathbf{Z} = \mathbf{A}^T \mathbf{W} = \begin{bmatrix} \leftarrow & \vec{a}_1^T & \rightarrow \\ \leftarrow & \vec{a}_2^T & \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow \\ \vec{w}_1 & \vec{w}_2 \\ \downarrow & \downarrow \end{bmatrix}$$

$$\begin{array}{c} \text{row1 col1} \\ \text{row2 col1} \\ (-1 \times 3) + (-2 \times 4) \\ -3 + -8 \\ -11 \end{array} = \begin{bmatrix} \vec{a}_1^T \vec{w}_1 & \vec{a}_2^T \vec{w}_1 \\ \vec{a}_2^T \vec{w}_1 \end{bmatrix} \begin{array}{c} \text{row1 col2} \\ \text{row1 col2} \\ = [11 \quad 17] \end{array}$$

Each row of A (transpose A) is timed by both columns

Matrix Multiplication Rules

Thursday, 28 September 2023 07:35

Matrix Multiplication:

Matrix multiplication rules

$A = \begin{bmatrix} 1 & -1 & 0.1 \\ 2 & -2 & 0.2 \end{bmatrix}$ $A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix}$ $W = \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix}$ $Z = A^T W = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix} \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix}$

$\vec{a}_1^T \vec{w}_1 = (1 \times 3) + (2 \times 4) = 11$ **3 by 4 matrix**

$\vec{a}_3^T \vec{w}_2 = (0.1 \times 5) + (0.2 \times 6) = 1.7$
 $0.5 + 1.2$

$\vec{a}_2^T \vec{w}_3 = (-1 \times 7) + (-2 \times 8) = -23$
 $-7 + -16$

In matrix multiplication results, each row is a Activation value (transposed) matrix, each column is a different weighting matrix

When doing matrix multiplication:

$$A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix} \quad W = \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix}$$

$1 \times 1 + (-1) \times 4 = 11$

Your matrix must be the same LENGTH for dot product to work correctly

Results will always have the same number of rows as Activation Values

Results will always have the same number of columns as W

Matrix Multiplication Code

Thursday, 28 September 2023 08:09

Matrix multiplication in NumPy

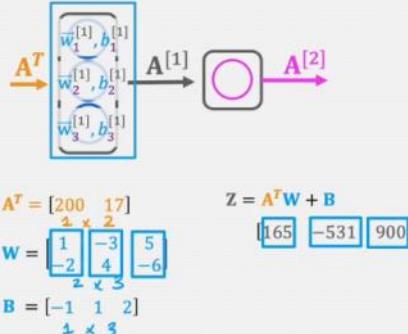
$$A = \begin{bmatrix} 1 & -1 & 0.1 \\ 2 & -2 & 0.2 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix} \quad W = \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix} \quad Z = A^T W = \begin{bmatrix} 11 & 17 & 23 & 9 \\ -11 & -17 & -23 & -9 \\ 1.1 & 1.7 & 2.3 & 0.9 \end{bmatrix}$$

`A=np.array([[1,-1,0.1],
[2,-2,0.2]])` `W=np.array([[3,5,7,9],
[4,6,8,0]])` `Z = np.matmul(AT,W)`
`AT=np.array([[1,2],
[-1,-2],
[0.1,0.2]])` `result = [[11,17,23,9],
[-11,-17,-23,-9],
[1.1,1.7,2.3,0.9]]`
`AT=A.T transpose`

You must transpose your Activation Values before giving it to matmul

Dense Layer Vectorized:

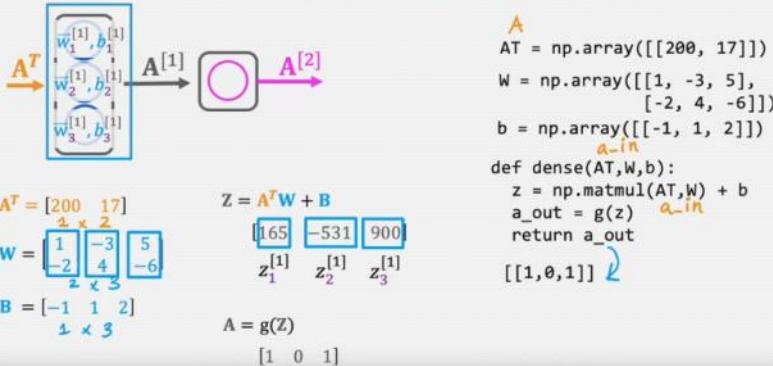
Dense layer vectorized



Each activation value is multiplied by each column of W matrix to get Z

Vectorized Forward Prop:

Dense layer vectorized

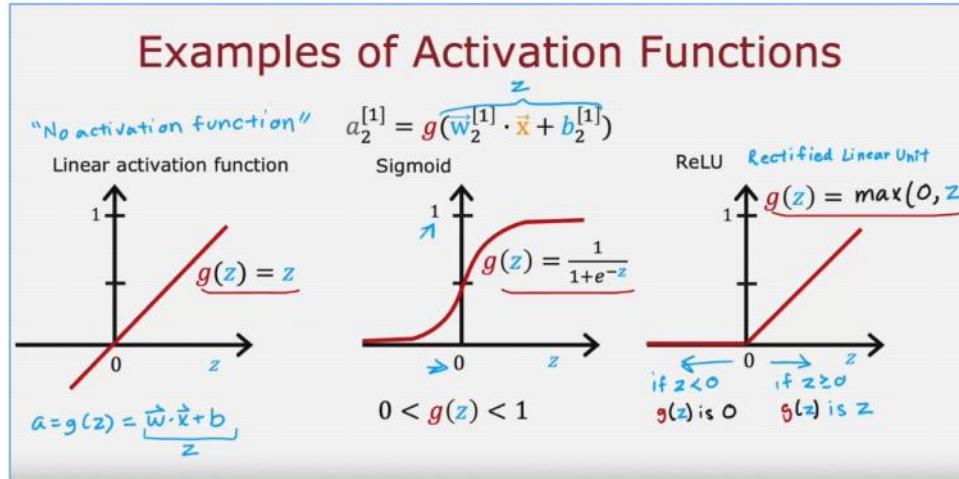


Using matmul to vectorize forward prop

Alternatives To The Sigmoid Activation

Friday, 29 September 2023 07:21

Choices for activation functions:



Linear Activation Function

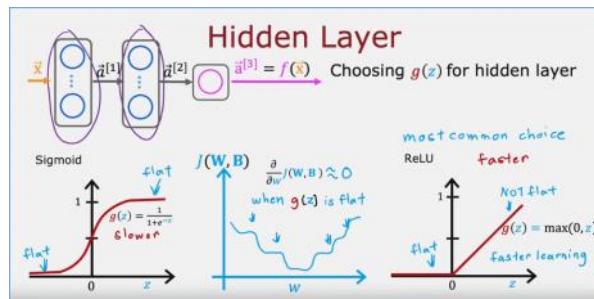
ReLU function

Sigmoid Function

Choosing Activation Functions

Friday, 29 September 2023 07:43

Hidden Layers Activation Functions:

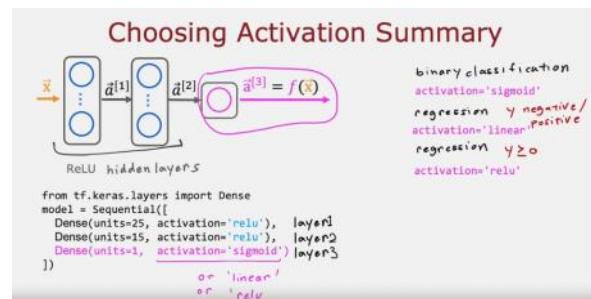


We used to use Sigmoid Functions, now we use ReLU most of the time (apart from output layer)

ReLU is faster than sigmoid when executing

ReLU less flat than sigmoid function, so much faster to train with gradient descent

Summary of Activation functions:



Output layers:

When classifying: use sigmoid

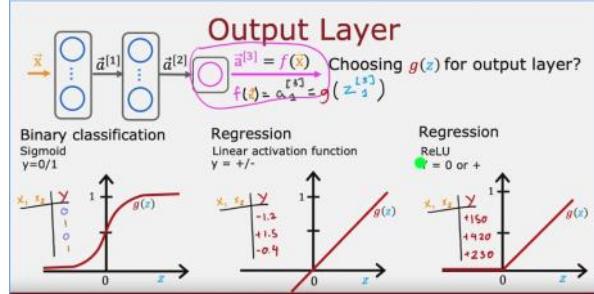
When doing positive or negative: use linear activation function

When doing only positive: use ReLU

Hidden Layers before output:

Just use ReLU

You change your output activation function depending on requirements:



Why Do We Need Activation Function?

Friday, 29 September 2023 08:18

There's no point making a neural network for linear activation

Linear Example

one feature x $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$

w is scalar

$a^{[1]}$

$w_1^{[1]}, b_1^{[1]}$

' a ' is scalar

$g(z) = z$

$a^{[1]} = \underbrace{w_1^{[1]} x}_{\downarrow} + b_1^{[1]}$

$a^{[2]} = \underbrace{w_1^{[2]} a^{[1]}}_{\downarrow} + b_1^{[2]}$

$= w_1^{[2]} (w_1^{[1]} x + b_1^{[1]}) + b_1^{[2]}$

$\vec{a}^{[2]} = (\underbrace{\vec{w}_1^{[2]} \vec{w}_1^{[1]}}_{w}) x + \underbrace{w_1^{[2]} b_1^{[1]} + b_1^{[2]}}_{b}$

$\vec{a}^{[2]} = w x + b$

$f(x) = wx + b$ linear regression

A solely linear activation neural network cannot learn

Never use the linear activation function in hidden layers

ReLU Activation (Lab)

Saturday, 30 September 2023 11:46



C2_W2_Rel
u

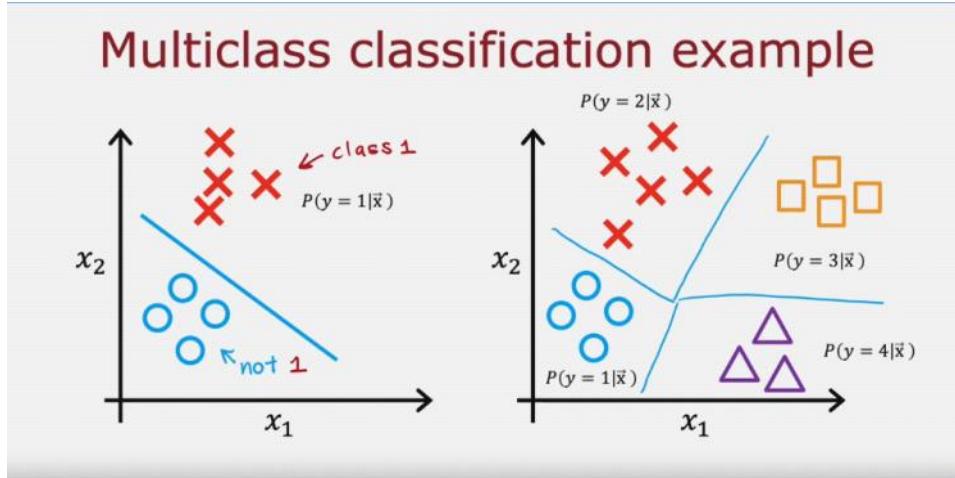


C2_W2_Rel
u - Jupyter...

Multiclass

Friday, 29 September 2023 09:01

Multiclass classification is when Y can be equal to lots of different classes



$P(y = 3 | \vec{x})$

This means what is the chances $Y = 3$ given X etc

Softmax (Z probabilities)

Friday, 29 September 2023 09:11

Softmax regression is based on Linear Regression (sigmoid function)

With Softmax you should calculate Z (input same as Sigmoid) for each of its inputs

Logistic regression (2 possible output values) $z = \vec{w} \cdot \vec{x} + b$		Softmax regression (4 possible outputs) $y = 1, 2, 3, 4$	
$\times \quad a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1 \vec{x})$	0.11	$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$	$P(y=1 \vec{x}) \quad 0.30$
$\circ \quad a_2 = 1 - a_1 = P(y=0 \vec{x})$	0.89	$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$	$P(y=2 \vec{x}) \quad 0.20$
<hr/>		$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$	$P(y=3 \vec{x}) \quad 0.15$
$\square \quad z_4 = \vec{w}_4 \cdot \vec{x} + b_4$		$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$	$P(y=4 \vec{x}) \quad 0.35$
<hr/>		$a_1 + a_2 + \dots + a_N = 1$	

Basically, Softmax calculates the odds that an input belongs in each class

SoftMax formula:

$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j = 1, \dots, N$$

Calculate all your 'Z' inputs for each class first

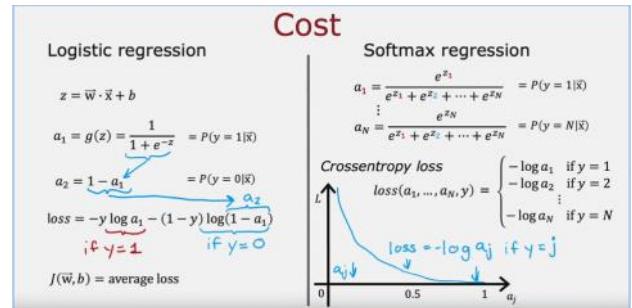
$$\text{parameters } \vec{w}_1, b_1, \vec{w}_2, b_2, \dots, \vec{w}_N, b_N$$

