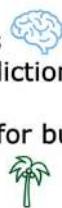


Advanced learning algorithms

Neural Networks
inference (prediction)
training



Practical advice for building machine learning systems



Decision Trees



决策树与神经网络相比。
©DeepLearning.AI decision trees compared to neural networks.

Andrew Ng

Neural networks

Origins: Algorithms that try to mimic the brain.

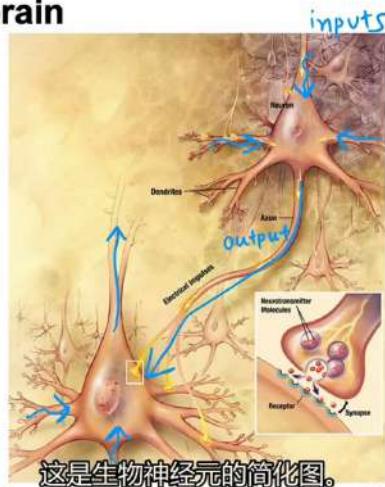


Used in the 1980's and early 1990's.
Fell out of favor in the late 1990's.

Resurgence from around 2005.

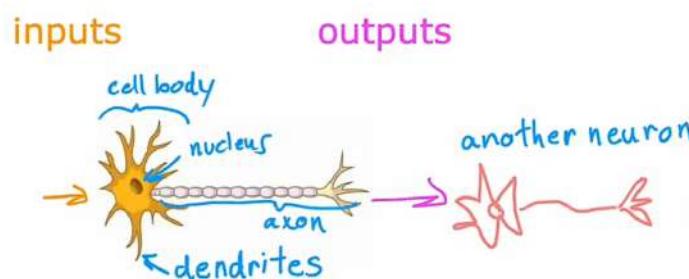
speech → images → text (NLP) → ...

Neurons in the brain

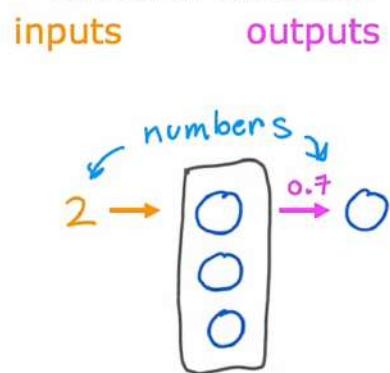


这是生物神经元的简化图。

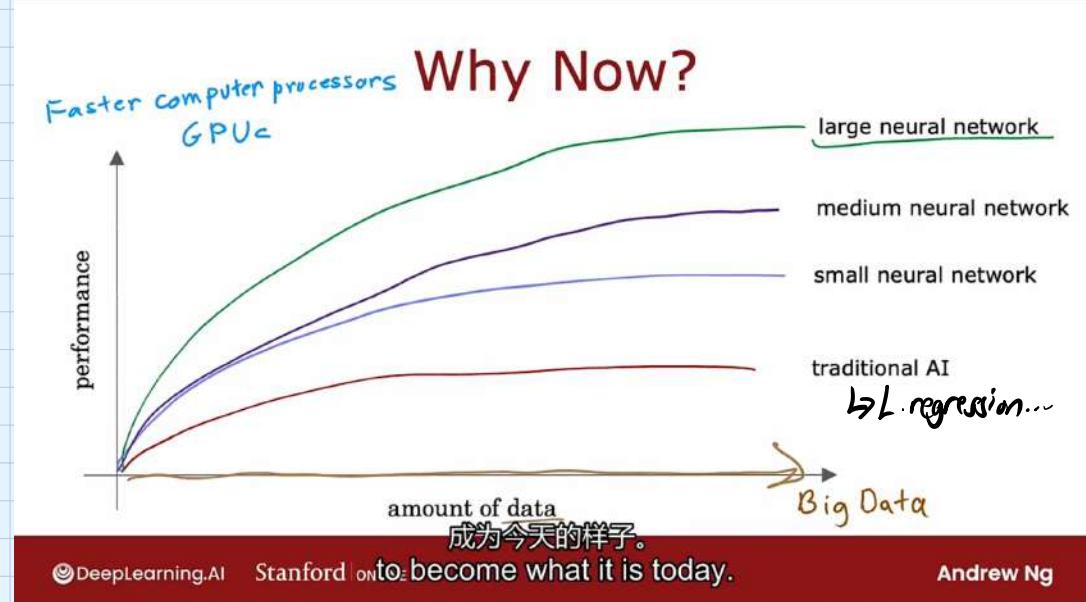
Biological neuron



Simplified mathematical model of a neuron



Why now ???

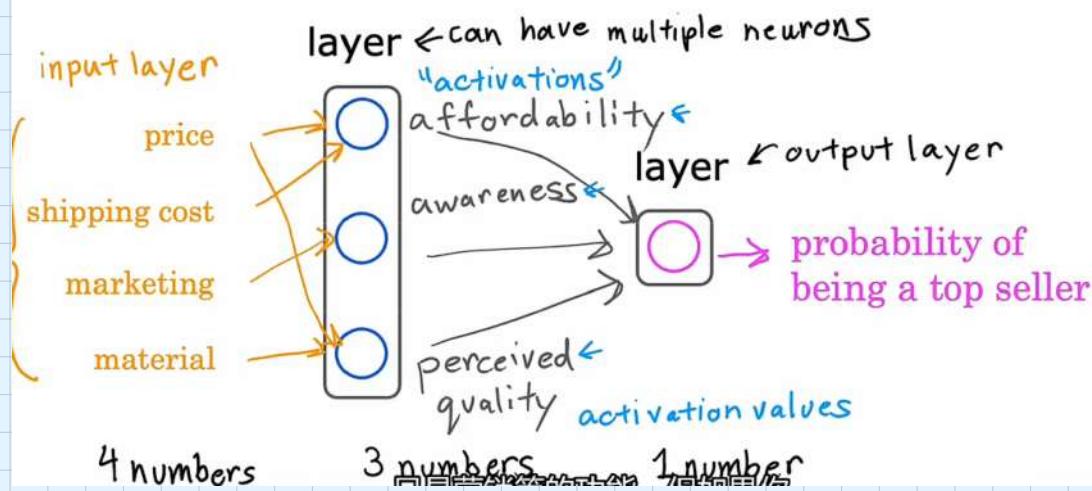
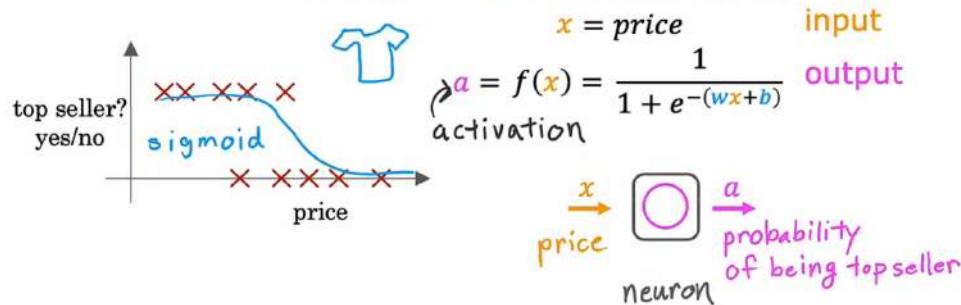


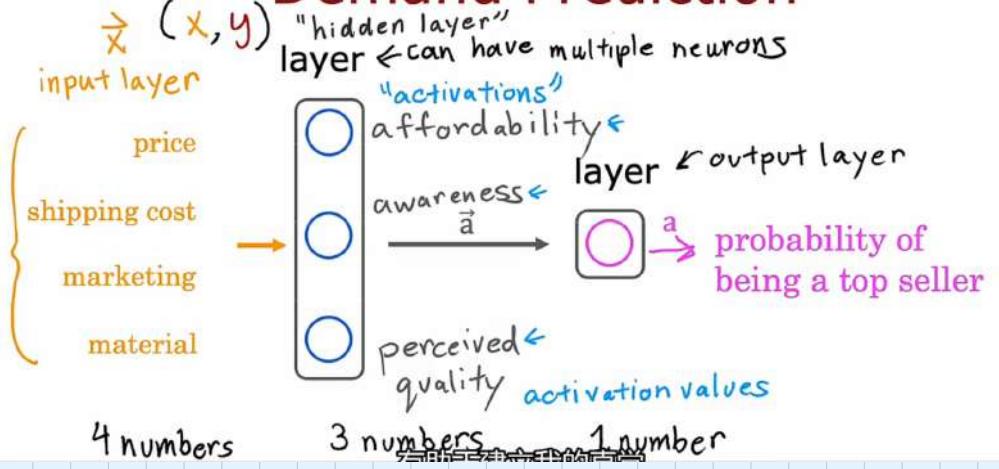
old AI can't implement

语音识别
图像处理
NLP ...

Demand Prediction

Demand Prediction





hidden layer make it strong!

这里的feature应该理解为特征而不是功能

Demand Prediction

feature engineering $x_1 x_2$

$\vec{x} (x, y)$ "hidden layer" layer ← can have multiple neurons

input layer

price
shipping cost
marketing
material

4 numbers → 3 numbers → 1 number

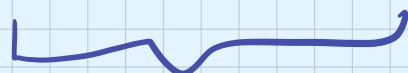
"activations"
affordability
awareness
perceived quality activation values

layer ← output layer

probability of being a top seller

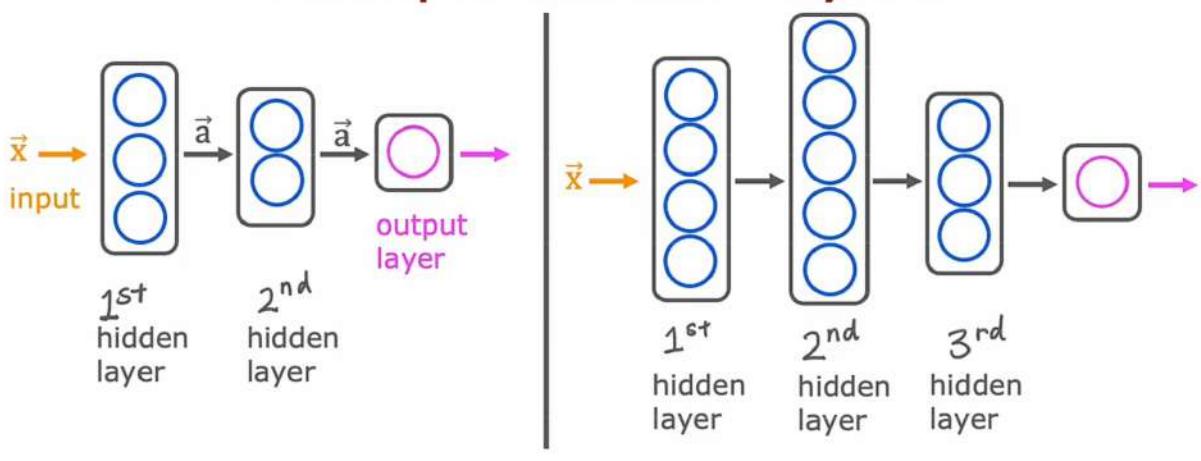
本身它想在这个隐藏层

DeepLearning.AI Stanford | it itself what are the features it Andrew Ng



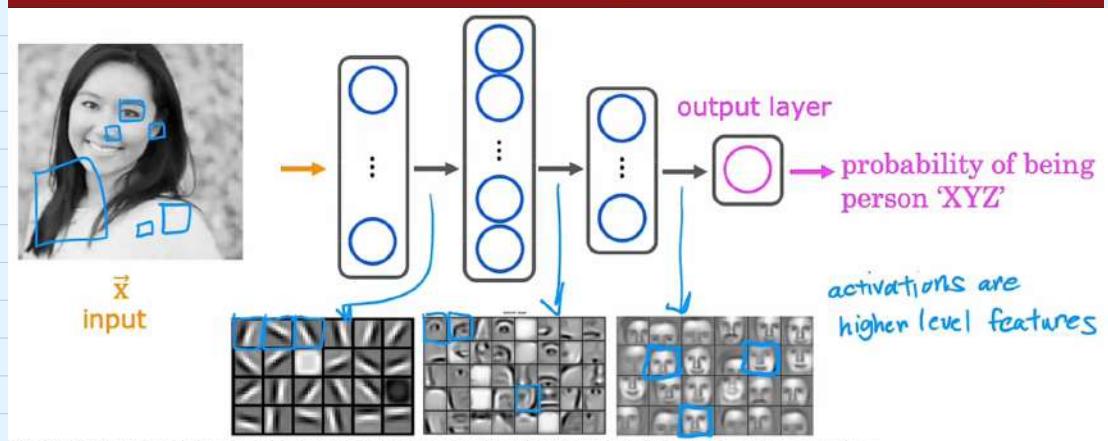
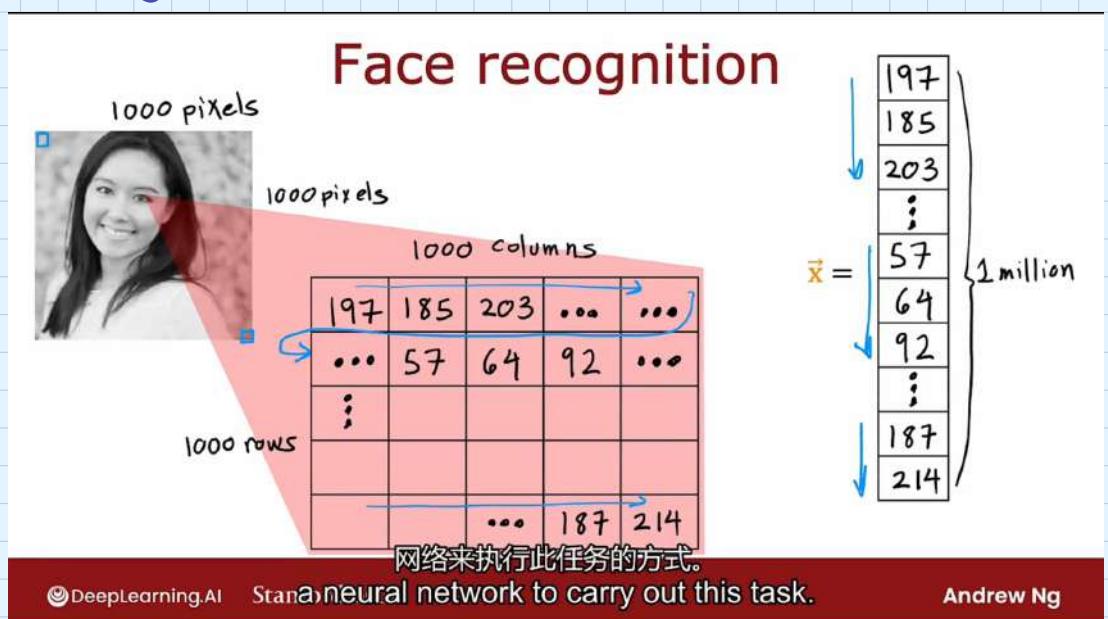
just like Logistic Regression