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j|\vec{x})$$

Then input them into the above softmax formula

E = numpy.exp()

Softmax Regression Cost:



Example:

If Y = 2 (example):

Then we get the log of a2

The closer to 2 our prediction was, the less loss

Softmax cost function calculates the loss between the correct category and one the model selected

Neural Network with SoftMax output

Saturday, 30 September 2023 12:13

Softmax calculates all the Z probabilities at once for all activation values

TesnorFlow with Softmax output (there's a better version coming:

MNIST with softmax

① specify the model
 $f_{\bar{w}, b}(\bar{x}) = ?$

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
```

② specify loss and cost
 $L(f_{\bar{w}, b}(\bar{x}), \bar{y})$

```
from tensorflow.keras.losses import
SparseCategoricalCrossentropy
model.compile(loss= SparseCategoricalCrossentropy() )
```

③ Train on data to minimize $J(\bar{w}, \bar{b})$

Note: better (recommended) version later.
Don't use the version shown here!

Improved implantation of Softmax (recommended) (Roundoff)

Saturday, 30 September 2023 12:27

When doing calculations on a computer you get round off error because computer doesn't have unlimited memory

Example is Logistic Regression:

Numerical Roundoff Errors

More numerically accurate implementation of logistic loss: $\frac{1}{1+e^{-z}} - \frac{1}{10,000}$

Logistic regression:

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

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=1, activation='sigmoid')
])
model.compile(loss=BinaryCrossEntropy())
```

Original loss

$$\text{loss} = -y \log(g(z)) - (1-y) \log(1-g(z))$$

```
model.compile(loss=BinaryCrossEntropy(from_logits=True))
```

More accurate loss (in code)

$$\text{loss} = -y \log\left(\frac{1}{1+e^{-z}}\right) - (1-y) \log\left(1 - \frac{1}{1+e^{-z}}\right)$$

logit := z

`model.compile(loss=BinaryCrossEntropy(from_logits=True))`

Softmax Full Model:

MNIST (more numerically accurate)

```
model import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='linear')
])
loss from tensorflow.keras.losses import SparseCategoricalCrossentropy
model.compile(..., loss=SparseCategoricalCrossentropy(from_logits=True))
fit model.fit(X, Y, epochs=100)
predict logits = model(X) ← not a1...a10
f_x = tf.nn.softmax(logits) ← is z1...z10
```

Activation Function Changed ^

Logistic Regression:

logistic regression (more numerically accurate)

```
model model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='linear')
])
loss from tensorflow.keras.losses import BinaryCrossentropy
model.compile(..., BinaryCrossentropy(from_logits=True))
fit model.fit(X, Y, epochs=100)
predict logit = model(X) ← z
f_x = tf.nn.sigmoid(logit)
```

Activation Function Changed ^

Stopping roundoff error during cost function:

Allows TensorFlow to calculate Z as intermediate value reducing roundoff error

Example for Softmax:

More numerically accurate implementation of softmax

Softmax regression

$$(a_1, \dots, a_{10}) = g(z_1, \dots, z_{10})$$

$$\text{Loss} = L(\hat{a}, y) = \begin{cases} -\log(a_1) & \text{if } y = 1 \\ -\log(a_{10}) & \text{if } y = 10 \end{cases}$$

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
model.compile(loss=SparseCategoricalCrossEntropy())
```

More Accurate

$$L(\hat{a}, y) = \begin{cases} -\log\left(\frac{e^{z_1}}{e^{z_1} + \dots + e^{z_{10}}}\right) & \text{if } y = 1 \\ -\log\left(\frac{e^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}}\right) & \text{if } y = 10 \end{cases}$$

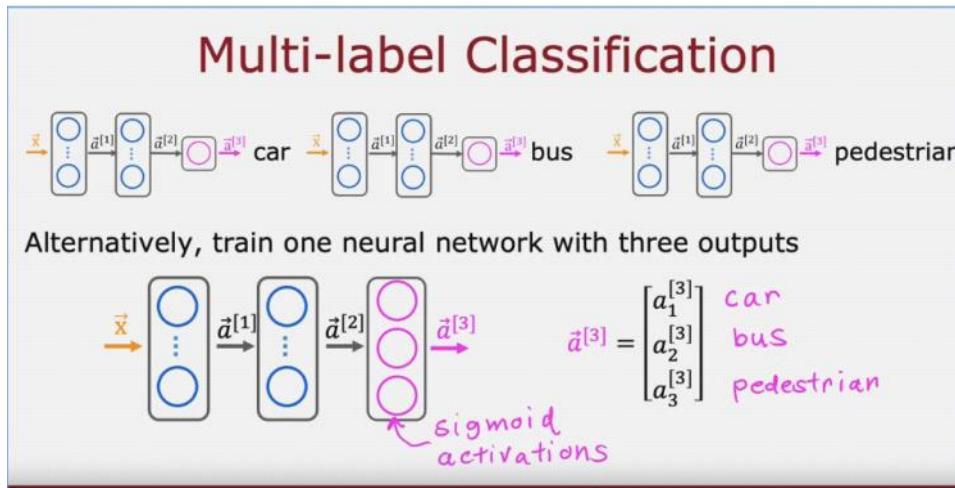
```
model.compile(loss=SparseCategoricalCrossEntropy(from_logits=True))
```

(from_logits=True) Stops roundoff errors due to small numbers with Z

Classification with multiple outputs

Tuesday, 3 October 2023 06:37

Muti label classification is when Y could fit multiple categories

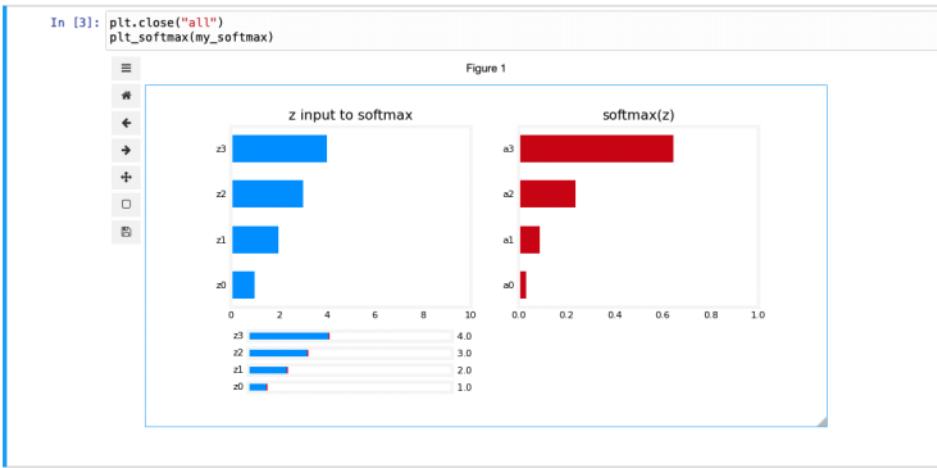


Use a sigmoid function for binary classification of output layer

3 neurons in output layer means 3 outputs overall

Softmax (Lab + Notes)

Tuesday, 3 October 2023 07:04



In Softmax classification your total probabilities must add to 1



C2_W2_Sof
tMax - Ju...



C2_W2_Sof
tMax

Mutliclass (Lab)

Tuesday, 3 October 2023 07:27



C2_W2_Mu
lticlass_TF



C2_W2_Mu
lticlass_T...

Advanced Optimization (Learning rate + bias)

Tuesday, 3 October 2023 07:29

Adam Algorithm Intuition

Adam: Adaptive Moment estimation not just one α

$$\begin{aligned} w_1 &= w_1 - \alpha_1 \frac{\partial}{\partial w_1} J(\vec{w}, b) \\ &\vdots \\ w_{10} &= w_{10} - \alpha_{10} \frac{\partial}{\partial w_{10}} J(\vec{w}, b) \\ b &= b - \alpha_{11} \frac{\partial}{\partial b} J(\vec{w}, b) \end{aligned}$$

The Adam Algorithm gives each feature a separate learning rate including bias

If a parameter (feature) keeps moving in same direction, increase learning rate

If a parameter is "bouncing" decrease learning rate

Code for Adam Algorithm:

MNIST Adam

model

```
model = Sequential([
    tf.keras.layers.Dense(units=25, activation='sigmoid'),
    tf.keras.layers.Dense(units=15, activation='sigmoid'),
    tf.keras.layers.Dense(units=10, activation='linear')
])
```

compile

```
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3),
              loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))
```

$$\alpha = 10^{-3} = 0.001$$

fit

```
model.fit(X, Y, epochs=100)
```

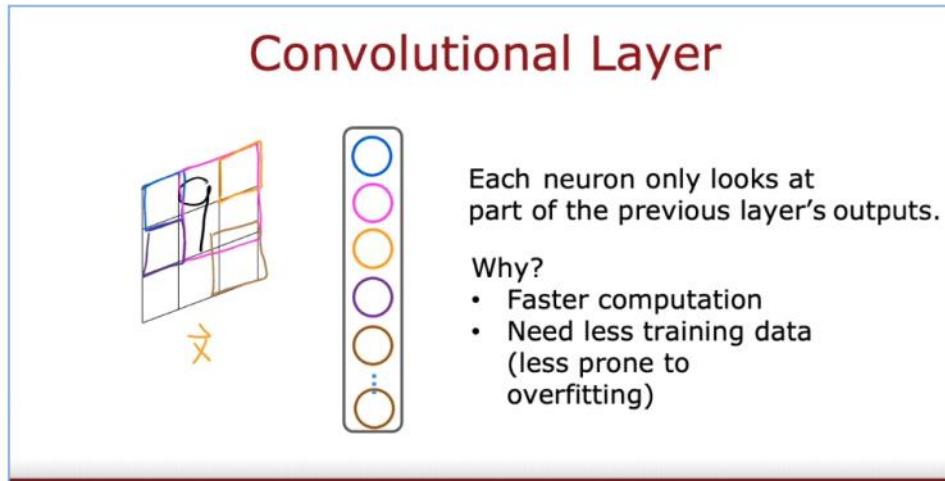
Above we set the initial learning rate

Additional Layer Types (CNN VS RNN)

Tuesday, 3 October 2023 08:11

RNN have the hidden layers who's inputs and outputs can't be controlled.

Convolutional Layer neurons can be limited to only use part of the input data



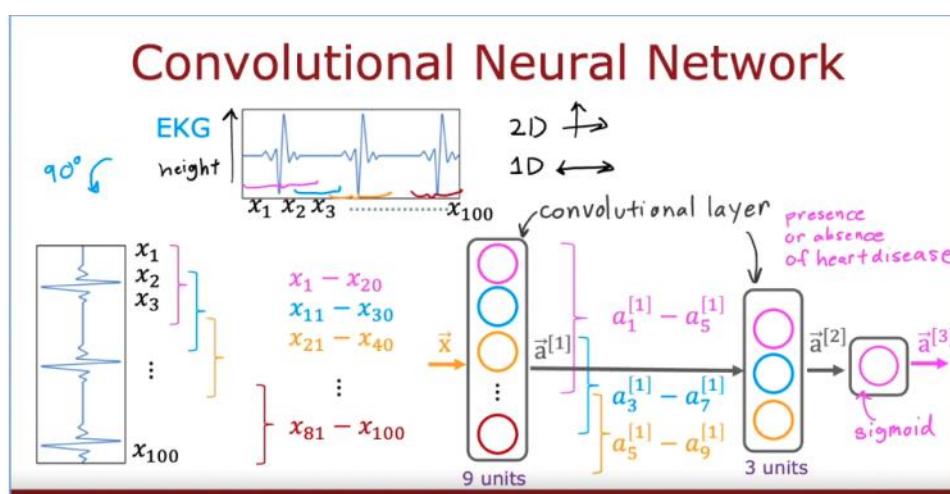
This means:

Faster computation

Less training data required

Less prone to overfitting

Convolutional Neural Network:



Gives us more choice to design, because we decide what areas we want each neuron to focus on

This is done with filters, we decide what pattern each neuron looks for.

Then each neuron gives a score (array of numbers) of how closely the image (example) makes the filter, going over the entire image

As you go deeper into the network the filters get more abstract, first layer might look for a door pattern, 5th layer filter might look for a house.

Make them good at looking for patterns.

What is a derivative? (Epsilon)

Tuesday, 3 October 2023 09:10

A Derivative is a measure of the impact of changing the parameter on a function's output

Epsilon - a very small number we change our function value by to find derivative

More Derivative Examples

$w = 3 \quad J(w) = w^2 = 9 \quad w \uparrow 0.001 \quad J(w) = J(3.001) = 9.006001 \quad \frac{\partial}{\partial w} J(w) = 6$
 $w = 2 \quad J(w) = w^2 = 4 \quad w \uparrow 0.001 \quad J(w) = J(2.001) = 4.004001 \quad \frac{\partial}{\partial w} J(w) = 4$

$w = 2 \quad J(w) = w^2 = 4 \quad w \uparrow 0.001 \quad J(w) = J(2.001) = 4.004001 \quad J(w) \uparrow 4 \times 0.001$

Increasing W by:

$w \uparrow 0.001$

Made J(W) go up by:

≈ 4.004001

Making the derivative 4 because the impact of the parameter change is 4

Code for calculating Derivative:

```
In [2]: J, w = sympy.symbols('J,w')
In [3]: J = w**3
Out[3]: w**3
In [4]: dJ_dw = sympy.diff(J,w)
Out[4]: 3*w**2
In [5]: dJ_dw.subs([(w,2)])
Out[5]: 4
```

Model Function output = J

Weight test = W

This code measures the difference

SSSS

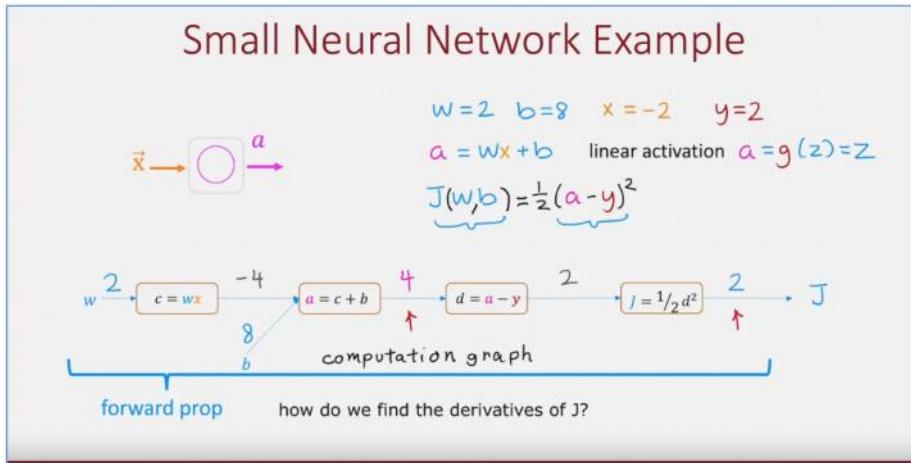
$$\frac{\partial}{\partial w} J(w) \quad w \uparrow \varepsilon \quad J(w) \uparrow k \times \varepsilon$$

Derivative means: How much does J(W)'s output change when W goes up by Epsilon

Computation Graph (back prop)

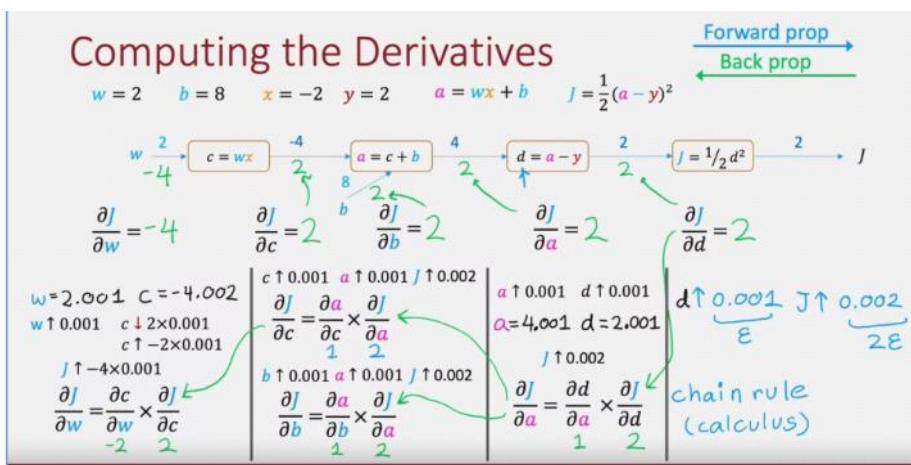
Friday, 6 October 2023 06:31

A computation graph is a set of nodes connected up (in the case of back prop, cost function result)



Back Prop:

Calculate all the different derivatives going backwards (right to left)



Calculate how much impact each calculation had on the next one (going right to left)

Back prop works by working out how much each layer contributed (derivative) the error produced by the cost function

Then we can correct each layer during training to minimize cost

Derivatives (Lab) + Info

Friday, 6 October 2023 06:53



C2_W2_De
rivatives



C2_W2_De
rivatives -...

Propagation Using A Computation Graph (Lab) + Info

Friday, 6 October 2023 06:58



C2_W2_Bac
kprop



C2_W2_Bac
kprop - Ju...

Lab

Thursday, 5 October 2023 16:30

For Softmax calculate Z dominator first always

```
### START CODE HERE ###
tf.keras.Input(shape=(400,)),
```

Above is required when creating a tensorflow model

Note: It is also possible to add an input layer that specifies the input dimension of the first layer. For example:
tf.keras.Input(shape=(400,), #specify input shape
We will include that here to illuminate some model sizing.

Shape is the expected matrix size, 400 dimensional vectors



C2_W2_Ass
ignment



C2_W2_Ass
ignment

Deciding What To Try Next (Fixing)

Friday, 6 October 2023 07:06

There are lots of things you can do to fix a bad model:

Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

$$\rightarrow J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

But it makes unacceptably large errors in predictions. What do you try next?

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features ($x_1^2, x_2^2, x_1x_2, etc$)
- Try decreasing λ ↑ ↑
- Try increasing λ

Running a diagnostic test will tell you what's wrong

Machine learning diagnostic

Diagnostic:

A test that you run to gain insight into what is/isn't working with a learning algorithm, to gain guidance into improving its performance.

Diagnostics can take time to implement
but doing so can be a very good use of your time.

Evaluating a model

Friday, 6 October 2023 07:12

Various ways to evaluate your model:

Splitting your dataset:

Evaluating your model	
Dataset:	
size	price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
30%	
1427	199
1380	212
1494	243

training set → $(x^{(1)}, y^{(1)})$
 \vdots
 $(x^{(m_{train})}, y^{(m_{train})})$
 test set → $(x_{test}^{(1)}, y_{test}^{(1)})$
 \vdots
 $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

m_{train} = no. training examples = 7
 m_{test} = no. test examples = 3

Evaluating a classification problem:

For classification you calculate the fraction of the test set that is mis-classified:

Train/test procedure for classification problem

fraction of the test set and the fraction of the train set that the algorithm has misclassified.

$$\hat{y} = \begin{cases} 1 & \text{if } f_{\vec{w}, b}(\vec{x}^{(i)}) \geq 0.5 \\ 0 & \text{if } f_{\vec{w}, b}(\vec{x}^{(i)}) < 0.5 \end{cases}$$

count $\hat{y} \neq y$

$J_{test}(\vec{w}, b)$ is the fraction of the test set that has been misclassified.

$J_{train}(\vec{w}, b)$ is the fraction of the train set that has been misclassified.

E.g. 61/200

Training set - For teaching the model

Test Set - For testing the model

80/20 split

Calculating Test and Training Error (Linear Regression):

Train/test procedure for linear regression (with squared error cost)

Fit parameters by minimizing cost function $J(\vec{w}, b)$

$$\rightarrow J(\vec{w}, b) = \frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m_{train}} \sum_{j=1}^n w_j^2$$

Compute test error:

$$J_{test}(\vec{w}, b) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (f_{\vec{w}, b}(\vec{x}_{test}^{(i)}) - y_{test}^{(i)})^2$$

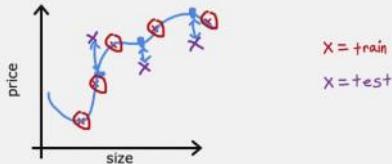

Compute training error:

$$J_{train}(\vec{w}, b) = \frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}_{train}^{(i)}) - y_{train}^{(i)})^2$$

Test Error - Different between the correct answer and your predictions (much more accurate prediction for general error)

Training Error - Different between the training set and the correct answer on data its seen before

Train/test procedure for linear regression (with squared error cost)



$J_{train}(\vec{w}, b)$ will be low

$J_{test}(\vec{w}, b)$ will be high

Model Selection And Training / Cross Validation (choosing)

Friday, 6 October 2023 07:47

You can't just use the test error to see how good a model is, because its often lower / more optimistic than the real error rate

Cross validation Error can be used to pick network design as well once model is chosen

Cross Validation:

Training/cross validation/test set	
size	price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
1380	212
1494	243

validation set
development set
dev set

+ training set 60% → $(x^{(1)}, y^{(1)})$ $m_{train} = 6$
 \vdots
 $(x^{(m_{train})}, y^{(m_{train})})$

cross validation → $(x_{cv}^{(1)}, y_{cv}^{(1)})$ $m_{cv} = 2$
 \vdots
 $(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

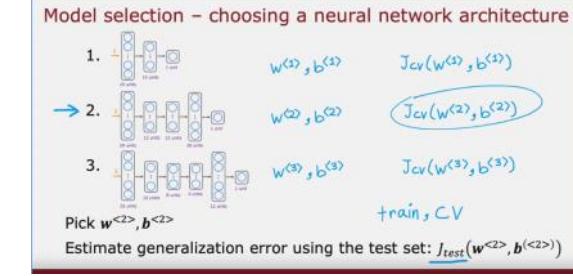
test set 20% → $(x_{test}^{(1)}, y_{test}^{(1)})$ $m_{test} = 2$
 \vdots
 $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Split your dataset into 3

Training set - what you train your model on

Cross validation - for cross checking how valid your model is

Test Set - For testing your model on new data



Pick the one with the lowest (JCV)

Evaluating a model:

Calculate:

Training/cross validation/test set	
Training error:	$J_{train}(\bar{w}, b) = \frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} (f_{\bar{w}, b}(\bar{x}^{(i)}) - y^{(i)})^2$
Cross validation error:	$J_{cv}(\bar{w}, b) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (f_{\bar{w}, b}(\bar{x}_{cv}^{(i)}) - y_{cv}^{(i)})^2$ (validation error, dev error)
Test error:	$J_{test}(\bar{w}, b) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (f_{\bar{w}, b}(\bar{x}_{test}^{(i)}) - y_{test}^{(i)})^2$

Pick the model with the LOWEST Cross Validation error (Jcv)

Model selection	
$d=1$	1. $f_{\bar{w}, b}(\bar{x}) = w_1 x + b$ $w^{(1)}, b^{(1)}$ → $J_{cv}(w^{(1)}, b^{(1)})$
$d=2$	2. $f_{\bar{w}, b}(\bar{x}) = w_1 x + w_2 x^2 + b$ → $J_{cv}(w^{(2)}, b^{(2)})$
$d=3$	3. $f_{\bar{w}, b}(\bar{x}) = w_1 x + w_2 x^2 + w_3 x^3 + b$ \vdots
\vdots	
$d=10$	10. $f_{\bar{w}, b}(\bar{x}) = w_1 x + w_2 x^2 + \dots + w_{10} x^{10} + b$ $J_{cv}(w^{(10)}, b^{(10)})$
→	Pick $w_1 x + \dots + w_4 x^4 + b$ $(J_{cv}(w^{<4>}, b^{<4>}))$
Estimate generalization error using test the set: $J_{test}(w^{<4>}, b^{<4>})$	

Above we are choosing what model to use given there polynomials by there Jcv / Cross Validation Error

Model Evolution And Selection (Lab)

Friday, 6 October 2023 17:33



C2W3_Lab_
01_Mode...



C2W3_Lab_
01_Mode...

Cross validation techniques (10 fold, tuning parameters)

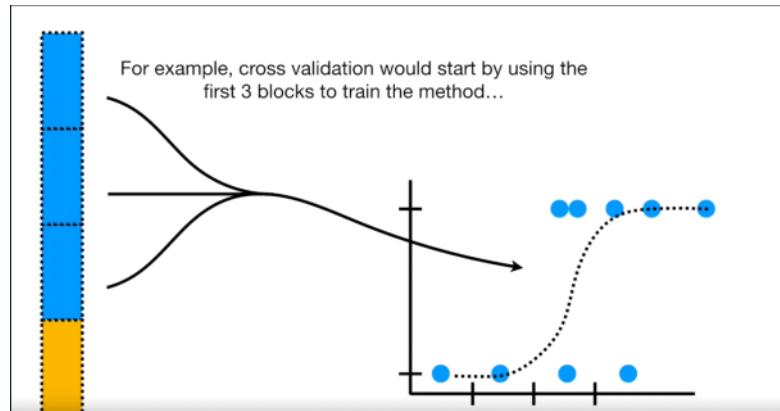
19 November 2024 20:45

Cross validation is a way of testing to see if your model will work well on new data.

Never reuse your training data to test your model it won't work.

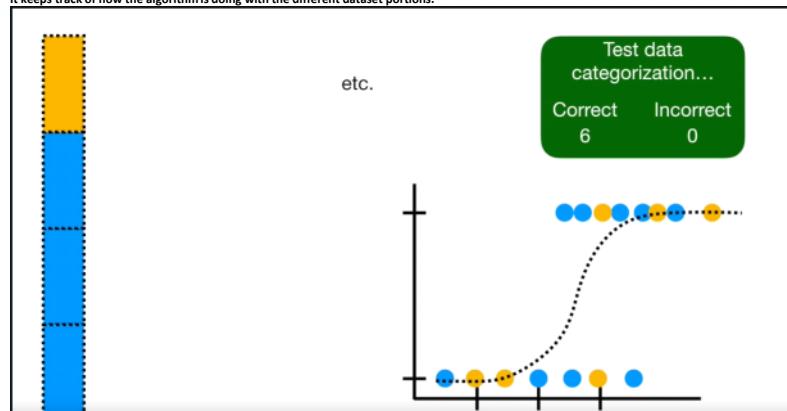
When deciding what data to use for testing you use cross validation.

Cross validation uses all the data one at a time, then summarizes the results at the end



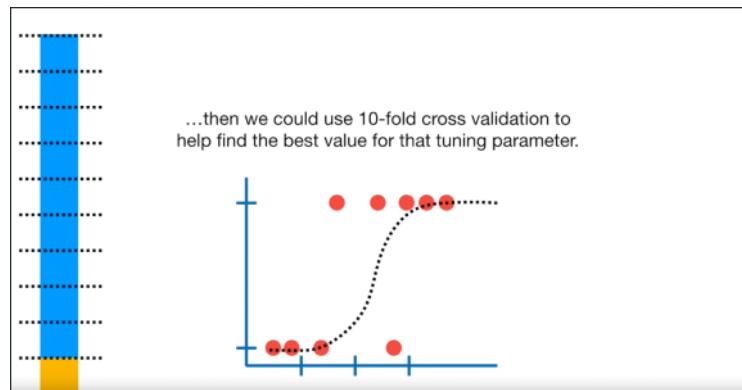
the data is divided into blocks, cross validations might start by using the first 75% for training and the last 25% (4 fold cross validation) for testing as above. Then it will keep changing which blocks it is using.

It keeps track of how the algorithm is doing with the different dataset portions.



Tuning Parameters:

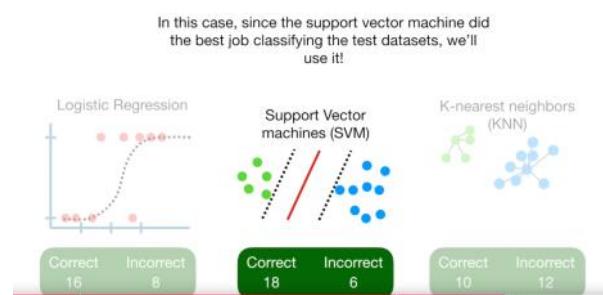
If you are using an algorithm that has parameters that aren't trained, you can use cross validation to compare a choice of values



We can use this to decide which algorithm and/or which configuration of model to use

The score from cross validation depends on how well the model does on all portions (switching testing and training blocks)

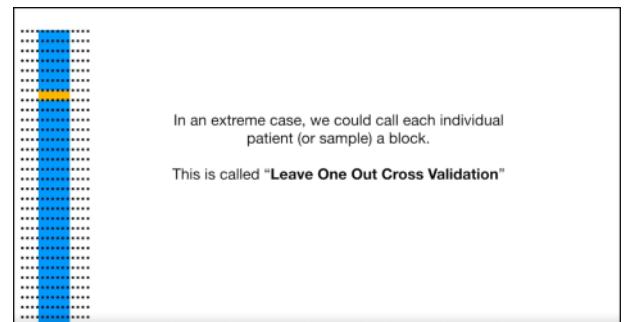
Example for picking a model:



How many blocks the dataset is divided into:

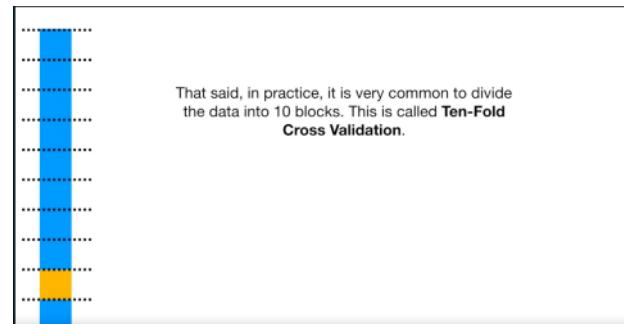
In example to the right, we are doing "Four Fold Cross validation" - Because we have 4 blocks

"Leave One Out Cross validation" - Each sample is tested individually



But most of the time you will use:

"Ten Fold Cross validation" - When the data is divided into 10 blocks



That said, in practice, it is very common to divide the data into 10 blocks. This is called **Ten-Fold Cross Validation**.

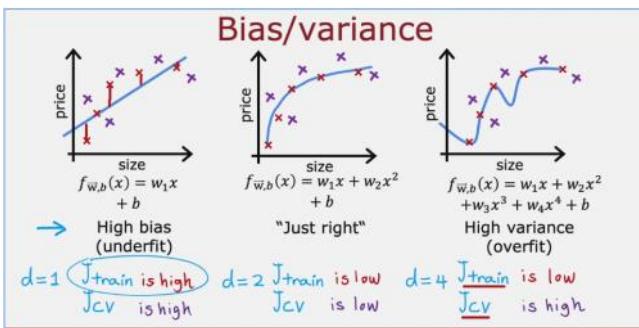
K-folds cross validation

19 November 2024 20:44

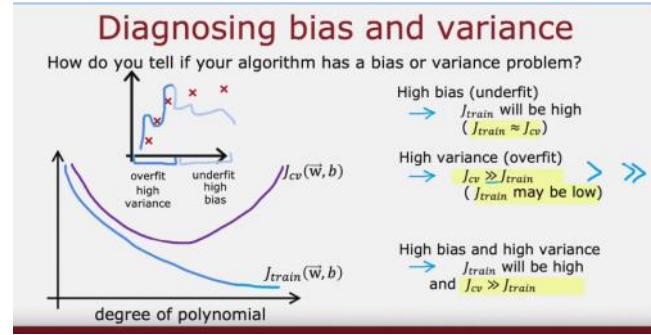
Diagnosing Bias and Variance (check overfit etc)

Friday, 6 October 2023 08:50

This is about working out if your model is under or over fit



Summary:



Underfit:

Train Error: High

Cross Validation: High

"Just right":

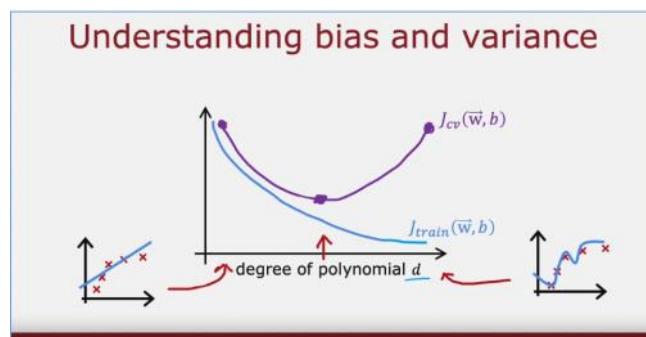
Train Error: Low

Cross Validation: Low

Overfit:

Train Error: Low

Cross Validation: High



J_{train} : As model gets more polynomials it fits dataset better

J_{cv} : Fits better for a while, then gets overfit

Choose the point both lines are closest

Regularization And Bias / Variance (regularization) (New Thing added)

Friday, 6 October 2023 09:28

To fix overfitting (High Variance) DO REGULARIZATION)