Multiple hidden layers



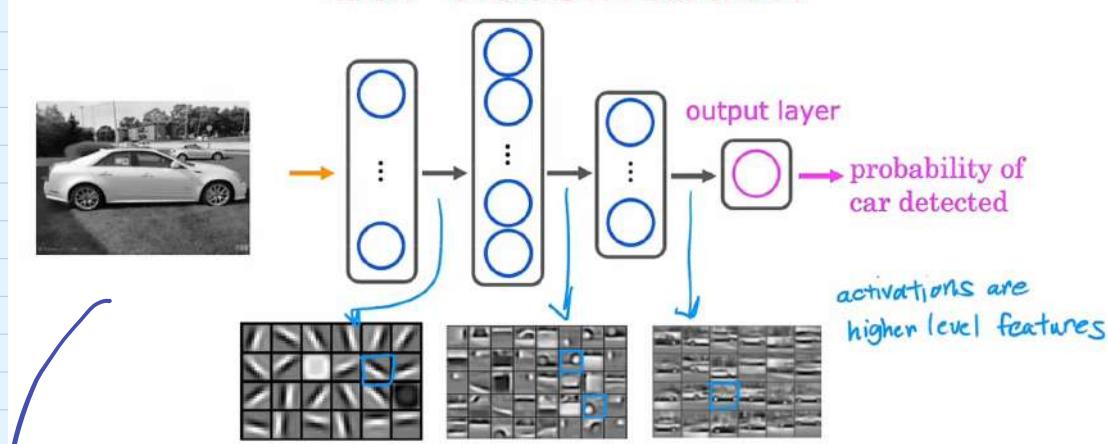
neural network architecture

Face recognition



source: Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations by Honglak Lee, Roger Grosse, Ranganath, Andrew Y Ng, 2008

Car classification



source: Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations by Honglak Lee, Roger Grosse, Rajen Ganesh, Andrew Y. Ng. [View paper](#)

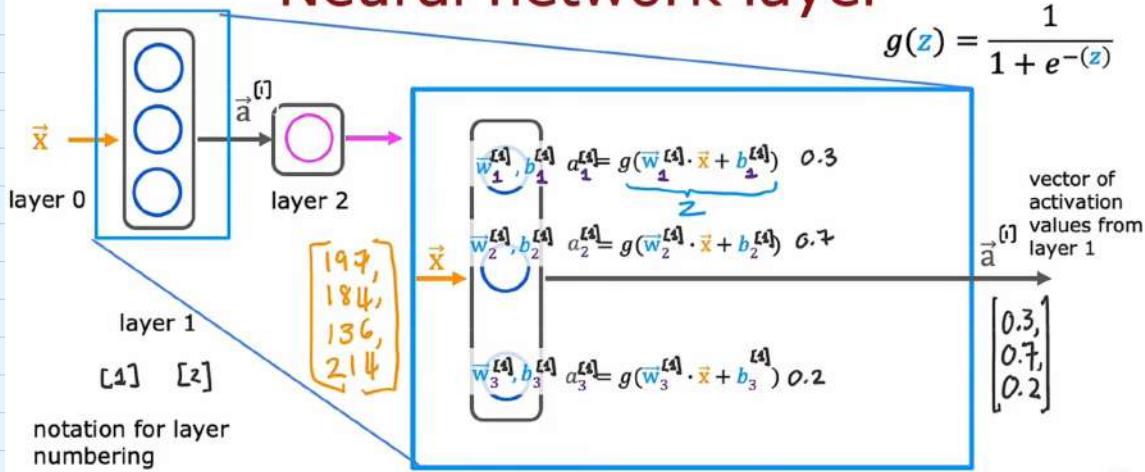
→ 不管丢什么进去，它会自动识别

Layer

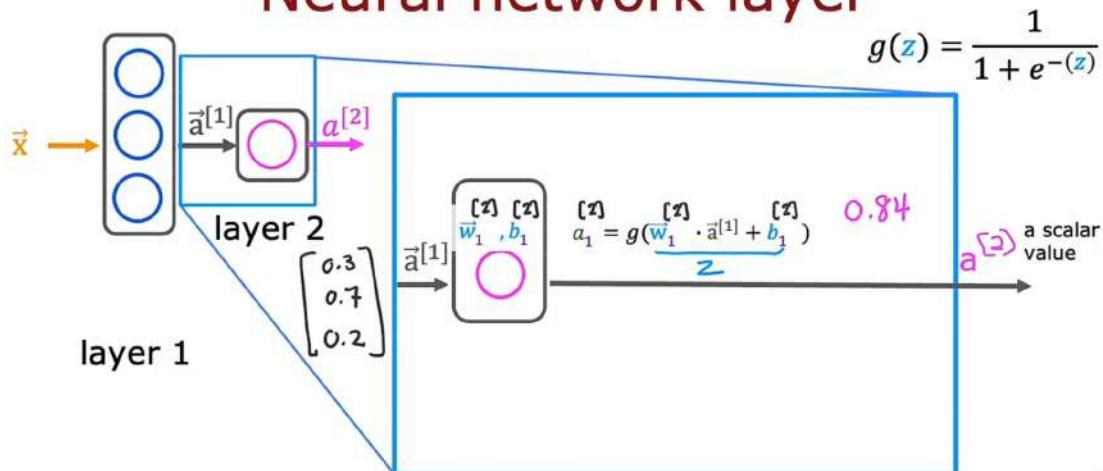
3 - 22 神经元向量表示

张量

Neural network layer



Neural network layer



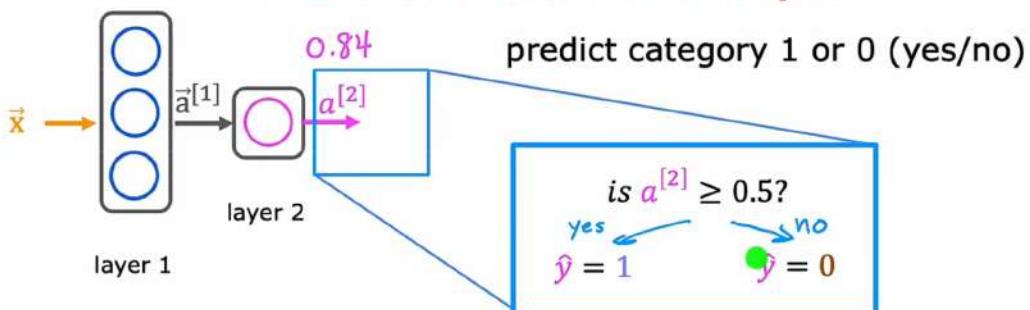
与神经网络的第 2 层相关联。

DeepLearning.AI

associated with layer 2 of the neural network.

Andrew Ng

Neural network layer



如果你愿意，这会给你最终的预测

DeepLearning.AI

If you wish, this then gives you

Andrew Ng

Complex neural network

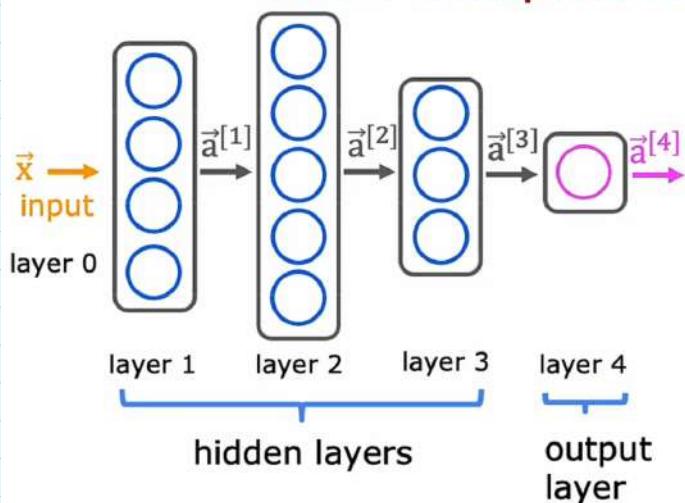
我大课签到

2024. 期初卡

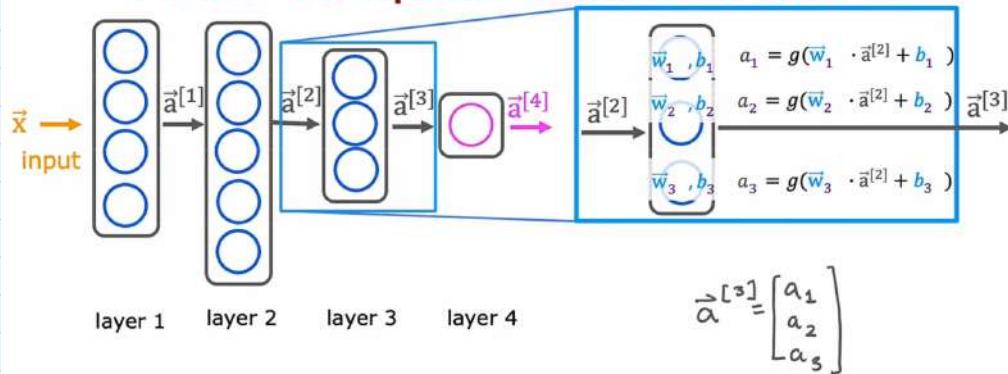
级梦研打卡

电子厂上班自学打卡

More complex neural network

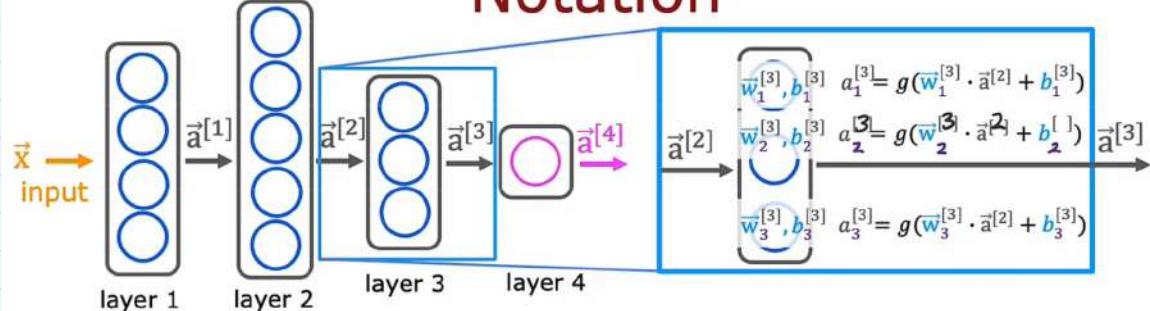


More complex neural network

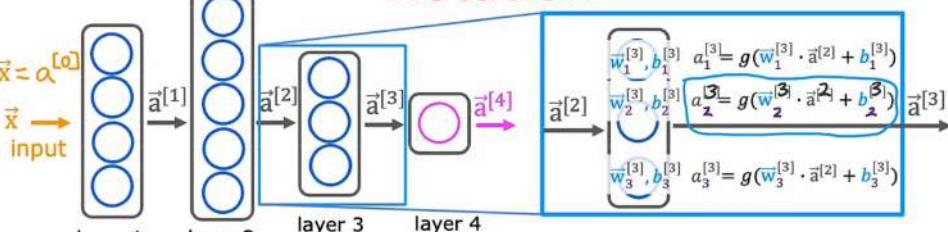


阶段一

Notation



INFORMATION



Activation value of
layer l , unit(neuron) j

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

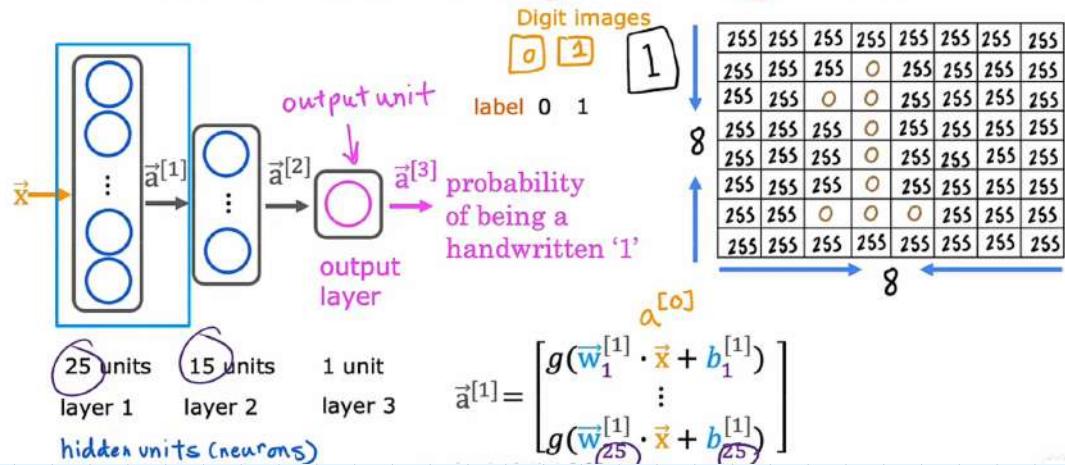
sigmoid
“activation function”