We know when our data is overfit because of training error and cross validation patterns (previous note)

Choosing the regularization parameter:

Choosing the regularization parameter λ

Model: $f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

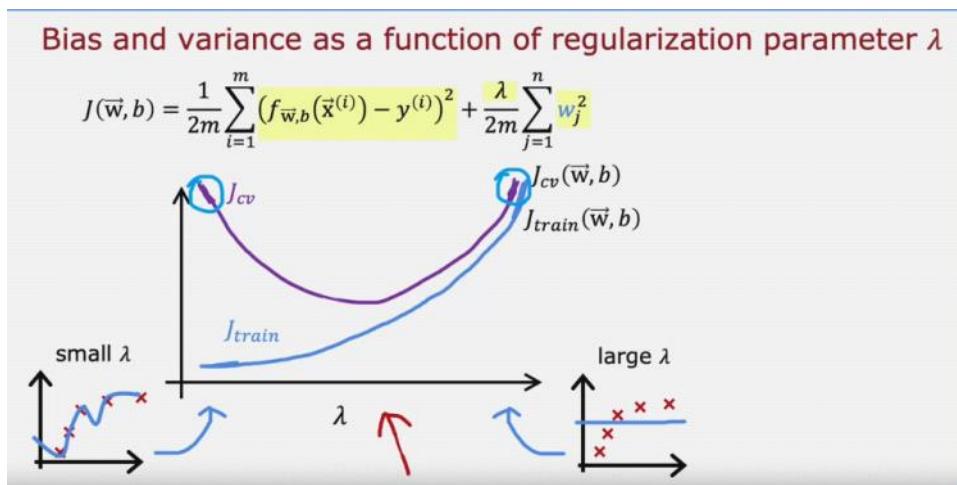
→ 1. Try $\lambda = 0$ → $\min_{\vec{w}, b} J(\vec{w}, b)$ → $w^{<1>} b^{<1>}$ → $J_{cv}(w^{<1>} b^{<1>})$
→ 2. Try $\lambda = 0.01$ → $w^{<2>} b^{<2>}$ → $J_{cv}(w^{<2>} b^{<2>})$
→ 3. Try $\lambda = 0.02$ → $w^{<3>} b^{<3>}$ → $J_{cv}(w^{<3>} b^{<3>})$
→ 4. Try $\lambda = 0.04$ → $w^{<4>} b^{<4>}$ → $J_{cv}(w^{<4>} b^{<4>})$
→ 5. Try $\lambda = 0.08$ → $w^{<5>} b^{<5>}$ → $J_{cv}(w^{<5>} b^{<5>})$
⋮
→ 12. Try $\lambda \approx 10$ → $w^{<12>} b^{<12>}$ → $J_{cv}(w^{<12>} b^{<12>})$

Pick $w^{<5>} b^{<5>}$

Report test error: $J_{test}(w^{<5>} b^{<5>})$

Try different regularization parameters and pick the one with the lowest cross validation

Then calculate the Test set error afterwards to confirm



As lambda (regularization) increases, Jtrain error will increase

Regularization can cause under and over fitting if too big or too small, that's why J_{cv} is the U shape

Pick the one where the lines are closest

Establishing a baseline level of performance (also high bias, high variance gaps)

Monday, 9 October 2023 06:36

Side Note:

To summarise, A model with a high bias error underfits data and makes very simplistic assumptions on it. A model with a high variance error overfits the data and learns too much from it. A good model is where both Bias and Variance errors are balanced.

Performance baseline - what you can expect your model to achieve

Types of baseline:

Establishing a baseline level of performance

What is the level of error you can reasonably hope to get to?

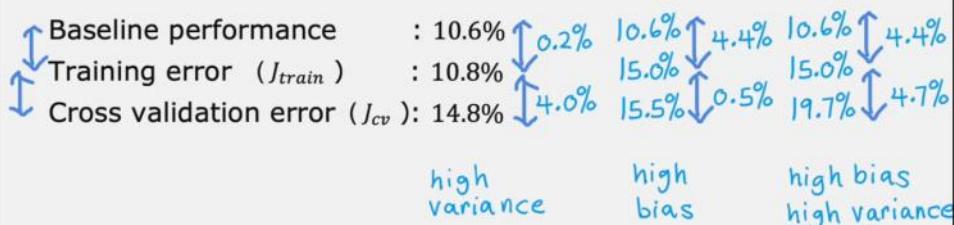
- • Human level performance
- Competing algorithms performance
- Guess based on experience

Human level performance - how well a human does this task

Competing algorithms performance - how well others models have done

Guess based on experience - given previous models you have made

Bias/variance examples



High variance - overfit

High bias - underfit

If gap between baseline and training is smaller than gap between JCV and training - underfit

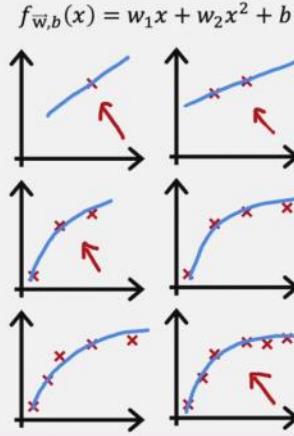
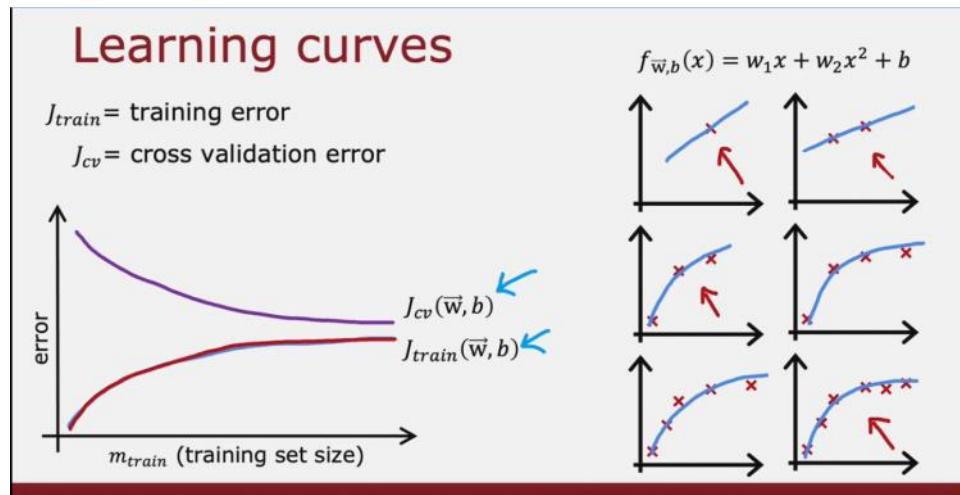
If JCV bigger than gap between baseline and training - overfit

If all 3 are large - high bias + high variance (fits part of the data)

Learning Curves (More data + Jtrain increase)

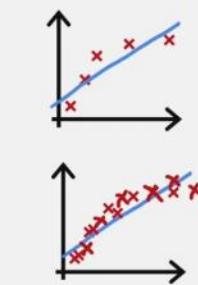
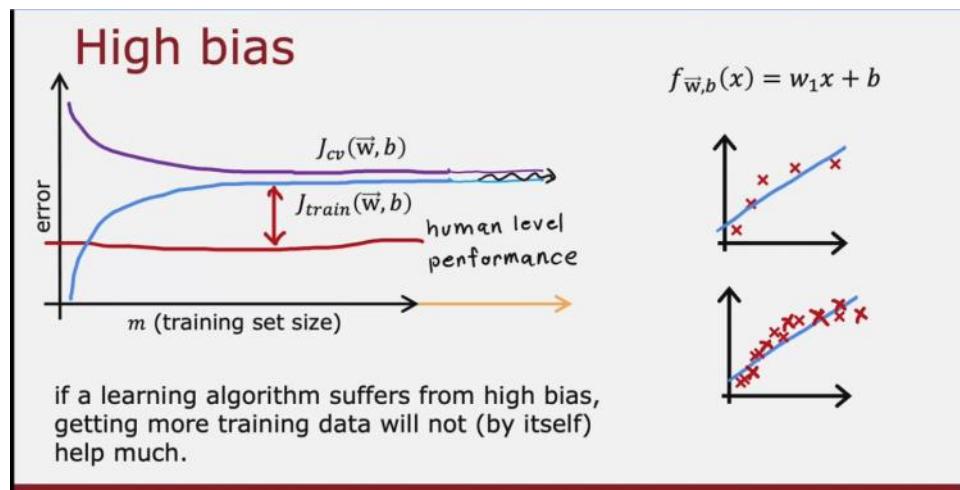
Monday, 9 October 2023 07:02

Learning Curves - Tell you how well your model model is doing given its experience



As you get more training data - your J_{train} gets bigger because harder to fit all the points (high variance)

High bias Learning Curve:



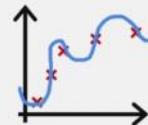
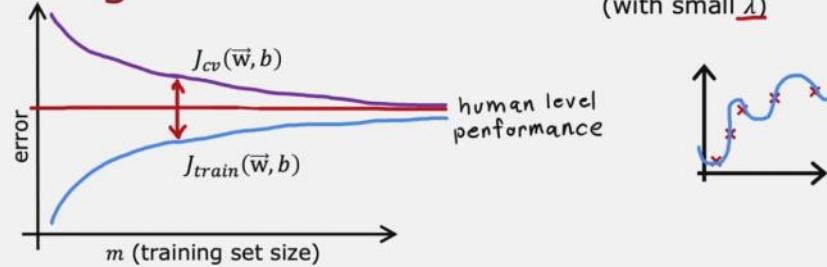
If your model is underfit (high bias) - getting more training data will not fix it alone

High Variance:

High variance

$$f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

(with small λ)



If a learning algorithm suffers from high variance,
getting more training data is likely to help.

If your model has high variance (overfitting) - adding more data can help

Deciding What To Try Next Revisited (Removing data)

Monday, 9 October 2023 07:40

Things to try:

Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

But it makes unacceptably large errors in predictions. What do you try next?

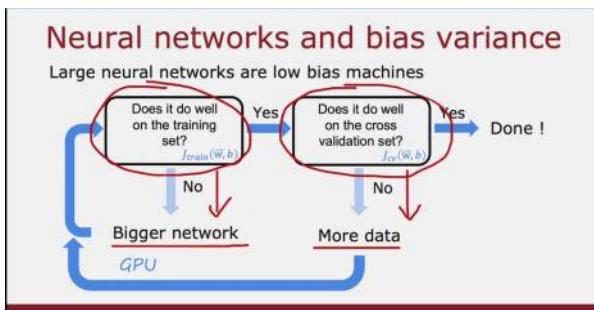
- Get more training examples fixes high variance
- Try smaller sets of features $x, x^2, \cancel{x}, \cancel{x}, \cancel{y} \dots$ fixes high variance
- Try getting additional features fixes high bias
- Try adding polynomial features $(x_1^2, x_2^2, x_1 x_2, \text{etc})$ fixes high bias
- Try decreasing λ fixes high bias
- Try increasing λ fixes high variance

Don't remove training data to fix high bias

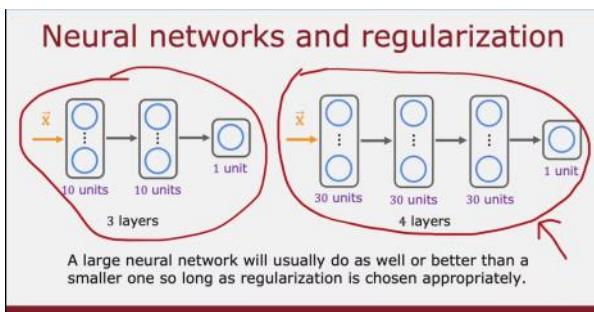
Bias / Variance and Neural Networks (New Line added)

Monday, 9 October 2023 08:02

The "simple recipe":



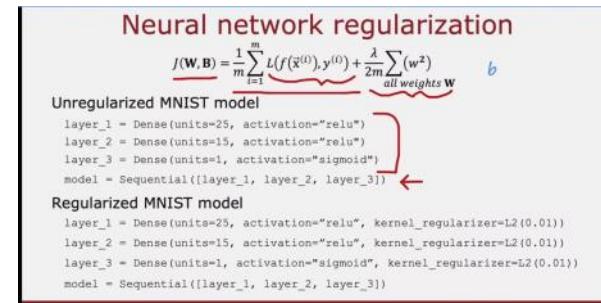
In large Neural Nets, you will be fighting variance (overfitting) mostly



As long as you choose the right neural network larger the better, apart from more expensive hardware wise

Larger Neural Network is better as long as you regularize correctly

Regularizing in TensorFlow:

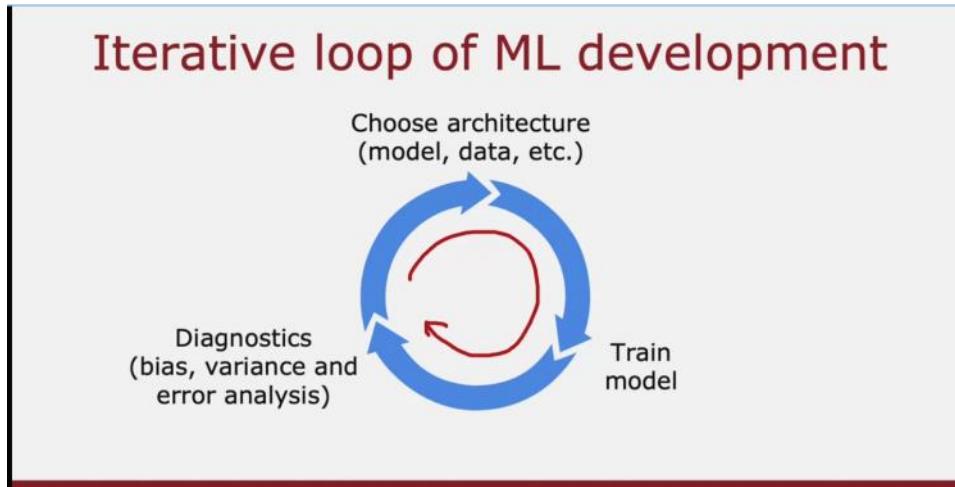


To fix overfitting (High Variance) DO REGULARIZATION

Iterative loop of ML Development

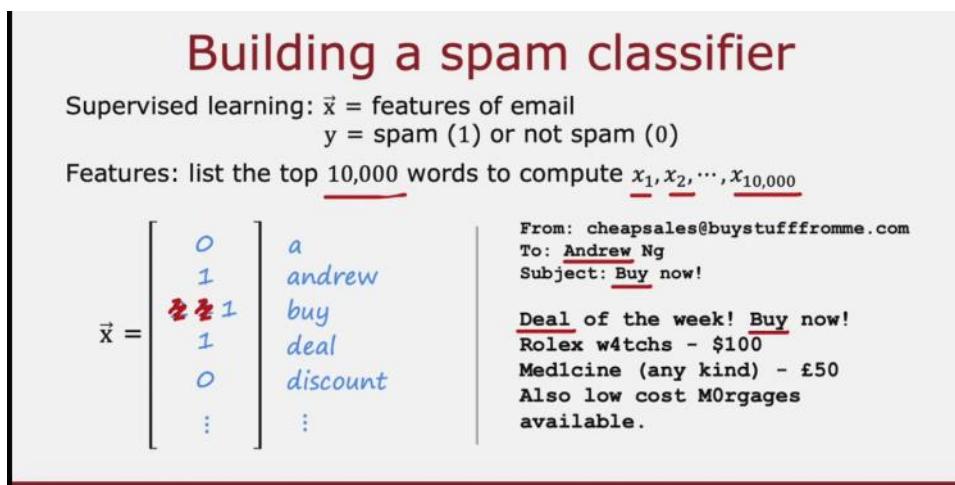
Monday, 9 October 2023 08:37

Development Loop:



Keep going round this

Building Spam Classifier:



Training each feature...

Building a spam classifier

How to try to reduce your spam classifier's error?

- Collect more data. E.g., "Honeypot" project.
- Develop sophisticated features based on email routing (from email header).
- Define sophisticated features from email body.
E.g., should "discounting" and "discount" be treated as the same word.
- Design algorithms to detect misspellings.
E.g., w4tches, med1cine, m0rtgage.

Options for improving (diagnosis)

Error Analysis

Monday, 9 October 2023 09:35

Bias and Variance is the most important thing for diagnosis, this is second.

Getting a sub-set of 100 misclassifications

Error analysis

$m_{cv} = \frac{500}{5000}$ examples in cross validation set.

Algorithm misclassifies $\frac{100}{1000}$ of them.

Manually examine 100 examples and categorize them based on common traits.

Pharma: 21	more data features
Deliberate misspellings (w4tches, med1cine): 3	
Unusual email routing: 7	
Steal passwords (phishing): 18	more data features
Spam message in embedded image: 5	

Allows you to analyze where your problems might be

If one particular category is being misclassified, add more features of it

This allows you to try solutions

Adding Data (data augmentation, data synthesis, data driven)

Wednesday, 18 October 2023 06:50

Adding more data of everything is good, but not always possible

So add data that Error Analysis has found would help

The more labelled data we add, the better

Data augmentation:

Data augmentation

Augmentation: modifying an existing training example to create a new training example.

Data augmentation by introducing distortions

Adjusting existing data to make more examples

Data augmentation for speech

Speech recognition example

- Original audio (voice search: "What is today's weather?")
- + Noisy background: Crowd
- + Noisy background: Car
- + Audio on bad cellphone connection

This can also be done with speech by adding in sounds

Data augmentation by introducing distortions

Distortion introduced should be representation of the type of noise/distortions in the test set.

When doing data augmentation, your data should still be similar to original

Data Synthesis:

Artificial data synthesis for photo OCR

Data synthesis is about creating "fake" data that looks real, like using different fonts (screenshots) to train model

Data centric approach:

Engineering the data used by your system

Conventional model-centric approach:
 $AI = \text{Code} + \text{Data}$
(algorithm/model)
work on this

Data-centric approach:
 $AI = \text{Code} + \text{Data}$
(algorithm/model)
work on this

Engineering your data is as important as the code

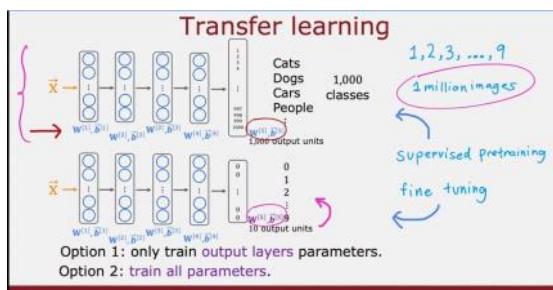
Transfer Learning: Using data from a different task (pre-train + fine tuning)

Wednesday, 18 October 2023 07:21

Transfer Learning:

Step 1 - Supervised pre-training - training an existing neural net trained from an existing dataset on a new dataset on a related task with same input type, just select your output layer

Step 2 - Fine tuning - Training an existing model on a new dataset for a more specific task



Summary:

Transfer learning summary

1. Download neural network parameters pretrained on a large dataset with same input type (e.g., images, audio, text) as your application (or train your own). *1 million images*
2. Further train (fine tune) the network on your own data. *1000 images* *50 images*

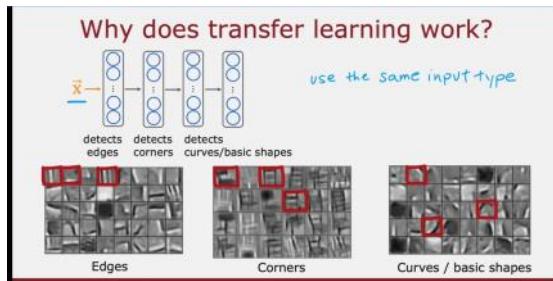
GPT 3 etc was pretrained

When pre-training you make your own output layer to control what you want as output

You can download pre-training models online

Fine tuning

Sometimes pre-training models are available online



The reason pre-training works is because some of the stages, like finding edges, are very common across all models, so can be transferred