Parameters w & b of layer l , unit j

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

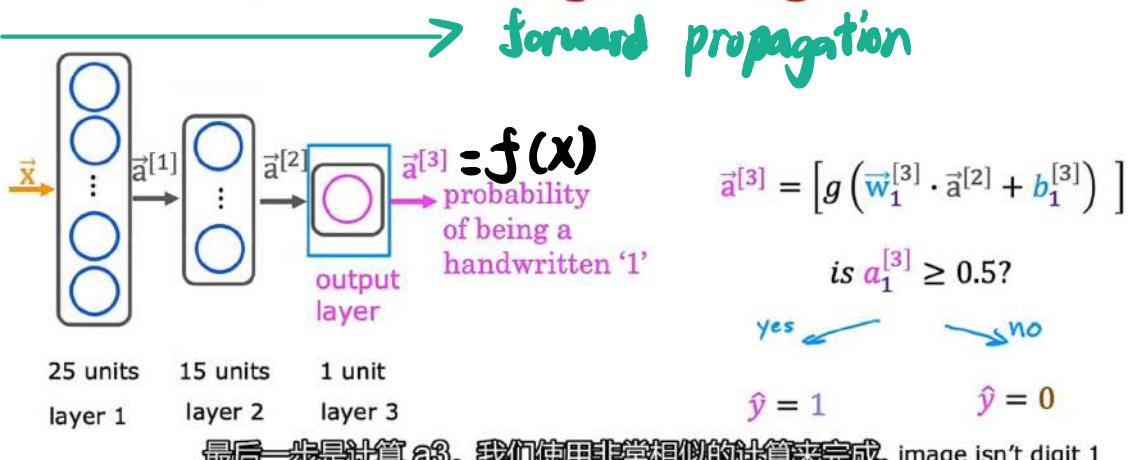
Handwritten digit recognition

Handwritten digit recognition



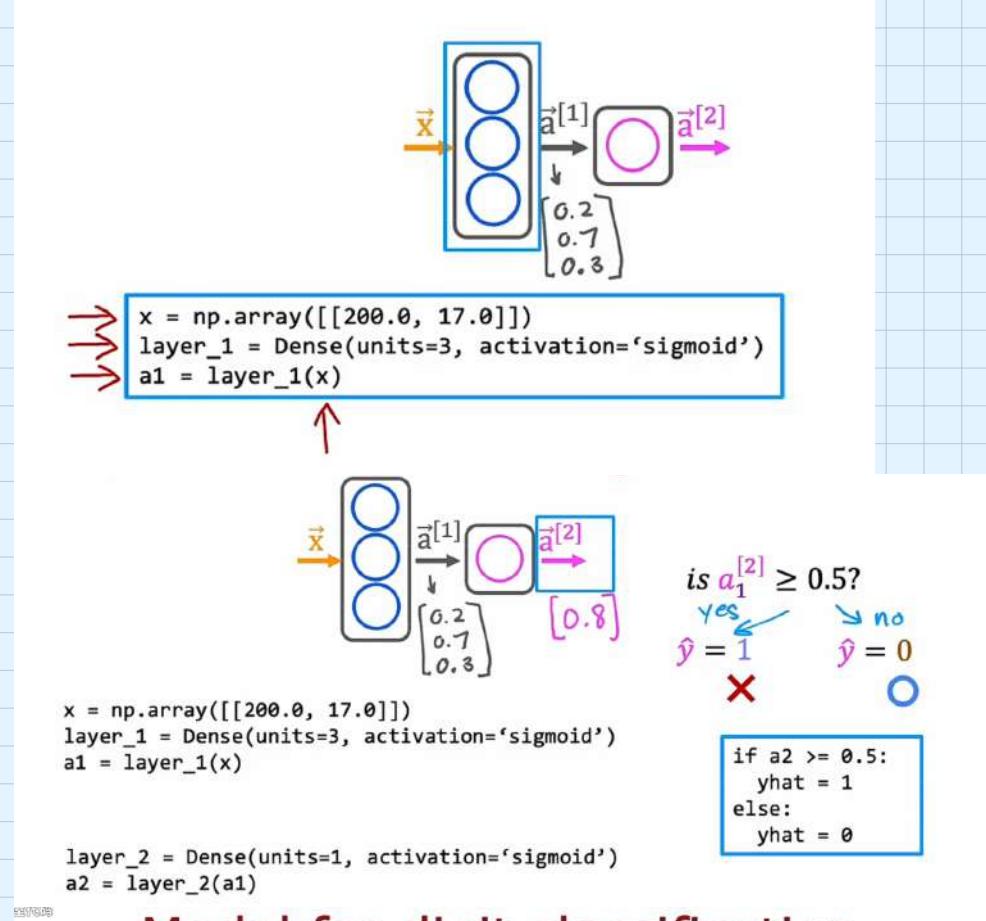
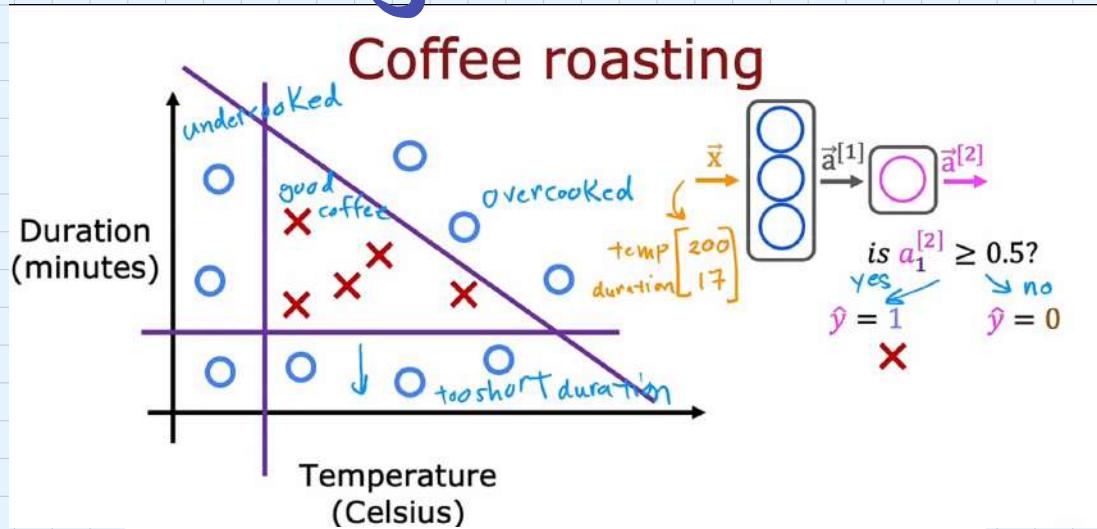
compute $\vec{a}^{[2]}$

Layer 2 $\Rightarrow \vec{a}^{[2]} = \begin{bmatrix} g(\vec{w}_1^{[2]} \cdot \vec{a}^{[1]} + b_1^{[2]}) \\ \vdots \\ g(\vec{w}_{15}^{[2]} \cdot \vec{a}^{[1]} + b_{15}^{[2]}) \end{bmatrix}$

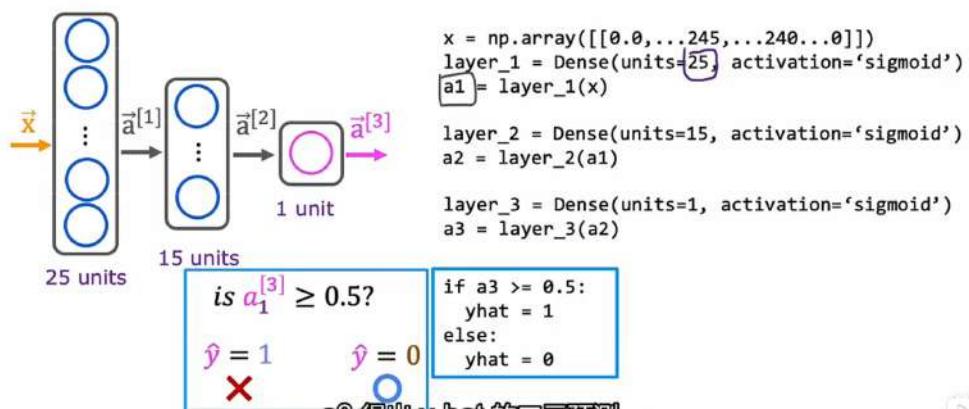


Tensorflow

~ Coffee roasting



Model for digit classification



Data in Tensorflow

Feature vectors

temperature (Celsius)	duration (minutes)	Good coffee? (1/0)
200.0	17.0	1
425.0	18.5	0
...

$x = np.array([[200.0, 17.0]]) \leftarrow$
[[200.0, 17.0]]

why? two?

Note about numpy arrays

3 columns
2 rows
2 x 3 matrix
 $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$

4 rows
2 columns
4 x 2 matrix
 $\begin{bmatrix} 0.1 & 0.2 \\ -3 & -4 \\ -0.5 & -0.6 \\ 7 & 8 \end{bmatrix}$

$x = np.array([[1, 2, 3], [4, 5, 6]])$
[[1, 2, 3],
[4, 5, 6]]

$x = np.array([[0.1, 0.2], [-3.0, -4.0], [-0.5, -0.6], [7.0, 8.0]])$
[[0.1, 0.2],
[-3.0, -4.0],
[-0.5, -0.6],
[7.0, 8.0]]]

2D array

2 x 3

4 x 2

1 x 2

2 x 1

$x = np.array([[200, 17]]) \rightarrow [200 \quad 17] \quad 1 \times 2$

$x = np.array([[200], [17]]) \rightarrow [200] \quad 2 \times 1$

$\rightarrow x = np.array([200, 17]) \quad \sim 1D$
"vector"

} matrix used
} in Tf

2D array

↓
More efficient

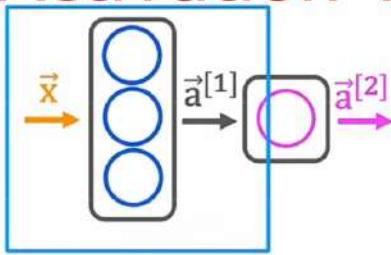
Feature vectors

temperature (Celsius)	duration (minutes)	Good coffee? (1/0)
200.0	17.0	1
425.0	18.5	0
...

$x = np.array([[200.0, 17.0]]) \leftarrow$
[[200.0, 17.0]]

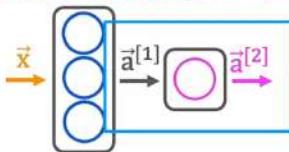
$\downarrow \quad \downarrow \quad 1 \times 2$
[200.0 17.0]

Activation vector



```
x = np.array([[200.0, 17.0]])
layer_1 = Dense(units=3, activation='sigmoid')
a1 = layer_1(x)
→ [[0.2, 0.7, 0.3]] 1 x 3 matrix
→ tf.Tensor([[[0.2 0.7 0.3]]], shape=(1, 3), dtype=float32)
→ a1.numpy()
array([[1.4661001, 1.125196 , 3.2159438]], dtype=float32)
```

Activation vector

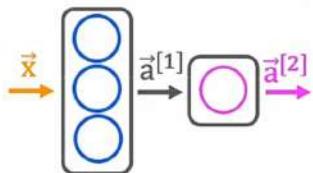


```
layer_2 = Dense(units=1, activation='sigmoid')
a2 = layer_2(a1)
→ [[0.8]] ←
→ tf.Tensor([[[0.8]]], shape=(1, 1), dtype=float32)
a2.numpy()
array([[0.8]], dtype=float32)
```

两种表示方式
Tensor } 对象
Numpy }

Building a neural network

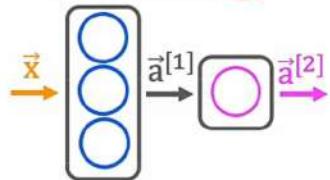
What you saw earlier



```
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)
```

Building a neural network architecture



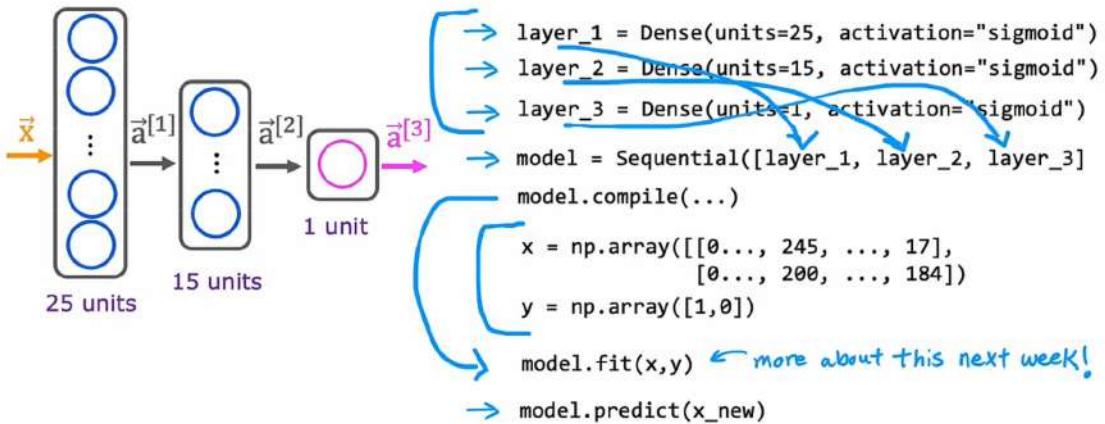
```
layer_1 = Dense(units=3, activation="sigmoid") ← layer1
layer_2 = Dense(units=1, activation="sigmoid") ← layer2
model = Sequential([layer_1, layer_2])
```

		y
200	17	1
120	5	0
425	20	0
212	18	1

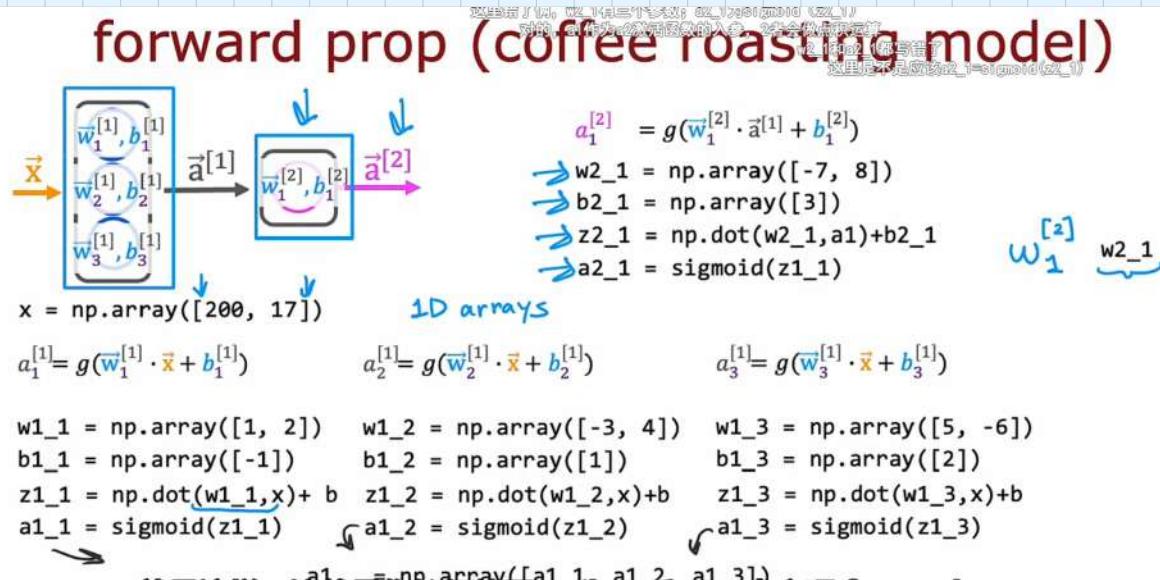
targets y = np.array([1,0,0,1])

```
model.compile(...)  
model.fit(x,y) ← more about this next week!  
model.predict(x_new) ←
```

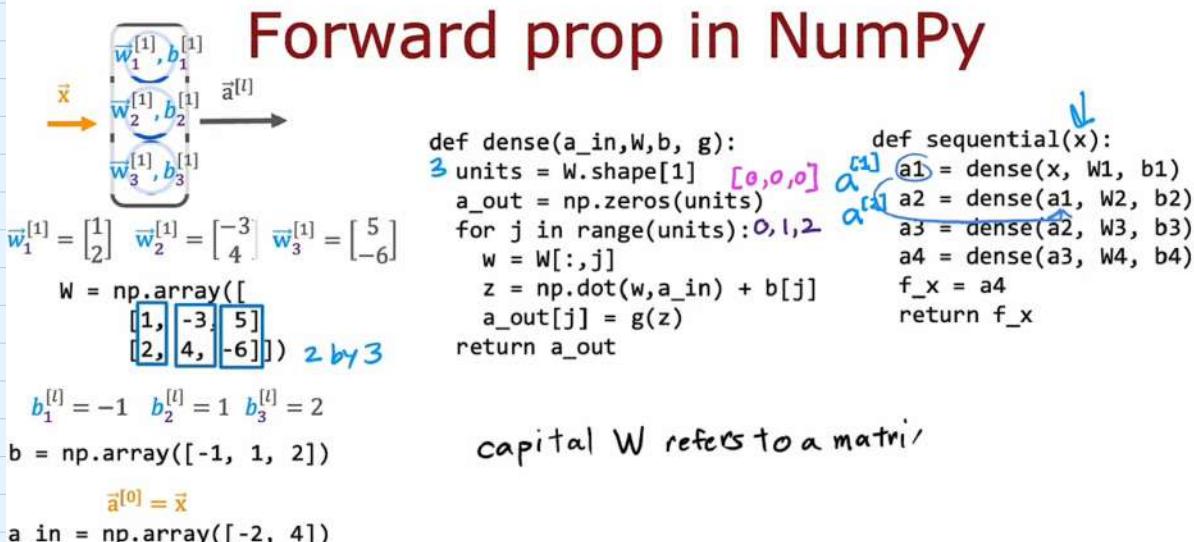
Digit classification model



Forward prop in a single layer (just using Python & Numpy)

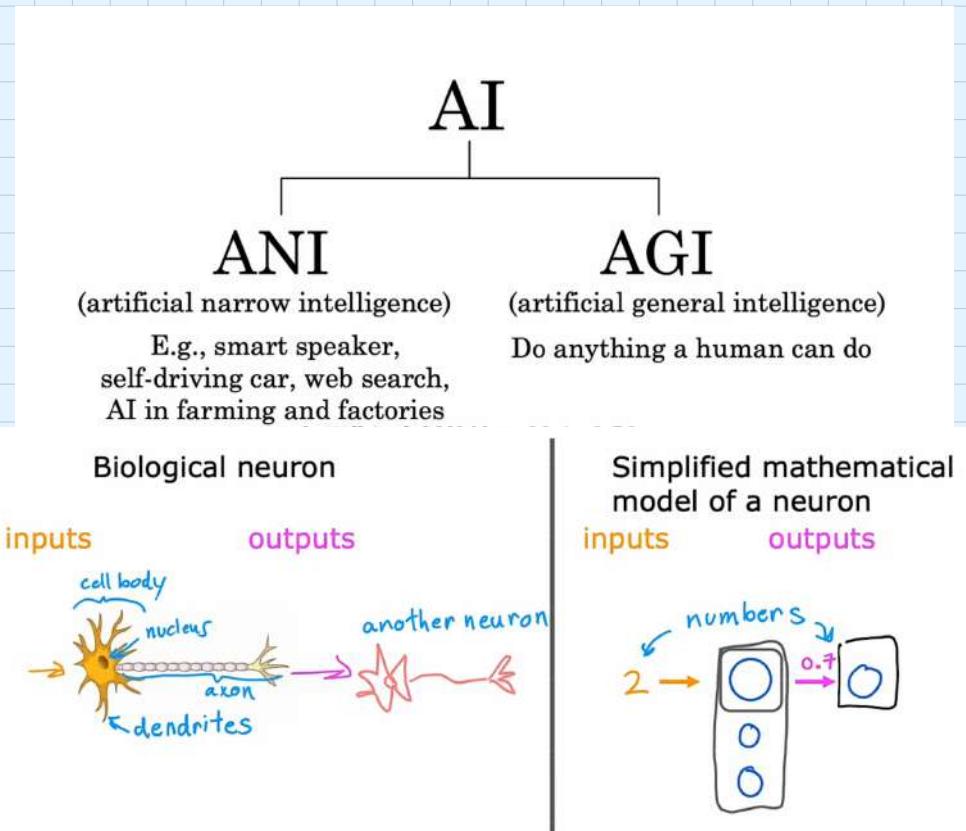


General implement of forward propagation



→ basic physical meanings behind Equation!