Transfer Learning: Using data from a different task (pre-train + fine tuning) 2

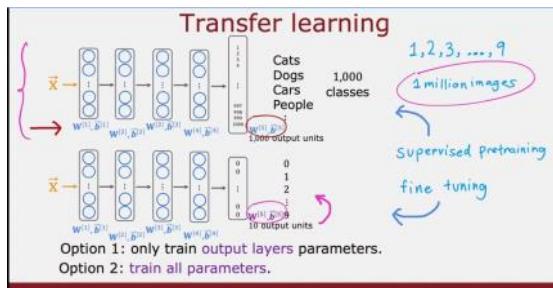
Wednesday, 18 October 2023 07:21

Transfer Learning:

Step 1 - Supervised pre-training - training an existing neural net trained from an existing dataset on a new dataset on a related task just select a new output layer

Step 2 - Fine tuning - Training an existing model on a new dataset for a more specific task - Can give noticeable improvements

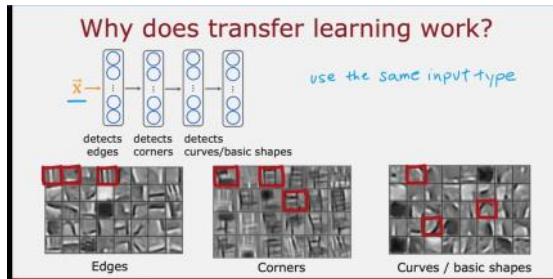
When doing transfer learning you 'freeze' the layers (part from output layer). Fine tuning is the process of unfreezing these are adapting them



When pre-training you make your own output layer to control what you want as output

You can download pre-training models online

Sometimes pre-training models are available online



The reason pre-training works is because some of the stages, like finding edges, are very common across all models, so can be transferred

Summary:

Transfer learning summary

1. Download neural network parameters pretrained on a large dataset with same input type (e.g., images, audio, text) as your application (or train your own). *1 million images*
2. Further train (fine tune) the network on your own data. *1000 images*

50 images

GPT 3 etc was pretrained

Fine tuning:

Fine tuning is when you train all layers in the model, not just change the output layer.

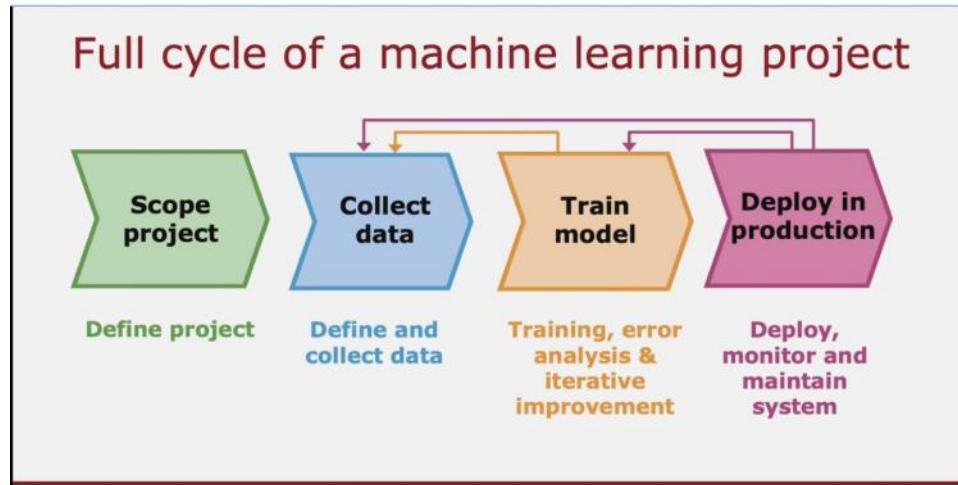
When doing transfer learning you 'freeze' all the layers of the model apart from output layer, ensuring they keep the same weights and bias's.

With fine tuning, once you have done some transfer learning, you unfreeze those layers. To allow them to 'tune' themselves to your new application.

This allows the model to learn the features of the new task while retaining useful knowledge from previous task.

Full Cycle Of A Machine Learning Project (Upgrade by use)

Wednesday, 18 October 2023 08:09



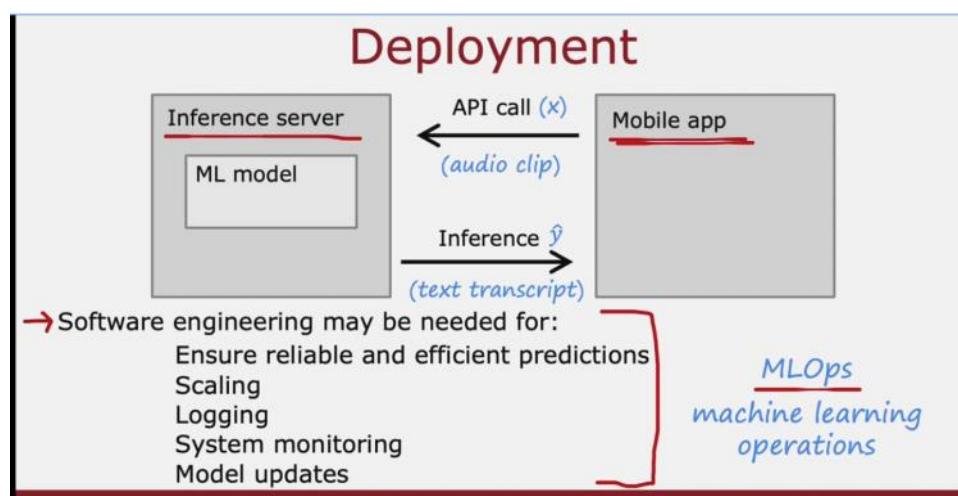
Scope Project - what you want to do

Collect Data - Label data and prepare it for model

Train Model - Train your model (including pre-training) in a loop until low cost

Deploy In Production - Deploy system and use data from it (if allowed) to prefect model

Deployment:



Inference Server - A REST API that utilizes your model

Model updates - taking data from usage attempts and training your model on it

MLOps - is to do all of this, Scaling, Logging, Model updates etc

Fairness, bias, and ethics

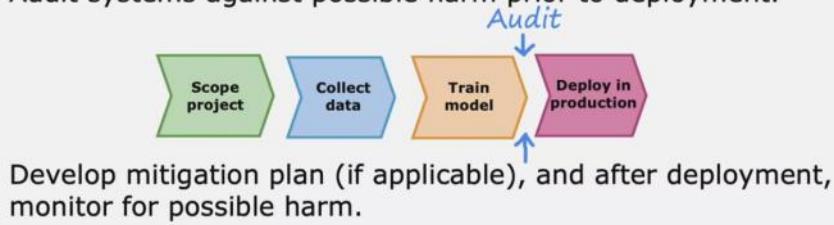
Wednesday, 18 October 2023 08:34

Guidelines

Get a diverse team to brainstorm things that might go wrong, with emphasis on possible harm to vulnerable groups.

Carry out literature search on standards/guidelines for your industry.

Audit systems against possible harm prior to deployment.



You should always carry out and Audit (checking for bias against a group etc) before deployment

Make a mitigation plan - A what if something goes wrong plan (software)

Transfer Learning & Fine tuning (v2)

21 November 2024 21:44

Transfer learning:

Transfer learning is the idea of taking features learned from one task and using them on a new, similar problem.

This is done because there will be common patterns between the problems as they are similar

Often done when you lack data for the new task

Workflow:

1. Take the layers from the existing model
2. Freeze those layers to not destroy information
3. Added new trainable layers on top of the frozen ones. They will learn to turn old features into predictions for the data.
4. Train the new layers on your dataset

Fine tuning:

Fine tuning is the process of unfreezing those pretrained layers and retraining them on the new data at a new low learning rate.

This can make noticeable improvements.

Error Metrics For Skewed Dataset (Precision & Recall)

Wednesday, 18 October 2023 08:50

When you have a dataset that isn't balanced very well, you can't use classification errors (train or test)

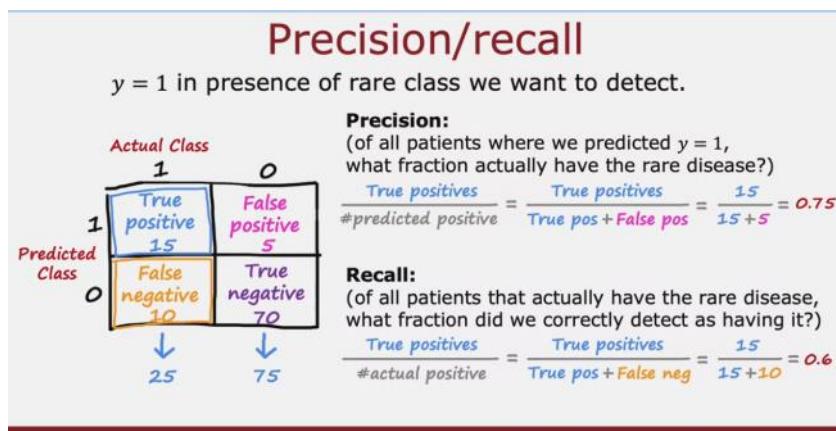
Precision and Recall - Allows you to know how accurate your model is when giving positive readings (Precision) & that the model is detecting most of them (Recall)

Precision - allows you to know how accurate your model is when giving positive readings

Recall - How good the model is at detecting the correct classification overall

The higher Precision and Recall the better

Precision / Recall:



Precision:

(of all patients where we predicted $y = 1$, what fraction actually have the rare disease?)

$$\frac{\text{True positives}}{\#\text{predicted positive}} = \frac{\text{True positives}}{\text{True pos} + \text{False pos}} = \frac{15}{15+5} = 0.75$$

Precision - what fraction of the data is positive ($y=1$)

Recall:

(of all patients that actually have the rare disease, what fraction did we correctly detect as having it?)

$$\frac{\text{True positives}}{\#\text{actual positive}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}} = \frac{15}{15+10} = 0.6$$

Recall - What fraction did our model correctly detect

Trading Off Precision and Recall (Thresholds, Confidence)

Wednesday, 18 October 2023 09:39

Pick the model with the highest Precision and Recall

Trading off precision and recall

Logistic regression: $0 < f_{\vec{w}, b}(\vec{x}) < 1$

- Predict 1 if $f_{\vec{w}, b}(\vec{x}) \geq 0.5$
- Predict 0 if $f_{\vec{w}, b}(\vec{x}) < 0.5$

Suppose we want to predict $y = 1$ (rare disease) only if very confident.

higher Precision, lower recall

Suppose we want to avoid missing too many case of rare disease (when in doubt predict $y = 1$)

lower precision, higher recall

More generally predict 1 if: $f_{\vec{w}, b}(\vec{x}) \geq \text{threshold}$.

precision = $\frac{\text{true positives}}{\text{total predicted positive}}$

recall = $\frac{\text{true positives}}{\text{total actual positive}}$

threshold = 0.99
threshold = 0.01

The higher the threshold, higher the Precision.

Precision is our confidence level when making predictions

The lower the threshold, the higher the Recall, because more examples are detected.

When both Precision and Recall are low, your model isn't very good.

F1 Score:

F1 Score works out how to trade off Precision and Recall

F1 score

How to compare precision/recall numbers?

	Precision (P)	Recall (R)	Average	F ₁ score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.501	0.0392

print("y=1")

~~Average = $\frac{P+R}{2}$~~

$$F_1 \text{ score} = \frac{1}{\frac{1}{2}(\frac{1}{P} + \frac{1}{R})} = 2 \frac{PR}{P+R}$$

The higher the F1 Score the better the model (more balanced the Precision and Recall are)

Practice Lab - Advice For Applying Machine Learning

Sunday, 12 November 2023 20:01

Remember Regularization fixes overfitting (High Variance)



C2_W3_Ass
ignment



C2_W3_Ass
ignment - ...

Decision Tree Model (Lowest...)

Wednesday, 15 November 2023 10:46

Cat Example:

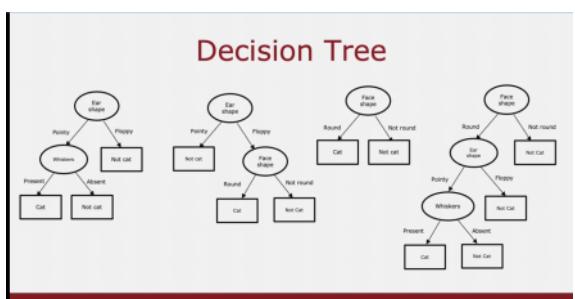
Cat classification example			
Ear shape (x_1)	Face shape (x_2)	Whiskers (x_3)	Cat
Pointy ↗	Round ↗	Present ↗	1
Floppy ↗	Not round ↗	Present ↗	1
Floppy	Round	Absent ↗	0
Pointy	Not round	Present	0
Pointy	Round	Present	1
Pointy	Round	Absent	1
Pointy	Not round	Absent	0
Floppy	Round	Absent	1
Floppy	Round	Absent	0
Floppy	Round	Absent	0

Categorical (discrete values)

X

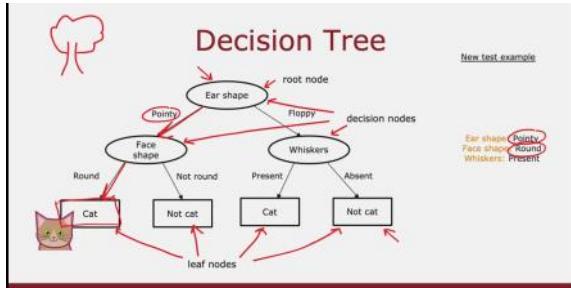
y

Other Tree's:



This is binary classification

The job of a decision tree model is to pick a tree that gets the best JCV, Training error, Testing Error etc



Decision Tree's start off at a root node

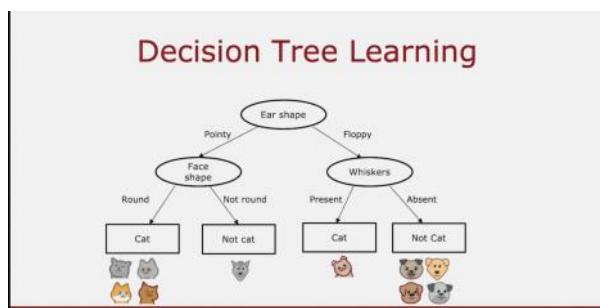
At each layer in the tree a Decision is made leading to a Decision Node

Then we reach our leaf nodes where a prediction is made

Learning Process (max purity, stop splitting)

Wednesday, 15 November 2023 11:05

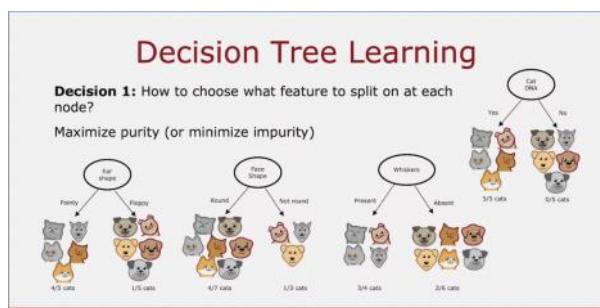
Below uses the example of the cat



The first step is to pick a root node feature (Lets say Ear shape of a cat)

At the next node, Decision node, We will pick another feature (Face shape)

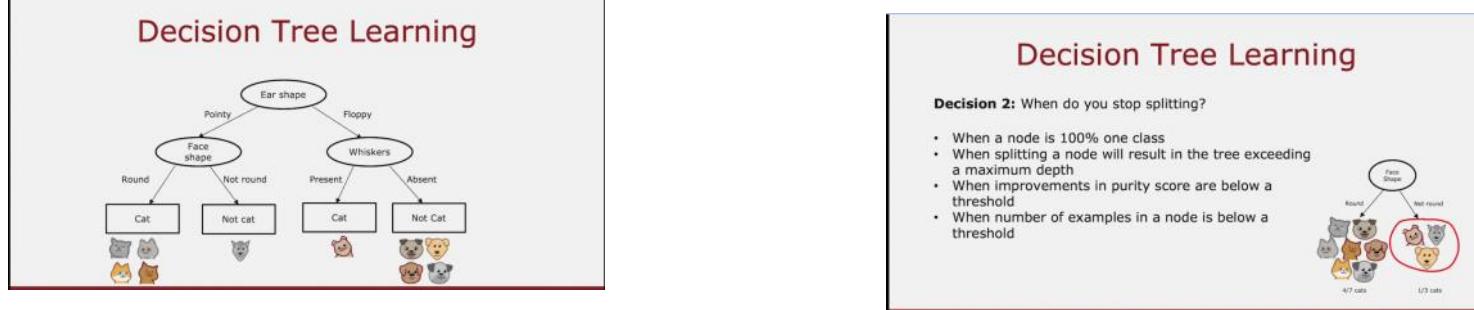
Once we are confident we have made our classification, create a leaf node (Cat, Not Cat)



Maximize Purity - How many of your dataset are the same class

When you reach high purity - Leaf Node (Make your prediction)

When to stop splitting (making choices):



When you reach 100% in one class

Or it reaches a depth limit (how many branches)

Keeping tree small helps limit overfitting

Very small gains in purity

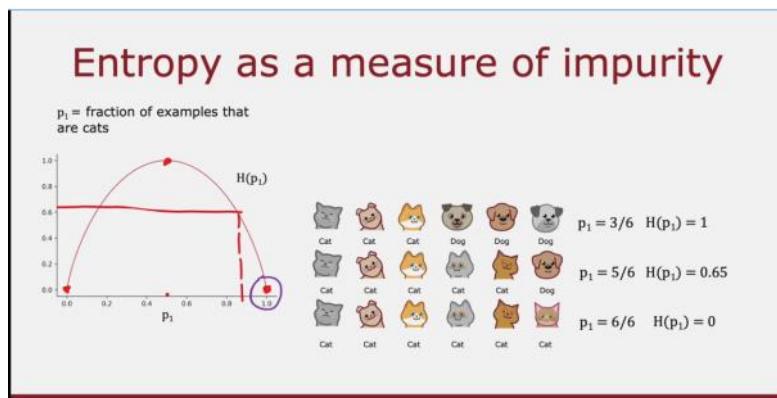
makes it more manageable overall

Number of examples at a node is below a threshold

Measuring Purity (Entropy, Weighted Entropy...each branch)

Wednesday, 15 November 2023 11:37

Entropy as a measure of impurity:

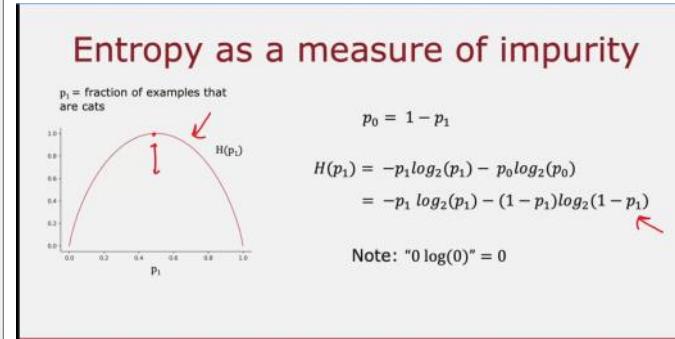


Purity increases the more of your data is of the same class (cats etc)

On the X axis (bottom axis) is fraction of "same" classified data

On Y axis (Up axis) is the Entropy as measure of impurity

The lower the Entropy the better



W = Dataset fraction

H = Entropy

P = Purity

When choosing where to split, use a weighted entropy of each branch to make an average

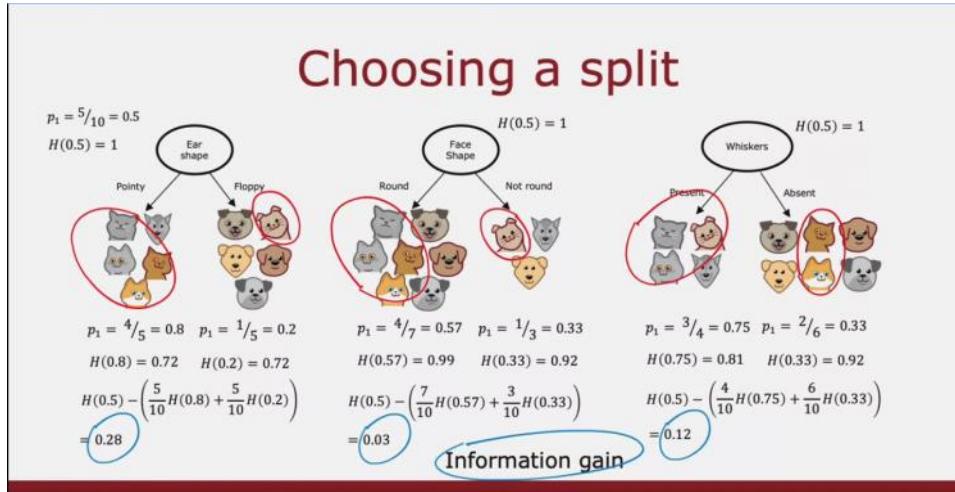
This will require you to calculate variance

Choosing a split: Information Gain (+impacts...)

Monday, 20 November 2023 06:43

The size of the dataset impacts where you split

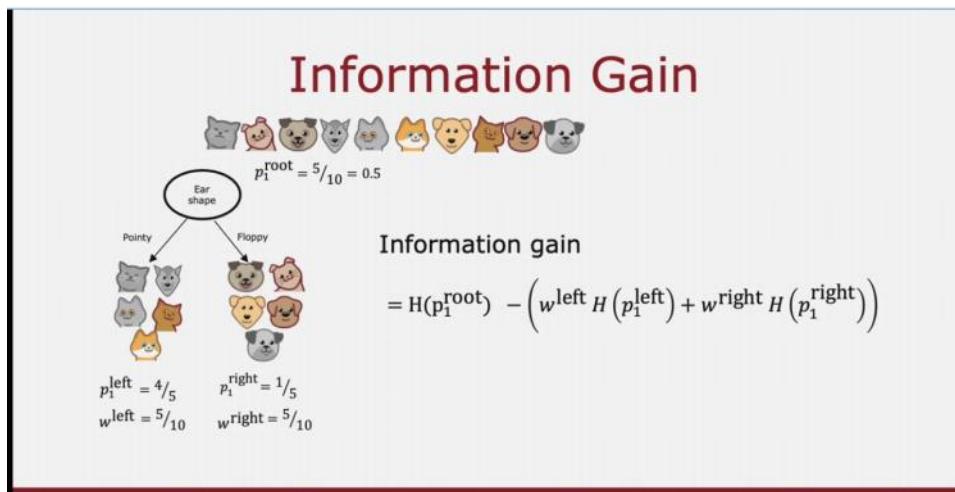
The larger the dataset with the greater the impurity the worse



Information Gain: Increase of purity given decision

Higher information Gain the better

Always choose the split with highest information Gain



W = Dataset fraction

H = Entropy

P = Purity

Root node purity -(dataset fractionLeft * Entropy scoresLeft + dataset fractionRight * Entropy scoresRight)

Putting It Together (info gain...building tree, max depth, Recursive...)

Monday, 20 November 2023 07:01

Decision Tree Process:

Decision Tree Learning

- Start with all examples at the root node
- Calculate information gain for all possible features, and pick the one with the highest information gain
- Split dataset according to selected feature, and create left and right branches of the tree
- Keep repeating splitting process until stopping criteria is met:
 - When a node is 100% one class
 - When splitting a node will result in the tree exceeding a maximum depth
 - Information gain from additional splits is less than threshold
 - When number of examples in a node is below a threshold

Choosing Max Depth:

Higher the max depth, bigger the tree your willing to build

Like picking to use more polynomials, you increase chance of overfitting

Open source libraries can decide this for you quite well

Start with all example at root node,

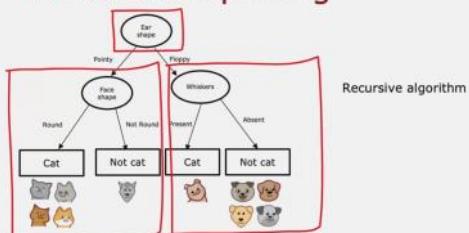
Caluculate information gain for all features, pick highest scoring one

Split dataset accordingly, creating left and right branches

Keep splitting (DOING ABOVE) until :

100% one class or depth limit

Recursive splitting



Recursive algorithm - Just means one tree leads to another

Using One-Hot Encoding For Categorical Features (binary...)

Monday, 20 November 2023 07:19

This approach can be used for "normal" neural networks with minor adjustments:

One hot encoding

Ear-shape	Pointy ears	Floppy ears	Oval ears	Face shape	Whiskers	Cat
Pointy	1	0	0	Round	Present	1
Oval	0	0	1	Not round	Present	1
Oval	0	0	1	Round	Absent	0
Pointy	1	0	0	Not round	Present	0
Oval	0	0	1	Round	Present	1
Pointy	1	0	0	Round	Absent	1
Floppy	0	1	0	Not round	Absent	0
Oval	0	0	1	Round	Absent	1
Floppy	0	1	0	Round	Absent	0
Floppy	0	1	0	Round	Absent	0

One hot encoding

If a categorical feature can take on k values, create k binary features (0 or 1 valued).

One hot encoding and neural networks

Pointy ears	Floppy ears	Round ears	Face shape	Whiskers	Cat
1	0	0	Round	Present	1
0	0	1	Not-round	Present	1
0	0	1	Round	Absent	0
1	0	0	Not-round	Present	0
0	0	1	Round	Present	1
1	0	0	Round	Absent	1
0	1	0	Not-round	Absent	1
0	0	1	Round	Absent	1
0	1	0	Round	Present	1
0	0	1	Round	Absent	0

Switch all features to numbers

Doing this allows to feed to neural networks because they expect numbers as inputs

When we have multiple categories for a feature, turn it into K binary features

The feature that equals "1" is the "hot-one"

We don't need to adjust anything else to make this a decision tree

Lab: Decision Trees

Monday, 20 November 2023 10:11



C2_W4_Lab
_01_Deci...



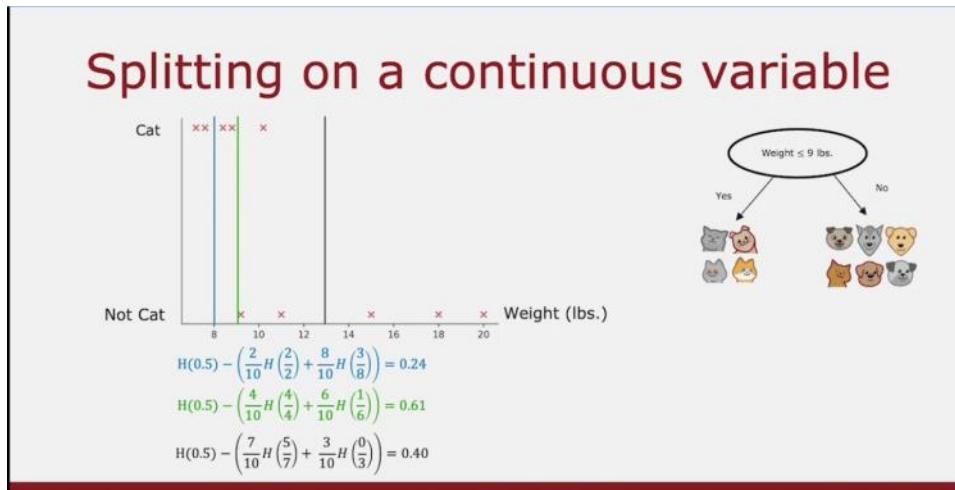
C2_W4_Lab
_01_Deci...

Continuous valued Features

Monday, 20 November 2023 07:50

Splitting on a Continuous Variable:

A continuous variable is a number not a classification but this is not regression because it still predicts a class overall



A good model will decide on a numerical threshold with the highest possible Information Gain

You still do the normal information gain calculation

- Choose the 9 mid-points between the 10 examples as possible splits, and find the split that gives the highest information gain.

You can find the best information gain using midpoints

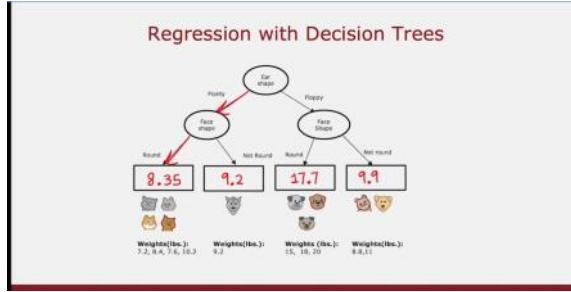
Regression Trees

Monday, 20 November 2023 08:05

Regression Trees predict numbers

Regression with Decision Trees: Predicting a number			
Ear shape	Face shape	Whiskers	Weight (lbs.)
Pointy	Round	Present	7.2
Floppy	Not round	Present	8.8
Floppy	Round	Absent	15
Pointy	Not round	Present	9.2
Pointy	Round	Present	8.4
Pointy	Round	Absent	7.6
Floppy	Not round	Absent	11
Pointy	Round	Absent	10.2
Floppy	Round	Absent	18
Floppy	Round	Absent	20

Predicting an animals weight (y) as a number

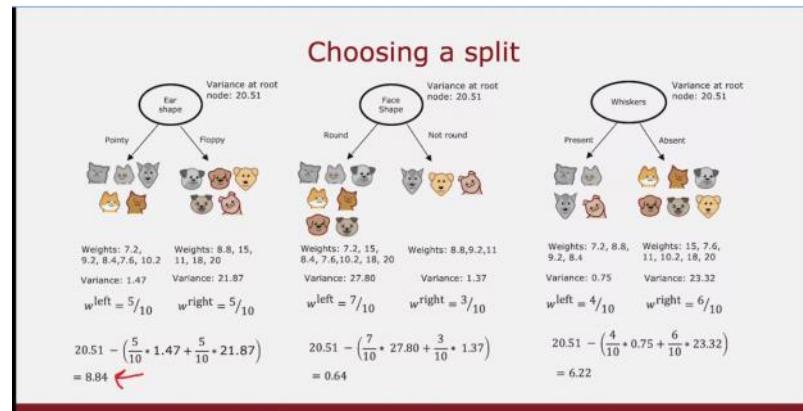


Regression trees take a numerical average of there leaf nodes to make predictions

Choosing a split for regression:

Variance is how much a set of number varies

The goal of our splits is to reduce variances of outputs (animal weights in this case)



As with classification tree, you must get weighted average (left and right)

You must also find root variance

Always pick greatest variance reduction (8.84 in this case)

Keep splitting until depth limit etc

Lab: Tree Ensembles

Tuesday, 21 November 2023 07:45


C2_W4_Lab
_02_Tree...


C2_W4_Lab
_02_Tree...

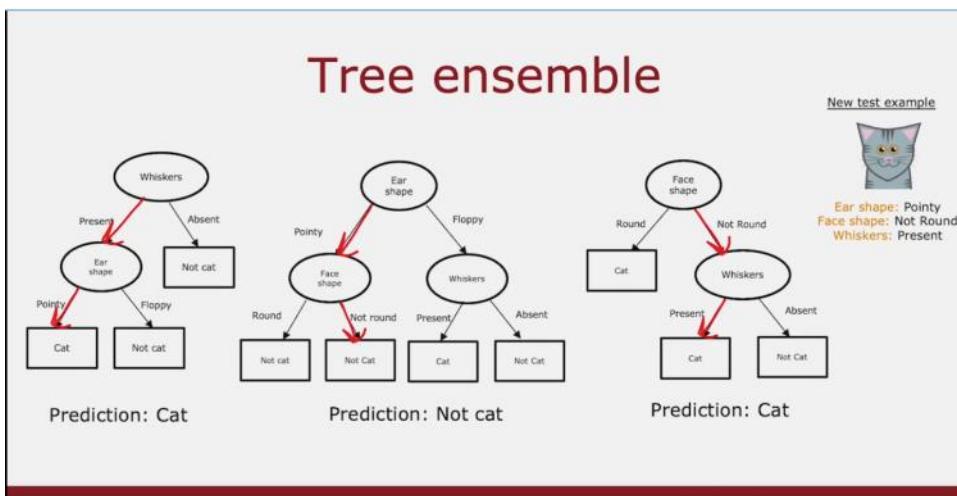
Using Multiple Decision Trees (sensitive, voting)

Monday, 20 November 2023 09:14

Decision Tree models are highly sensitive to changes in dataset



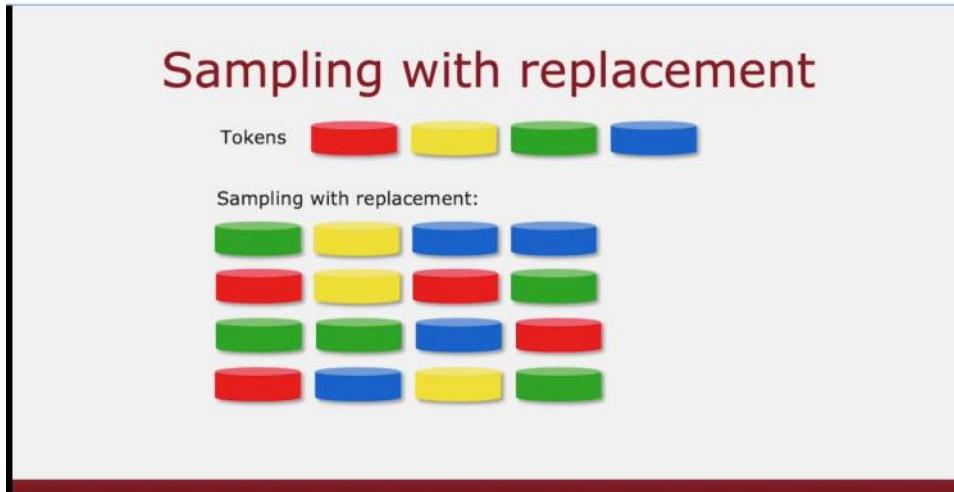
One change in data can lead to a completely different tree



Having a series of trees and having them vote, it makes your model less sensitive to one tree

Sampling With Replacement (pick a token...)

Monday, 20 November 2023 09:26



Sampling With Replacement Means - You put back your "token" after randomly picking it

Sampling with replacement

Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	1
Floppy	Not round	Absent	0
Pointy	Round	Absent	1
Pointy	Not round	Present	0
Floppy	Not round	Absent	0
Pointy	Round	Absent	1
Pointy	Round	Present	1
Floppy	Not round	Present	1
Floppy	Round	Absent	0
Pointy	Round	Absent	1

Using sampling with replacement you can make a new dataset to train a different tree

Due to replacement factor - some examples are the same (this is normal)

Random Forest Algorithm (using sampling...ensemble)

Monday, 20 November 2023 09:37

Generating a tree sample:

Generating a tree sample

Given training set of size m

For $b = 1$ to B

- Use sampling with replacement to create a new training set of size m
- Train a decision tree on the new dataset

Bagged decision tree

Using sampling with replacement you can make a dataset for an ensemble (series of) trees and have them vote

Increasing the amount of trees in the ensemble is always a good thing but does have decreasing returns

B = Recommended amount (64 to 128)

Randomizing the feature choice

At each node, when choosing a feature to use to split, if n features are available, pick a random subset of $k < n$ features and allow the algorithm to only choose from that subset of features.

$$k = \sqrt{n}$$

Random forest algorithm

To make your trees more random (and better) in your ensemble (making them better), force it to pick from a subset of features (k) out of total (n)

This is only required for larger trees

XGBoost

Tuesday, 21 November 2023 06:08

Boosted Tree Intuition (XGBoost):

Boosted trees intuition

Given training set of size m

For $b = 1$ to B :

- Use sampling with replacement to create a new training set of size m
- But instead of picking from all examples with equal $(1/m)$ probability, make it more likely to pick misclassified examples from previously trained trees

Train a decision tree on the new dataset

Eye shape	Fur length	Whiskers	Cat
Pokey	Round	Present	Yes
Flappy	Round	Absent	No
Flappy	Round	Present	No
Flatty	Round	Present	Yes
Flatty	Round	Absent	No
Flatty	Not round	Present	Yes
Flatty	Not round	Absent	No
Flatty	Not round	Present	Yes
Pokey	Not round	Absent	No
Pokey	Not round	Present	Yes
Flatty	Round	Absent	No
Flatty	Round	Present	Yes