Speculations on AGI (artificial general intelligence)



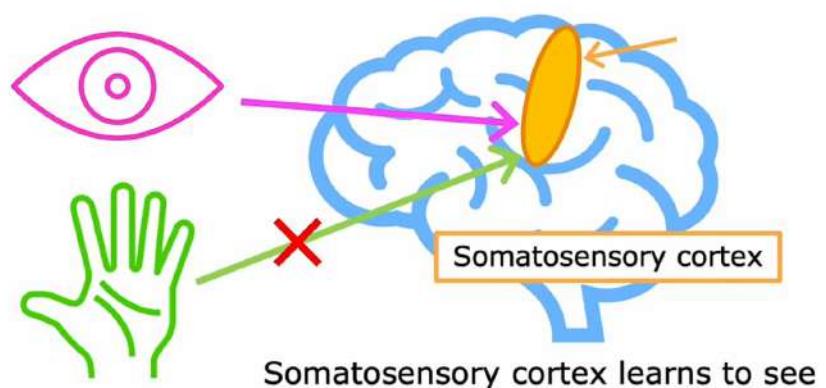
Neural network and the brain

Can we mimic the human brain?



We have (almost) no idea how the brain works

The "one learning algorithm" hypothesis



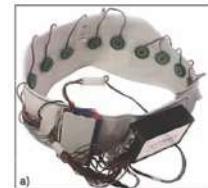
Sensor representations in the brain



Seeing with your tongue



Human echolocation (sonar)



Haptic belt: Direction sense



Implanting a 3rd eye

How neural networks are implemented efficiently
↳ caused Matrix multiplication

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):
    a_out = np.zeros(units)
    for j in range(units):
        w = W[:,j]
        z = np.dot(w,x) + b[j]
        a[j] = g(z)
    return a
```

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

def dense(A_in,W,B): all 2Darrays
    Z = np.matmul(A_in,W) + B
    A_out = g(Z) matrix multiplication
    return A_out
[[1,0,1]]
```

[1,0,1]

.....

matrix multiplication

Dot products

example $\begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 4 \end{bmatrix}$ $z = (1 \cdot 3) + (2 \cdot 4)$ $z = 3 + 8$ $z = 11$	in general $\begin{bmatrix} \uparrow \\ \vec{a} \end{bmatrix} \cdot \begin{bmatrix} \uparrow \\ \vec{w} \end{bmatrix}$ $z = \vec{a} \cdot \vec{w}$	transpose $\vec{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ $\vec{a}^T = [1 \ 2]$	vector vector multiplication $\begin{bmatrix} \leftarrow \vec{a}^T \rightarrow \\ 1 \times 2 \end{bmatrix} \begin{bmatrix} \uparrow \\ \vec{w} \end{bmatrix} 2 \times 1$ $z = \vec{a}^T \vec{w}$
---	---	---	---

equivalent

useful for understanding matrix multiplication

Vector matrix multiplication

$$\vec{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \vec{a}^T = [1 \ 2] \quad W = \begin{bmatrix} 3 & 5 \\ 4 & 6 \end{bmatrix} \quad Z = \vec{a}^T W \quad \left[\begin{array}{c|cc} \leftarrow \vec{a}^T & \rightarrow & \\ \hline & \vec{w}_1 & \vec{w}_2 \end{array} \right]$$

$$Z = [\vec{a}^T \vec{w}_1 \quad \vec{a}^T \vec{w}_2]$$

$$(1 * 3) + (2 * 4) \quad (1 * 5) + (2 * 6)$$

$$3 + 8 \quad 5 + 12$$

$$11 \quad 17$$

$$Z = [11 \ 17]$$

matrix matrix multiplication

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

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

rows columns

$$W = \begin{bmatrix} 3 & 5 \\ 4 & 6 \end{bmatrix}$$

$$Z = A^T W = \begin{bmatrix} \leftarrow & \overbrace{\vec{a}_1^T}^{\text{row 1}} & \rightarrow \\ \leftarrow & \overbrace{\vec{a}_2^T}^{\text{row 2}} & \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow & \vec{w}_1 & \uparrow \\ \downarrow & \vec{w}_2 & \downarrow \end{bmatrix}$$

$$\begin{array}{c} \text{row 1 col 1} \\ \text{row 2 col 1} \end{array} = \begin{bmatrix} \vec{a}_1^T \vec{w}_1 & \vec{a}_1^T \vec{w}_2 \\ \vec{a}_2^T \vec{w}_1 & \vec{a}_2^T \vec{w}_2 \end{bmatrix} \begin{array}{c} \text{row 1 col 2} \\ \text{row 2 col 2} \end{array}$$

$$\begin{array}{c} (-1 \times 3) + (-2 \times 4) \\ -3 + -8 \\ -11 \end{array} \quad \begin{array}{c} (-1 \times 5) + (-2 \times 6) \\ -5 + -12 \\ -17 \end{array}$$

$$= \begin{bmatrix} 11 & 17 \\ -11 & -17 \end{bmatrix}$$

general rules for
matrix multiplication
↳ next video!

Rules

Matrix multiplication rules

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

3×2 2×4

can only take dot products
of vectors that are same length

$$\begin{bmatrix} 0.1 & 0.2 \end{bmatrix} \quad \begin{bmatrix} 5 \\ 6 \end{bmatrix}$$

length 2 length 2

3 by 4 matrix
↳ same # rows as A^T
↳ same # columns as W

In Numpy

Matrix multiplication in NumPy

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

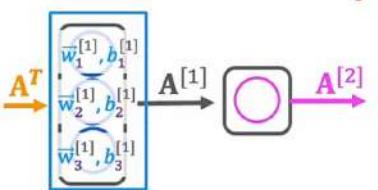
`A=np.array([1,-1,0.1], [2,-2,0.2]))` `W=np.array([3,5,7,9], [4,6,8,0]))` `Z = np.matmul(AT,W)`
`or`
`Z = AT @ W`

`AT=np.array([1,2], [-1,-2], [0.1,0.2])`

`AT=A.T` transpose

result
 $\begin{bmatrix} 11, 17, 23, 9 \\ -11, -17, -23, -9 \\ 1.1, 1.7, 2.3, 0.9 \end{bmatrix}$

Dense layer vectorized



$$\mathbf{A}^T = \begin{bmatrix} 200 & 17 \\ 1 & 2 \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} 1 & -3 & 5 \\ -2 & 4 & -6 \\ 2 & 3 & -6 \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} -1 & 1 & 2 \\ 1 & 2 & 3 \end{bmatrix}$$

$$\mathbf{Z} = \mathbf{A}^T \mathbf{W} + \mathbf{B}$$

$$\begin{bmatrix} 165 \\ z_1^{[1]} \end{bmatrix}, \begin{bmatrix} -531 \\ z_2^{[1]} \end{bmatrix}, \begin{bmatrix} 900 \\ z_3^{[1]} \end{bmatrix}$$

$$\mathbf{A} = g(\mathbf{Z})$$

$$\begin{bmatrix} 1 & 0 & 1 \end{bmatrix}$$

$\text{AT} = \text{np.array}([[200, 17]])$
 $\mathbf{W} = \text{np.array}([[1, -3, 5], [-2, 4, -6]])$
 $\mathbf{b} = \text{np.array}([[-1, 1, 2]])$

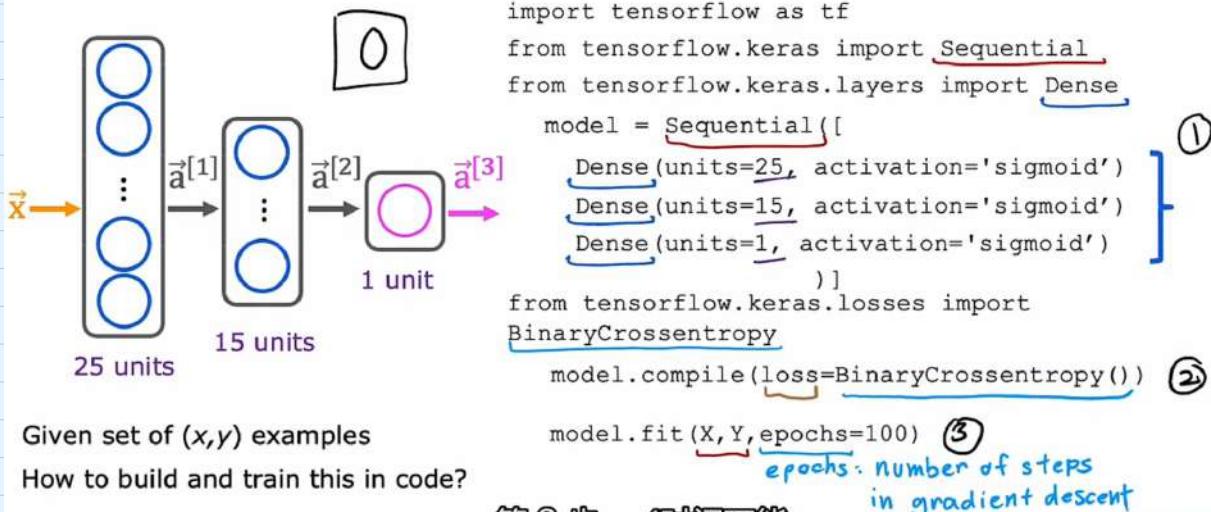
```
a_in
def dense(AT,W,b,g):
    z = np.matmul(AT,W) + b
    a_out = g(z)
    return a_out
[[1,0,1]]
```

Optional Lab

Second Week

TensorFlow implementation

Train a Neural Network in TensorFlow



Training detail

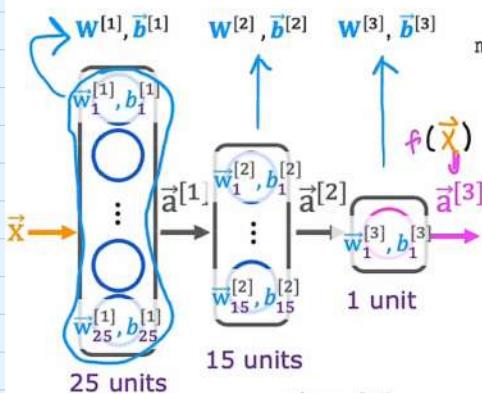
Model Training Steps

	Tensor Flow
① specify how to compute output given input x and parameters w, b (define model) $f_{\bar{w}, b}(\vec{x}) = ?$	logistic regression $z = \mathbf{np}.dot(w, x) + b$ $f_x = 1 / (1 + \mathbf{np}.exp(-z))$
② specify loss and cost $L(f_{\bar{w}, b}(\vec{x}), y)$ 1 example $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\bar{w}, b}(\vec{x}^{(i)}), y^{(i)})$	logistic loss $\text{loss} = -y * \mathbf{np}.log(f_x) - (1-y) * \mathbf{np}.log(1-f_x)$
③ Train on data to minimize $J(\vec{w}, b)$	neural network $\text{model} = \text{Sequential}([\text{Dense}(...), \text{Dense}(...), \text{Dense}(...)])$ binary cross entropy $\text{model.compile(loss=BinaryCrossentropy())}$ $\text{model.fit}(X, y, \text{epochs}=100)$

1. Create the model

define the model

$$f(\vec{x}) = ?$$



```

import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense

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

```

2. Loss and cost functions

Mnist digit classification problem

binary classification

$$L(f(\vec{x}), y) = -y \log(f(\vec{x})) - (1-y) \log(1-f(\vec{x}))$$

Compare prediction vs. target

logistic loss

also known as binary cross entropy

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)})$$

$\mathbf{W}^{[1]}, \mathbf{W}^{[2]}, \mathbf{W}^{[3]}$ $\mathbf{b}^{[1]}, \mathbf{b}^{[2]}, \mathbf{b}^{[3]}$

$$f_{\mathbf{W}, \mathbf{B}}(\vec{x})$$

model.compile(loss= BinaryCrossentropy())

regression

(predicting numbers and not categories) mean squared error

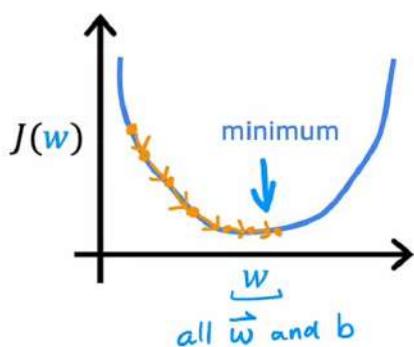
model.compile(loss= MeanSquaredError())

from tensorflow.keras.losses import
BinaryCrossentropy

Keras

from tensorflow.keras.losses import
MeanSquaredError

3. Gradient descent



repeat {

$$w_j^{[l]} = w_j^{[l]} - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b_j^{[l]} = b_j^{[l]} - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