```
graph TD; Root((Whiskers)) -- Present --> Cat1[Cat]; Root -- Absent --> EarShape((Ear shape)); EarShape -- Round --> Cat2[Cat]; EarShape -- Not round --> NotRound[Not round]
```

XGBoost in Code:

Using XGBoost

Classification

```
from xgboost import XGBClassifier
model = XGBClassifier()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

Regression

```
from xgboost import XGBRegressor
model = XGBRegressor()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

B = the size of our dataset segment

With Boosted Trees, we use sampling to build a dataset to train new trees, like before

But Boosted tree algorithm's are more likely to pick misclassified examples by previous trees in ensemble

XGBoost (extreme Gradient Boosting):

Open Source

Fast Efficient,

Good choice of default splitting criteria, start + stop

Built in regularization

Highly competitive algorithm for machine learning

When To Use Decision Trees

Tuesday, 21 November 2023 06:34

When to use Decision Trees & Tree ensembles:

Works well on tabular (structured) data

DO NOT use with unstructured data (images, audio etc)

Small Decision Tree's might be human readable

When to use Neural Networks:

Works well with all types of data (structured, unstructured)

May be slower than a decision tree

Works with transfer learning

Can be better paired with other multiple models compared with trees

Programming Assignment: Practice Lab: Decision Trees

Thursday, 23 November 2023 15:54



C2_W4_De
cision_Tr...

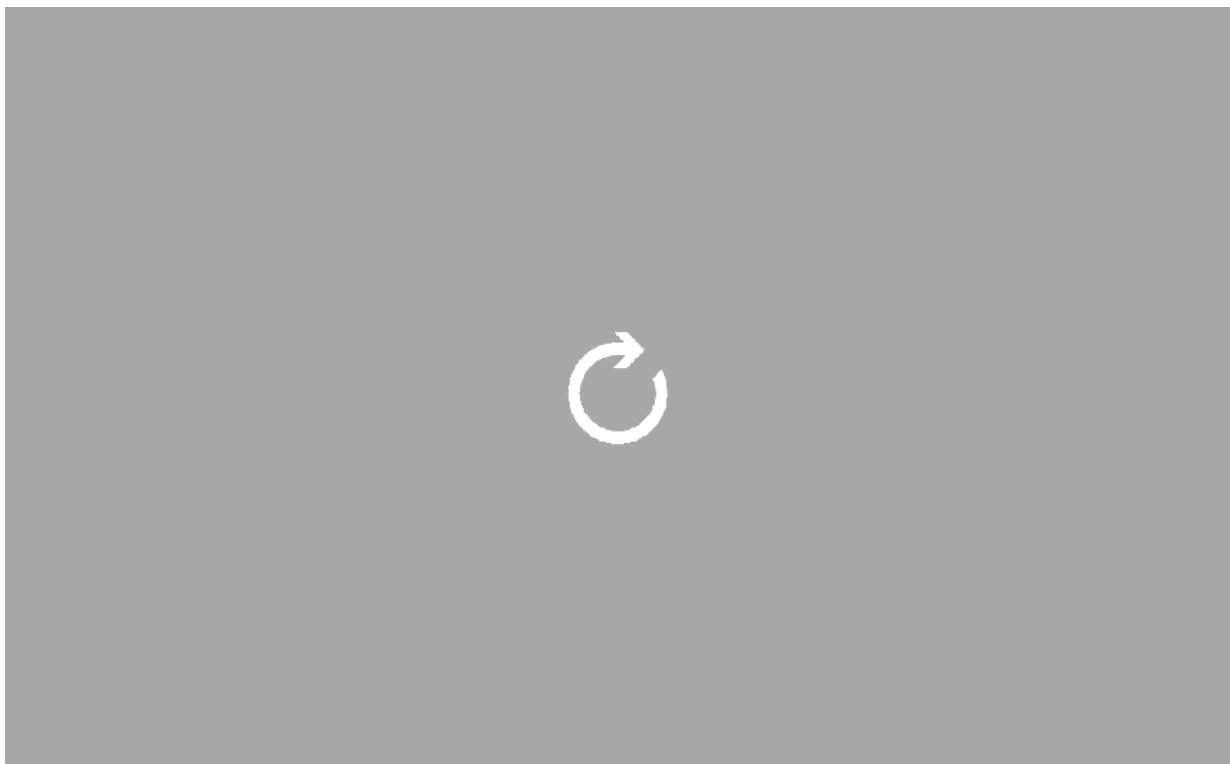


C2_W4_De
cision_Tr...

Video

Thursday, 23 November 2023 15:57

[index.mp4](#)



Losses functions for TensorFlow

Saturday, 30 September 2023 12:30

<https://keras.io/api/losses/>

Has the different loss functions for TensorFlow

Pandas

09 December 2024 23:55

https://pandas.pydata.org/Pandas_Cheat_Sheet.pdf