} Compute derivatives
for gradient descent
using "back propagation"

model.fit(X, y, epochs=100)

nowadays, 成熟了, 因此可以

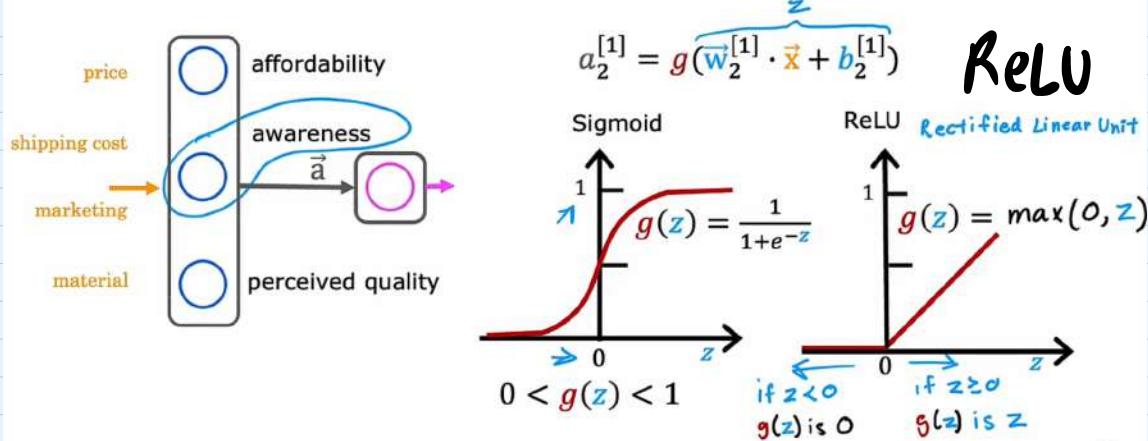
能夠替代

PyTorch → 主流

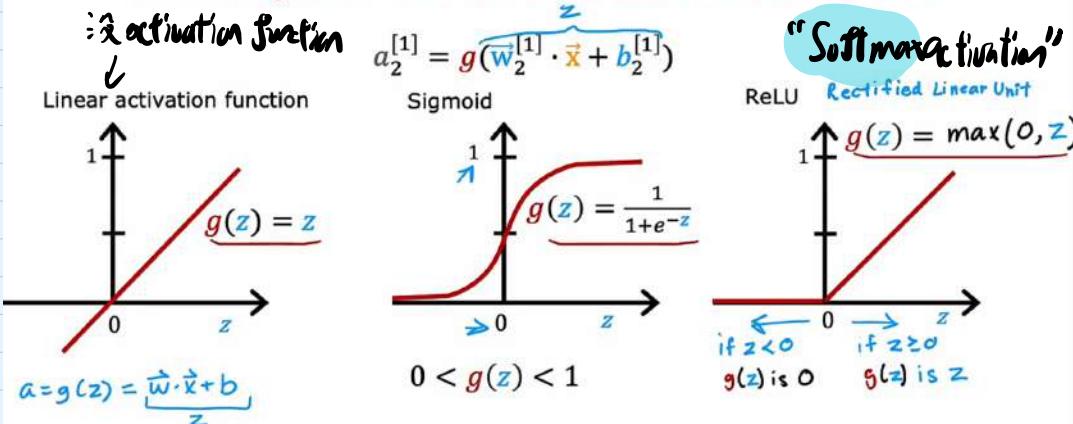
Alternative to the Sigmoid activation

Try probability

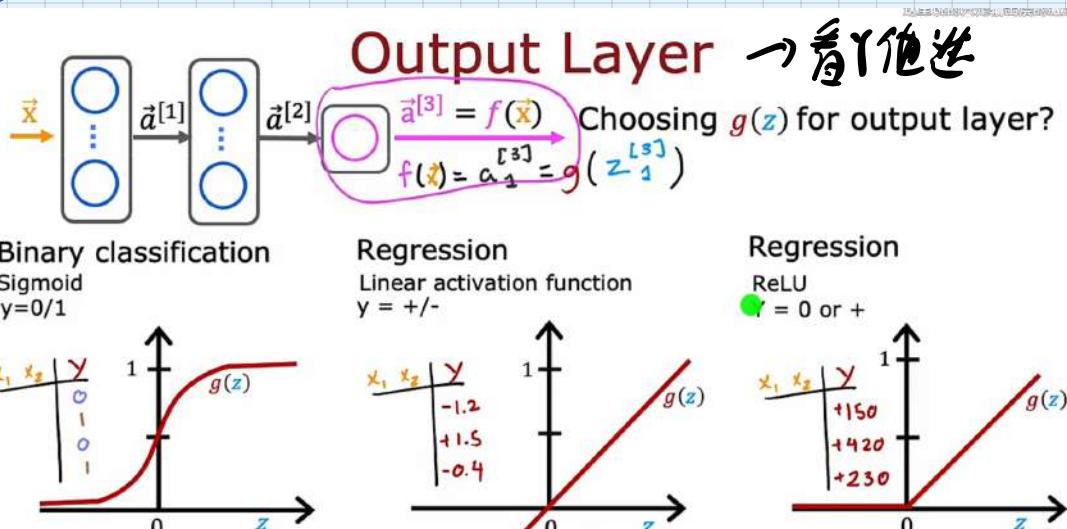
Demand Prediction Example

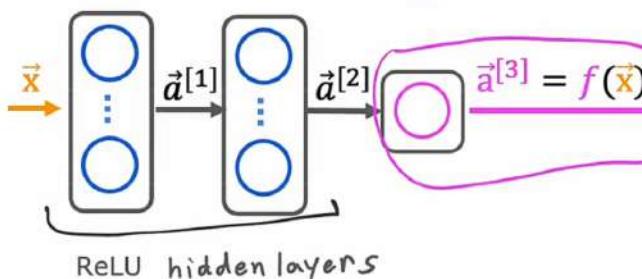
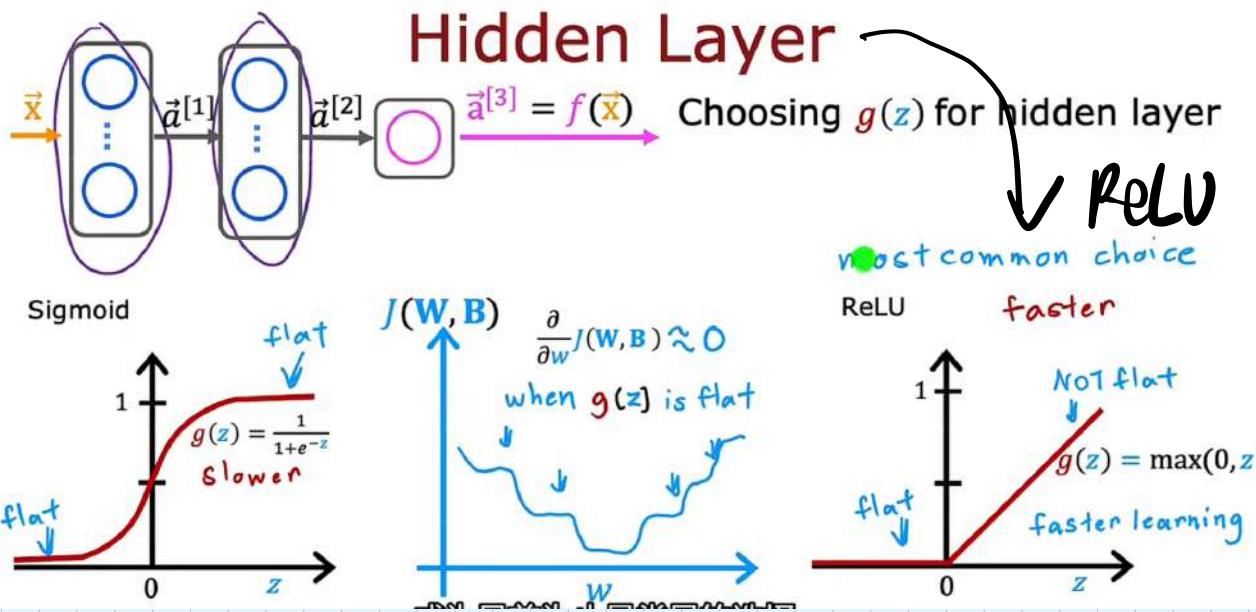


Examples of Activation Functions



Choosing activation functions





binary classification
activation='sigmoid'
regression y negative/
activation='linear'
regression $y \geq 0$
activation='relu'

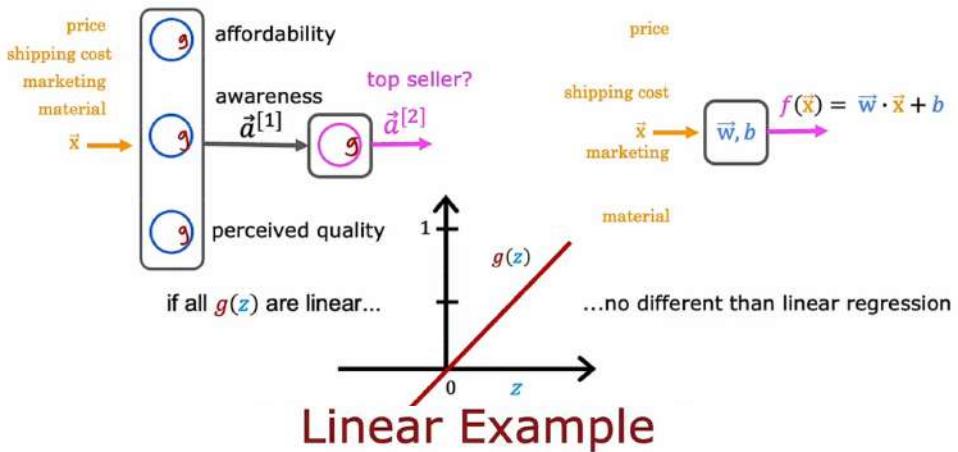
```
from tensorflow import keras
model = Sequential([
    Dense(units=25, activation='relu'), layer1
    Dense(units=15, activation='relu'), layer2
    Dense(units=1, activation='sigmoid') layer3
])
or 'linear'
or 'relu' or 'tanh' or 'softmax'
```

hidden \rightarrow ReLU (Soft Max)
output \rightarrow depends on y

如 tanh activation function or LeakyRelu or
Swish

Why do we need Activation function

Why do we need activation functions?



one feature
 x

w is scalar
 $a^{[1]}$
 $a^{[2]}$

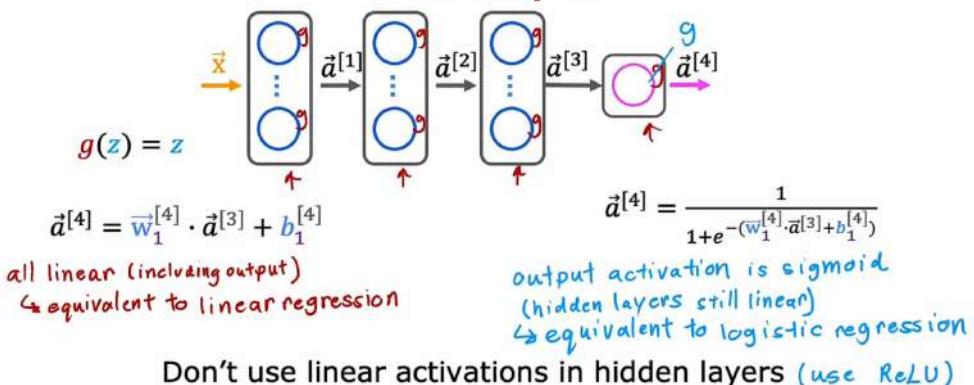
' a ' is scalar
 $g(z) = z$

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

$\vec{a}^{[2]} = \underline{w} \underline{x} + \underline{b}$
 $f(x) = wx + b$ linear regression

} 沒卵用

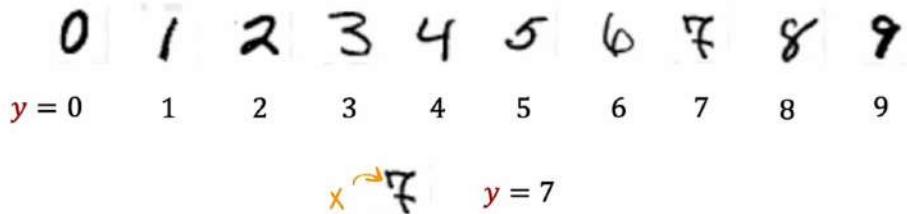
Example



ReLU 激活函数就可以了。

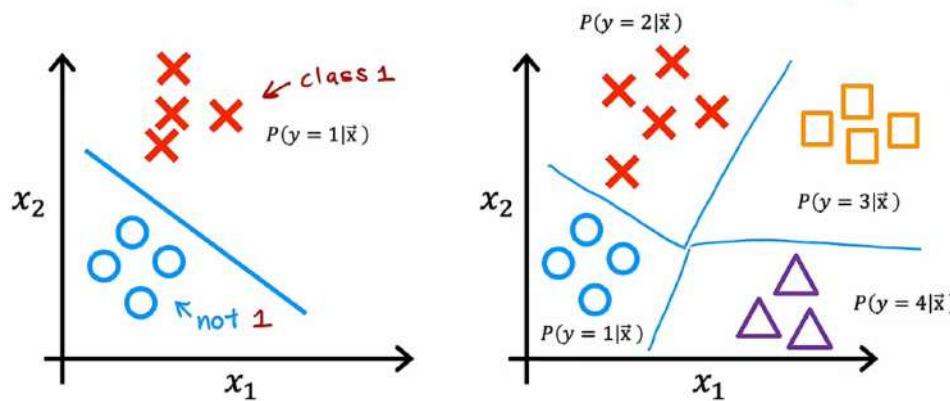
Multi Class

MNIST example



multiclass classification problem:
target y can take on more than two possible values

Multiclass classification example



Softmax

Logistic regression
(2 possible output values)

$$z = \vec{w} \cdot \vec{x} + b$$

$$\textcolor{red}{x} a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x}) \quad 0.11$$

$$\textcolor{blue}{o} a_2 = 1 - a_1 = P(y=0|\vec{x}) \quad 0.29$$

Softmax regression
(N possible outputs) $y=1, 2, 3, \dots, N$

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

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

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

$$\text{note: } a_1 + a_2 + \dots + a_N = 1$$

Softmax regression (4 possible outputs) $y=1, 2, 3, 4$

$$\textcolor{red}{x} z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \quad a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=1|\vec{x}) \quad 0.30$$

$$\textcolor{blue}{o} z_2 = \vec{w}_2 \cdot \vec{x} + b_2 \quad a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=2|\vec{x}) \quad 0.20$$

$$\square z_3 = \vec{w}_3 \cdot \vec{x} + b_3 \quad a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=3|\vec{x}) \quad 0.15$$

$$\triangle z_4 = \vec{w}_4 \cdot \vec{x} + b_4 \quad a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=4|\vec{x}) \quad 0.35$$

Cost

Logistic regression

$$z = \vec{w} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x})$$

$$a_2 = 1 - a_1 = P(y=0|\vec{x})$$

$$\text{loss} = -y \log a_1 - (1-y) \log(1-a_1)$$

$\text{if } y=1$ $\text{if } y=0$

$$J(\vec{w}, b) = \text{average loss}$$

Softmax regression

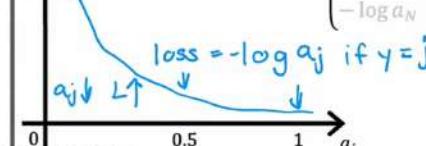
$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y=1|\vec{x})$$

$$\vdots$$

$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y=N|\vec{x})$$

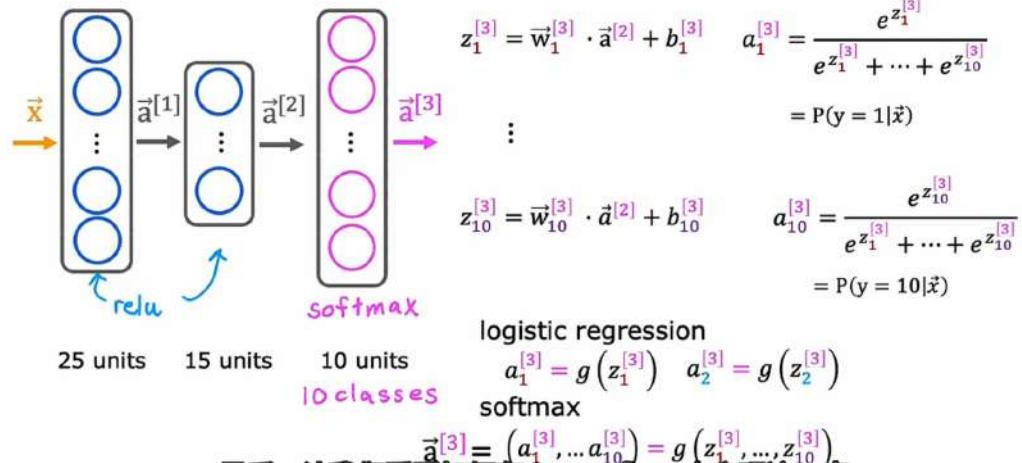
Crossentropy loss

$$\text{loss}(a_1, \dots, a_N, y) = \begin{cases} -\log a_1 & \text{if } y=1 \\ -\log a_2 & \text{if } y=2 \\ \vdots \\ -\log a_N & \text{if } y=N \end{cases}$$



Neural Network with Softmax output

Neural Network with Softmax output



① specify the model

$$f_{\vec{w}, b}(\vec{x}) = ?$$

② specify loss and cost

$$L(f_{\vec{w}, b}(\vec{x}), \vec{y})$$

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

MNIST with softmax

```
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')
])
from tensorflow.keras.losses import
SparseCategoricalCrossentropy
model.compile(loss= SparseCategoricalCrossentropy() )
model.fit(X, Y, epochs=100)
Note: better (recommended) version later.
Don't use the version shown here!
```

Improved implementation of softmax (better way)

Numerical Roundoff Errors

option 1

$$x = \frac{2}{10,000}$$

option 2

$$x = \left(1 + \frac{1}{10,000}\right) - \left(1 - \frac{1}{10,000}\right) =$$

```
jupyter Numerical roundoff errors (unSaved changes)
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Code Validate
```

```
In [1]: x1 = 2.0 / 10000
print(f"{x1:.18f}") # print 18 digits to the right of decimal point
0.00020000000000000000

In [2]: x2 = 1 + (1/10000) - (1 - 1/10000)
print(f"{x2: .18f}")
0.00019999999999978
```

Numerical Roundoff Errors

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

Logistic regression:

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

Original loss

$$loss = -y \log(\hat{a}) - (1-y) \log(1-\hat{a})$$

model = Sequential([

Dense(units=25, activation='relu')

Dense(units=15, activation='relu') 'linear'

Dense(units=10, activation='sigmoid')

model.compile(loss=BinaryCrossEntropy())

More accurate loss (in code)

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

logit: z

In softmax (10類)

More numerically accurate implementation of softmax

Softmax regression

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

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

model = Sequential([

Dense(units=25, activation='relu')

Dense(units=15, activation='relu')

Dense(units=10, activation='softmax')

'linear'

→ become linear

方法
方法

More Accurate

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

model.compile(loss=SparseCategoricalCrossentropy())

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 a, ..., a10
f_x = tf.nn.softmax(logits) ← is z1, ..., z10
```

↑
2

Classification with multiple outputs 标示多 label

Multi-label Classification



Is there a car?

$$\begin{matrix} \text{yes} \\ \text{no} \end{matrix} \quad \vec{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Is there a bus?

$$\begin{matrix} \text{no} \\ \text{no} \\ \text{yes} \end{matrix} \quad \vec{y} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

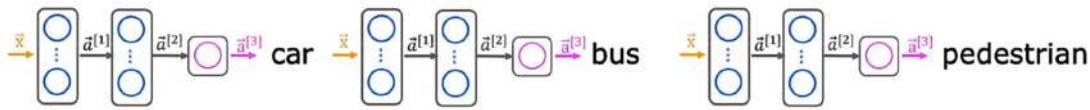
Is there a pedestrian?

$$\begin{matrix} \text{yes} \\ \text{yes} \\ \text{no} \end{matrix} \quad \vec{y} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

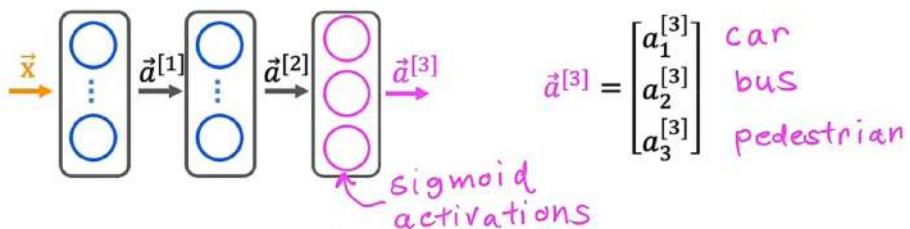
三个二分类问题

two ways \rightarrow 1) build three neuron network
2) train simultaneously

Multiple classes



Alternatively, train one neural network with three outputs



Multi-label vs multi-class

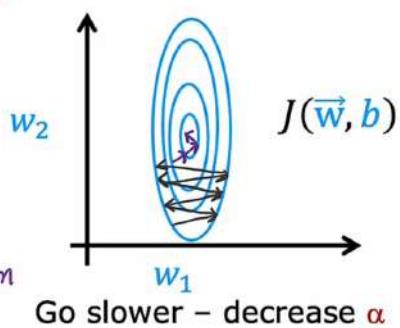
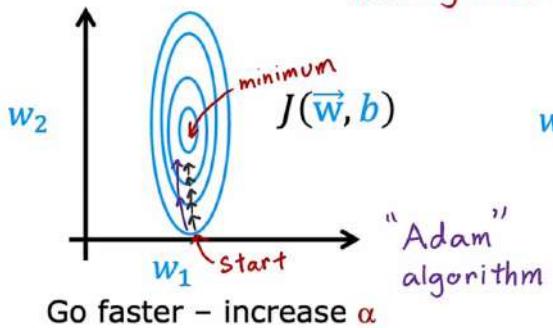
Advanced Optimization 高級最適化

→ train NN much faster than GD

Gradient Descent

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

learning rate



Adam Algorithm Intuition

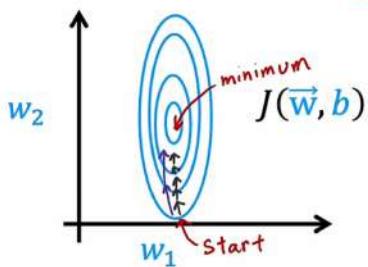
Adam: Adaptive Moment estimation not just one α

$$w_1 = w_1 - \underbrace{\alpha_1}_{\vdots} \frac{\partial}{\partial w_1} J(\vec{w}, b)$$

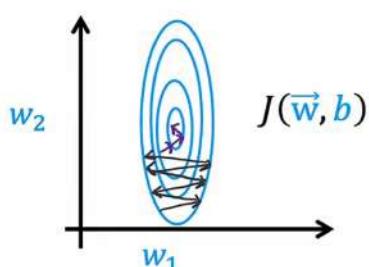
$$w_{10} = w_{10} - \underbrace{\alpha_{10}}_{\vdots} \frac{\partial}{\partial w_{10}} J(\vec{w}, b)$$

$$b = b - \underbrace{\alpha_{11}}_{\vdots} \frac{\partial}{\partial b} J(\vec{w}, b)$$

Adam Algorithm Intuition



If w_j (or b) keeps moving in same direction, increase α_j .



If w_j (or b) keeps oscillating, reduce α_j .

~: 学习率

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))
```

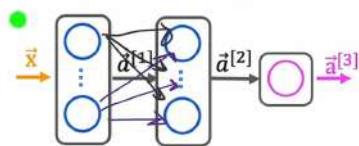
$\alpha = 10^{-3} = 0.001$

fit

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

Additional Layer Types

Dense Layer



Each neuron output is a function of all the activation outputs of the previous layer.

$$\vec{a}_1^{[2]} = g(\vec{w}_1^{[2]} \cdot \vec{a}^{[1]} + b_1^{[2]})$$

Convolutional Layer



Each Neuron only looks at part of the previous layer's inputs.

Why?

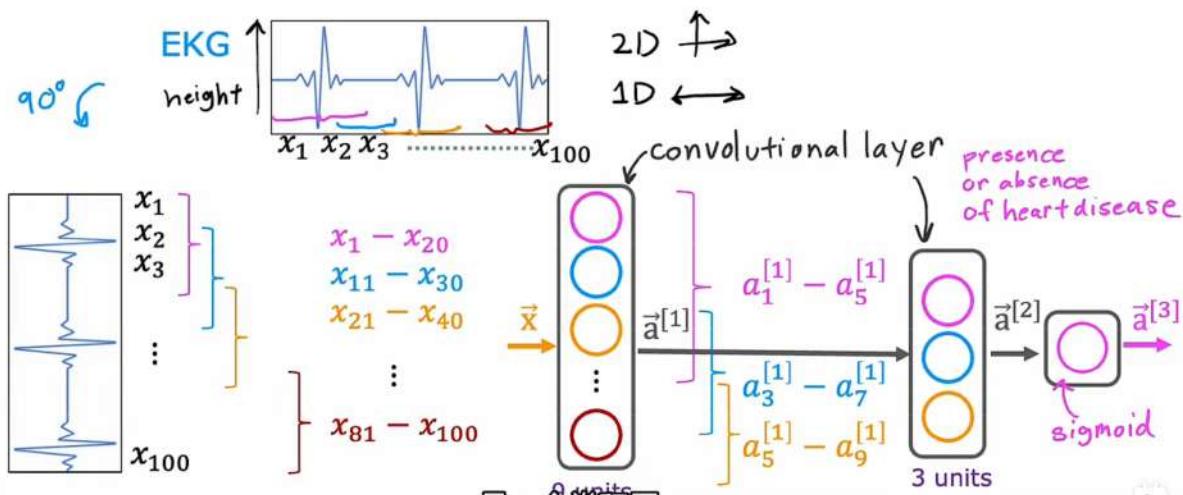
- Faster computation
- Need less training data (less prone to overfitting)

~:深度学习基础

→ 卷积层

→ 深度学习

Convolutional Neural Network



↳ 1D

Learn in Deep Learning

Optional Lab ~ see Jupyter

可选 \rightsquigarrow What is a derivative 5.1

可选 \rightsquigarrow Computation Graph 5.2

可选 \rightsquigarrow Larger Neural Network Example 5.3

\rightarrow 这里快速带过

Third week

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 ←
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features ($x_1^2, x_2^2, x_1x_2, \text{etc}$)
- Try decreasing λ
- Try increasing λ ↑↑

isn't

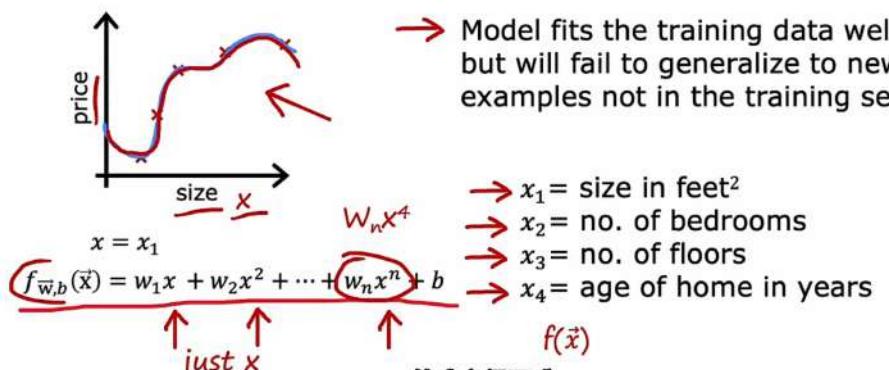
Machine learning diagnostic

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

Evaluating your model



Evaluating your model

Dataset:

	size	price			
70%	2104	400	{}	$(x^{(1)}, y^{(1)})$	$m_{train} =$
	1600	330		$(x^{(2)}, y^{(2)})$	no. training examples
	2400	369		\vdots	= 7
	1416	232		$(x^{(m_{train})}, y^{m_{train}})$	
	3000	540			
	1985	300			
30%	1534	315	{}	$(x^{(1)}, y^{(1)})$	$m_{test} =$
	1427	199		\vdots	no. test examples
	1380	212		$(x^{(m_{test})}, y^{m_{test}})$	= 3
	1494	243			

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

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

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

Compute test error:

$$J_{test}(\vec{w}, b) = \frac{1}{2m_{test}} \left[\sum_{i=1}^{m_{test}} (f_{\vec{w}, b}(\vec{x}_{test}^{(i)}) - y_{test}^{(i)})^2 \right] + \cancel{+ J_{train}(\vec{w}, b)}$$

Compute training error:

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

Train/test procedure for classification problem

Fit parameters by minimizing $J(\vec{w}, b)$ to find \vec{w}, b

E.g.,

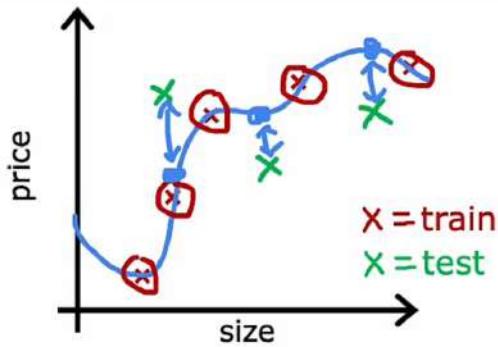
$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

Compute test error:

$$J_{test}(\vec{w}, b) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} \left[y_{test}^{(i)} \log(f_{\vec{w}, b}(\vec{x}_{test}^{(i)})) + (1 - y_{test}^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}_{test}^{(i)})) \right]$$

Compute train error:

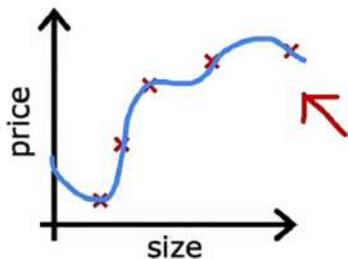
$$J_{train}(\vec{w}, b) = -\frac{1}{m_{train}} \sum_{i=1}^{m_{train}} \left[y_{train}^{(i)} \log(f_{\vec{w}, b}(\vec{x}_{train}^{(i)})) + (1 - y_{train}^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}_{train}^{(i)})) \right]$$



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

Model Selection / Cross-validation

Model selection (choosing a model)



$$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

Once parameters \vec{w}, b are fit to the training set, the training error $J_{train}(\vec{w}, b)$ is likely lower than the actual generalization error.

$J_{test}(\vec{w}, b)$ is better estimate of how well the model will generalize to new data than $J_{train}(\vec{w}, b)$.

Model selection (choosing a model)

- $d=1$ 1. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + b \rightarrow \vec{w}, b \rightarrow J_{test}(w^{<1>}, b^{<1>})$
- $d=2$ 2. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + b \rightarrow \vec{w}, b \rightarrow J_{test}(w^{<2>}, b^{<2>})$
- $d=3$ 3. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + w_3 x^3 + b \rightarrow \vec{w}^{<3>}, b^{<3>} \rightarrow J_{test}(w^{<3>}, b^{<3>})$
- \vdots
- $d=10$ 10. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + \dots + w_{10} x^{10} + b \rightarrow J_{test}(w^{<10>}, b^{<10>})$

Choose $w_1 x_1 + \dots + w_5 x^5 + b$ $d=5 \quad J_{test}(w^{<5>}, b^{<5>})$

How well does the model perform? Report test set error $J_{test}(w^{<5>}, b^{<5>})$?

The problem is $J_{test}(w^{<5>}, b^{<5>})$ is likely to be an optimistic estimate of generalization error. Ie: An extra parameter d (degree of polynomial) was chosen using the test set.

w, b

Training/cross validation/test set

size	price	validation set	development set
2104	400		
1600	330		
2400	369		
1416	232		
3000	540		
1985	300		
1534	315	$(x^{(1)}, y^{(1)})$	$m_{train} = 6$
1427	199	\vdots	
1380	212	$(x^{(m_{train})}, y^{(m_{train})})$	
1494	243	$(x_{cv}^{(1)}, y_{cv}^{(1)})$	$m_{cv} = 2$
		\vdots	
		$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$	
		$(x_{test}^{(1)}, y_{test}^{(1)})$	
		\vdots	
		$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$	$m_{test} = 2$

最左边的备注：日本交叉验证很便宜

Training/cross validation/test set

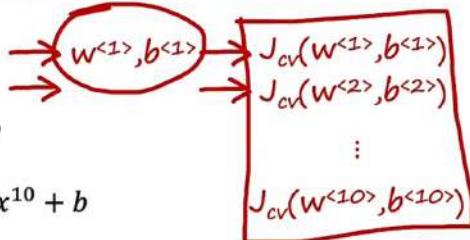
Training error: $J_{train}(\vec{w}, b) = \frac{1}{2m_{train}} \left[\sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 \right]$

Cross validation error: $J_{cv}(\vec{w}, b) = \frac{1}{2m_{cv}} \left[\sum_{i=1}^{m_{cv}} (f_{\vec{w}, b}(\vec{x}_{cv}^{(i)}) - y_{cv}^{(i)})^2 \right] \quad (\text{validation error, dev error})$

Test error: $J_{test}(\vec{w}, b) = \frac{1}{2m_{test}} \left[\sum_{i=1}^{m_{test}} (f_{\vec{w}, b}(\vec{x}_{test}^{(i)}) - y_{test}^{(i)})^2 \right]$

Model selection

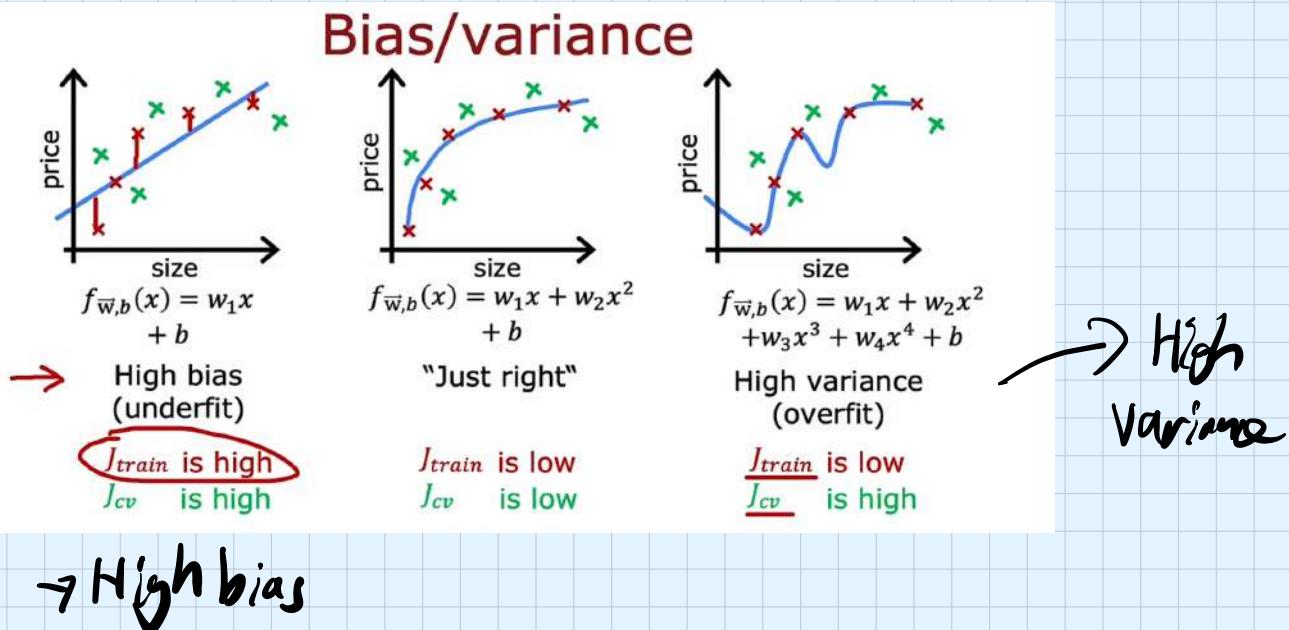
- $d=1$ 1. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + b$
- $d=2$ 2. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + b$
- $d=3$ 3. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + w_3 x^3 + b$
- \vdots
- $d=10$ 10. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + \dots + w_{10} x^{10} + b$



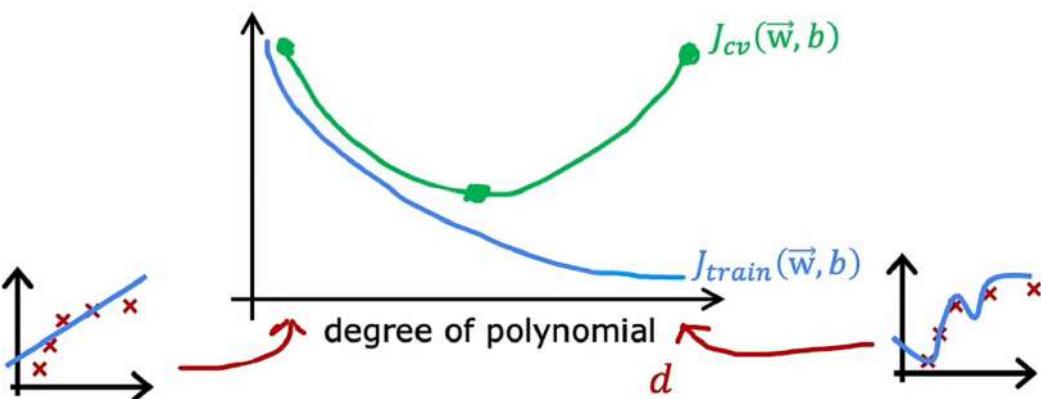
→ Pick $w_1 x_1 + \dots + w_4 x^4 + b$ $(J_{cv}(w^{<4>}, b^{<4>}))$

Estimate generalization error using test the set: $J_{test}(w^{<4>}, b^{<4>})$

Diagnosing bias and variance 偏差 & 方差



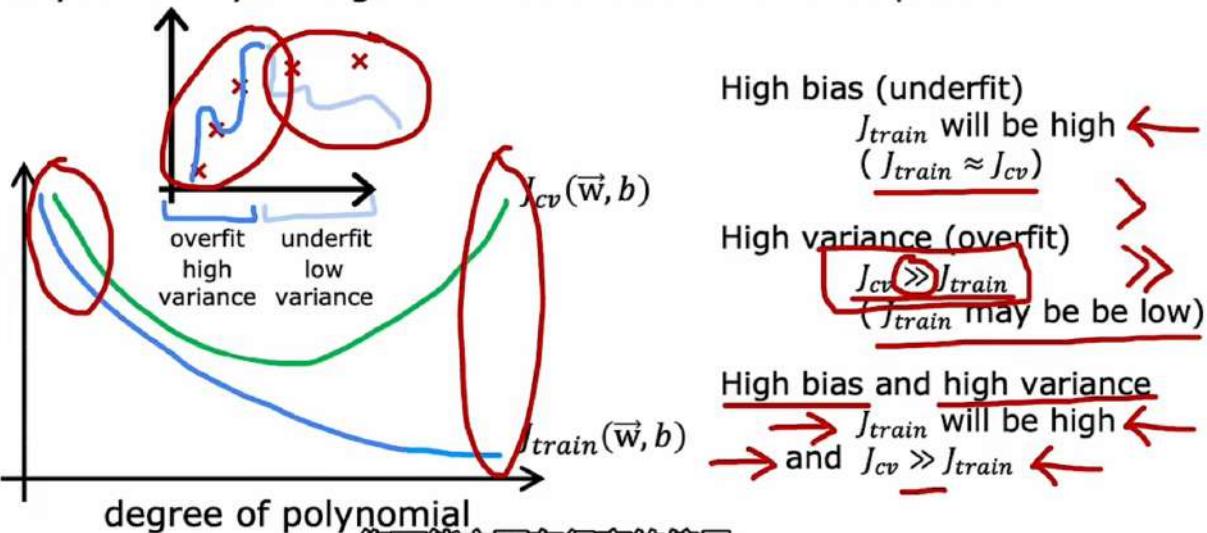
Understanding bias and variance



High
bias
Just
right

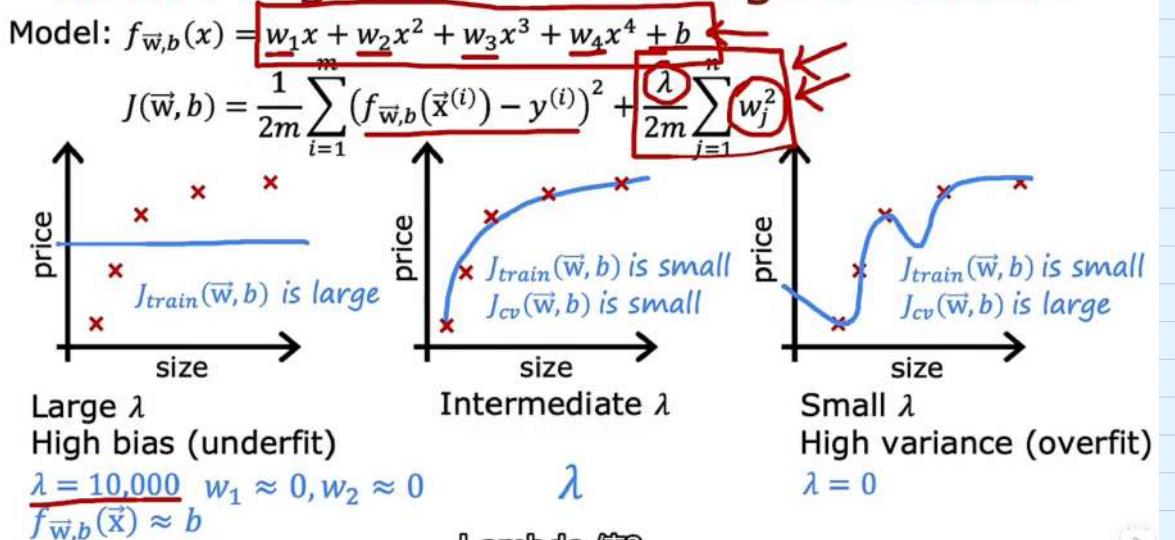
Diagnosing bias and variance

How do you tell if your algorithm has a bias or variance problem?



Regularization and bias variance 影响入取值

Linear regression with regularization



Choosing the regularization parameter λ

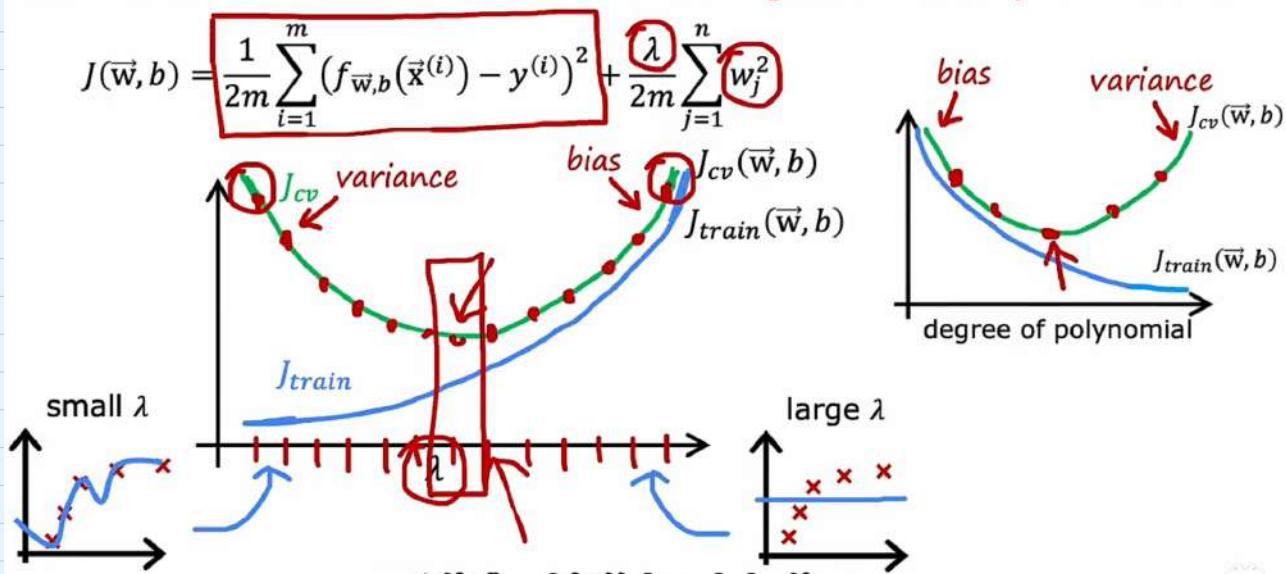
Model: $f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

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

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

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

Bias and variance as a function of regularization parameter λ



Establishing a baseline level of performance

→ 球打 (x12)

Speech recognition example



Human level performance

: 10.6%

Training error J_{train}

: 10.8% ↕ 0.2%

Cross validation error J_{cv}

: 14.8% ↕ 4.0%



又說了，要大量的數據之後

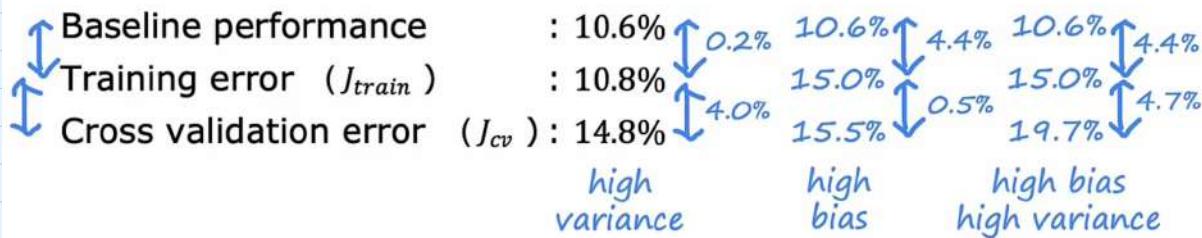
才

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

Bias/variance examples

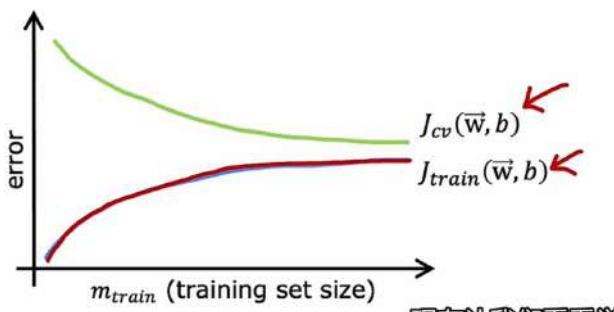


Learning curves (visualize) 并不常用

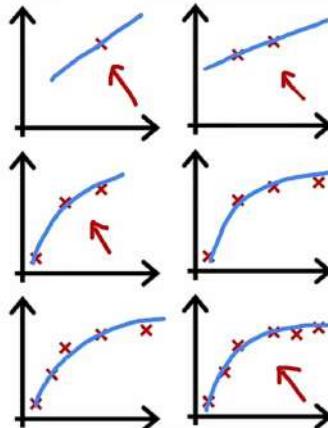
Learning curves

J_{train} = training error

J_{cv} = cross validation error



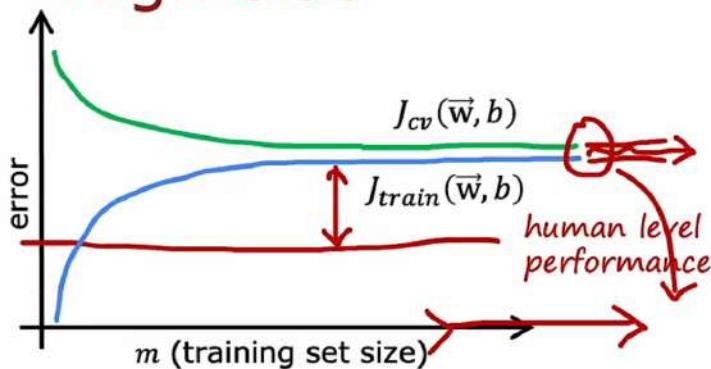
$$f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + b$$



欠拟合对新数据不够敏感，过拟合对新数据特别敏感

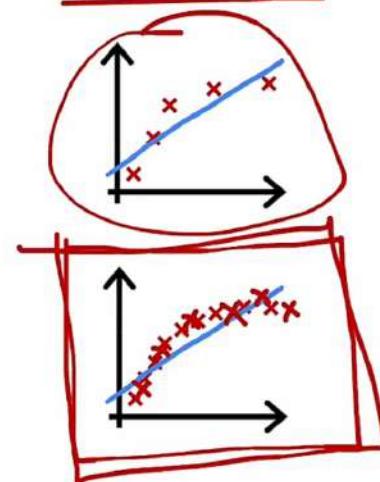
牛蛙，跟前面的知识串起来了

High bias



if a learning algorithm suffers from high bias, getting more training data will not (by itself) help much.

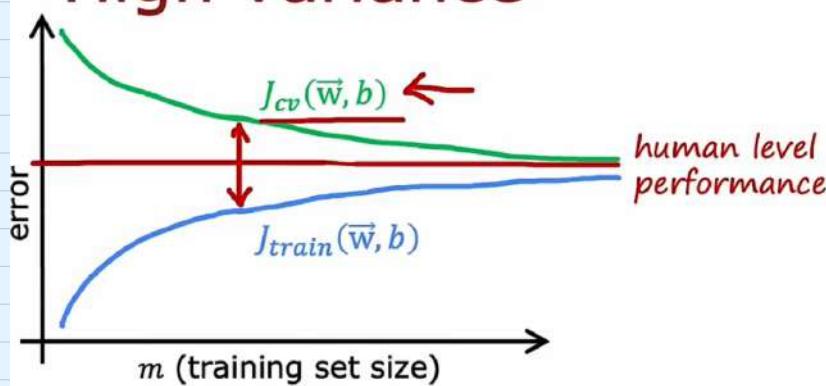
$$f_{\vec{w}, b}(x) = w_1 x + b$$



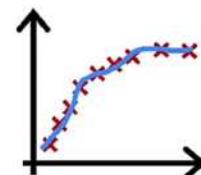
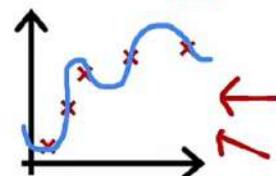
虽然是欠拟合，但是模型不是在训练中被学的

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

High variance



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



Deciding what to try next revisited

Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

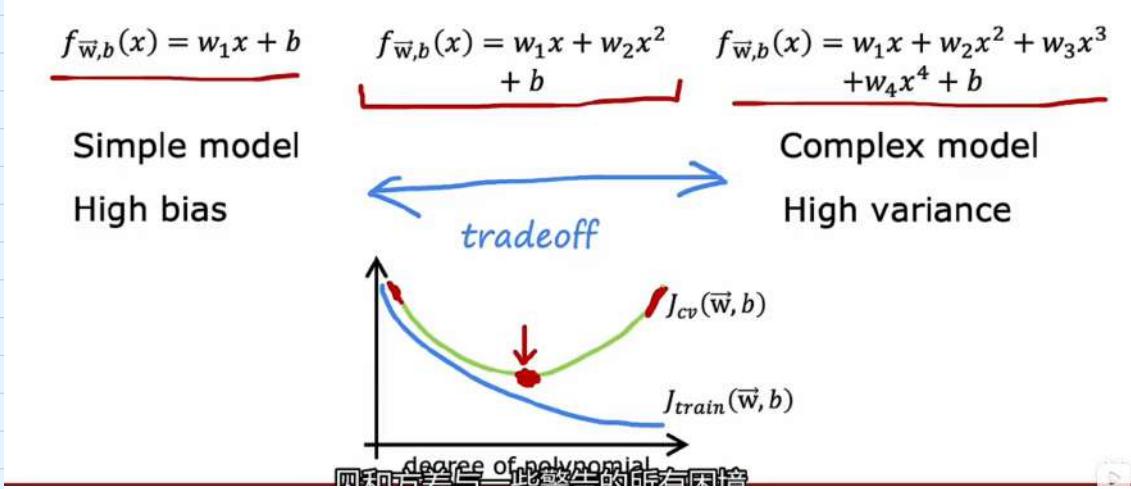
$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

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

- Get more training examples
 - Try smaller sets of features ~~x, x^2, x^3, \dots~~
 - Try getting additional features
 - Try adding polynomial features ~~$(x_1^2, x_2^2, x_1 x_2, \text{etc})$~~
 - Try decreasing λ
 - Try increasing λ
- | |
|---------------------|
| fixes high variance |
| fixes high variance |
| fixes high bias |
| fixes high bias |
| fixes high bias |
| fixes high variance |

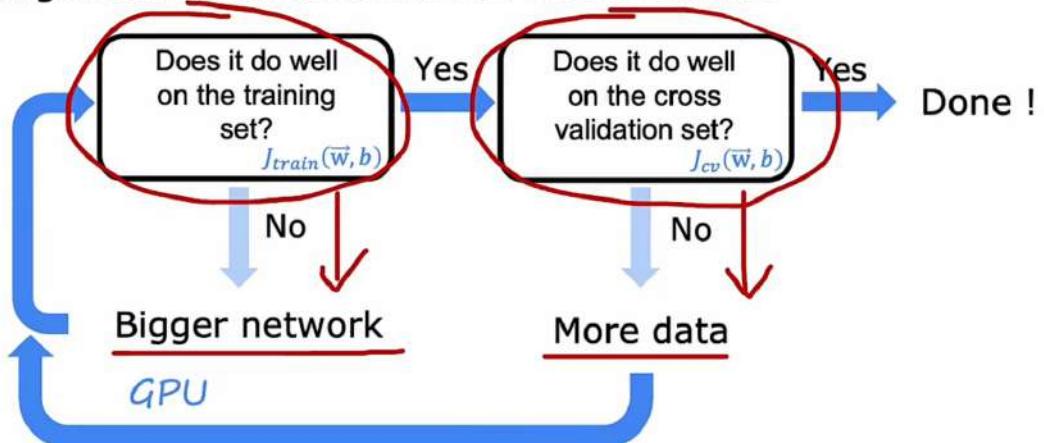
Bias and variance in neural networks

The bias variance tradeoff

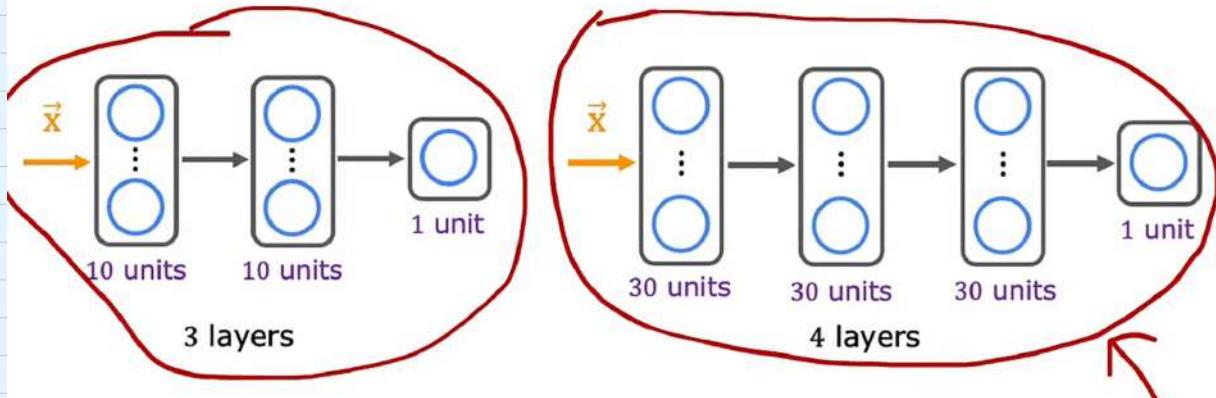


Neural networks and bias variance

Large neural networks are low bias machines



Neural networks and regularization



A large neural network will usually do as well or better than a smaller one so long as regularization is chosen appropriately.

Neural network regularization

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum_{\text{all weights } \mathbf{W}} (w^2)$$

Unregularized MNIST model

```
layer_1 = Dense(units=25, activation="relu")
layer_2 = Dense(units=15, activation="relu")
layer_3 = Dense(units=1, activation="sigmoid")
model = Sequential([layer_1, layer_2, layer_3])
```

Regularized MNIST model

```
layer_1 = Dense(units=25, activation="relu", kernel_regularizer=L2(0.01))
layer_2 = Dense(units=15, activation="relu", kernel_regularizer=L2(0.01))
layer_3 = Dense(units=1, activation="sigmoid", kernel_regularizer=L2(0.01))
model = Sequential([layer_1, layer_2, layer_3])
```

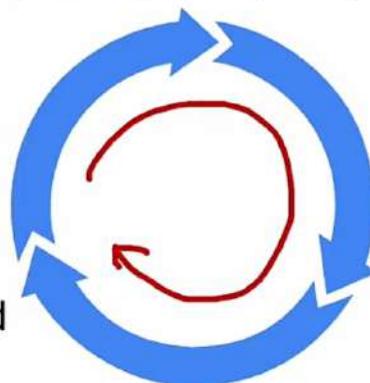
b

λ

Iterative loop of ML development

Iterative loop of ML development

Choose architecture
(model, data, etc.)



Spam classification example

From: cheapsales@buystufffromme.com
To: Andrew Ng
Subject: Buy now!

Deal of the week! Buy now!
Rolex w4tchs - \$100
Medlcine (any kind) - £50
Also low cost M0rgages
available.

From: Alfred Ng
To: Andrew Ng
Subject: Christmas dates?

Hey Andrew,
Was talking to Mom about plans
for Xmas. When do you get off
work. Meet Dec 22?
Alf

Building a spam classifier

Supervised learning: \vec{x} = features of email
 y = spam (1) or not spam (0)

Features: list the top 10,000 words to compute $x_1, x_2, \dots, x_{10,000}$

$$\vec{x} = \begin{bmatrix} 0 \\ 1 \\ \textcolor{red}{\cancel{2}} \textcolor{blue}{\cancel{2}} 1 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \quad \begin{array}{l} a \\ andrew \\ buy \\ deal \\ discount \\ \vdots \end{array}$$

From: cheapsales@buystufffromme.com
To: Andrew Ng
Subject: Buy now!

Deal of the week! Buy now!
Rolex w4tchs - \$100
Medlcine (any kind) - £50
Also low cost M0rgages
available.

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.

不是最有成果的方向。

Error analysis for case above

Error analysis

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

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

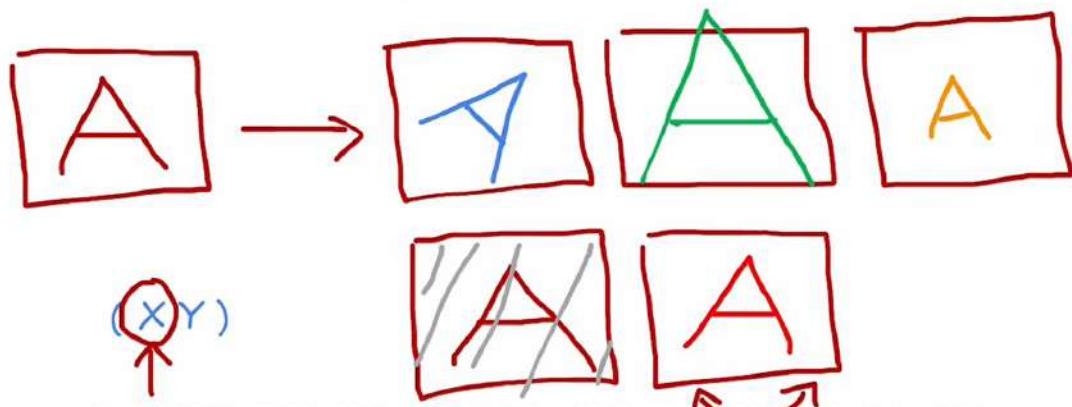
Manually examine $\frac{100}{1000}$ examples and categorize them based on common traits.

- 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

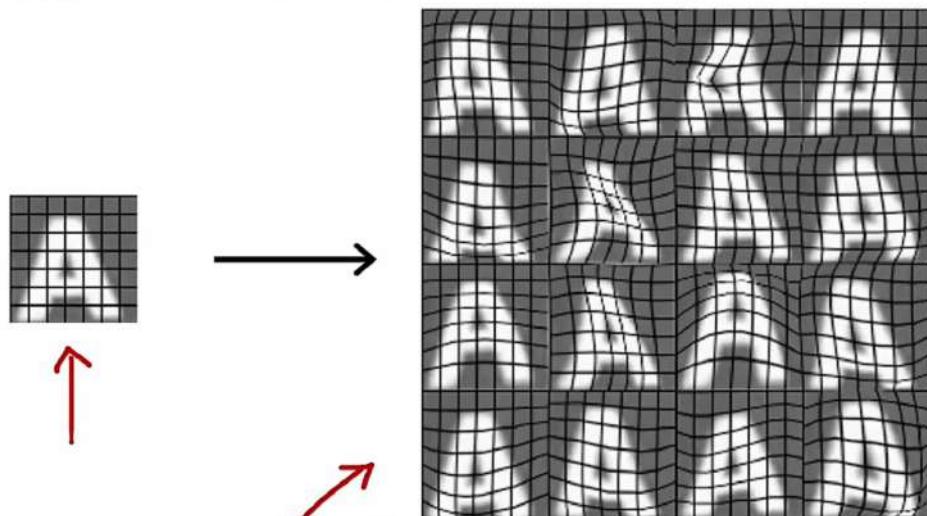
Adding data

Data augmentation

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



Data augmentation by introducing distortions



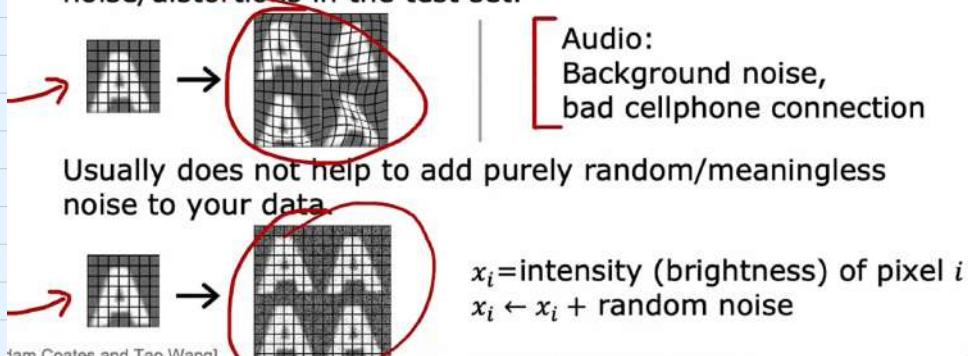
Data augmentation for speech

Speech recognition example

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

Data augmentation by introducing distortions

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



加噪
增加
捨棄性

Data synthesis

Synthesis: using artificial data inputs to create a new training example.

Artificial data synthesis for photo OCR



Artificial data synthesis for photo OCR



Real data



Synthetic data

Engineering the data used by your system

Conventional
model-centric
approach:

AI=Code + Data
(algorithm/model)

Work on this

Data-centric
approach:

AI=Code + Data
(algorithm/model)

Work on this

Transfer Learning 迁移学习

↳ 没那么简单

人的肩膀上

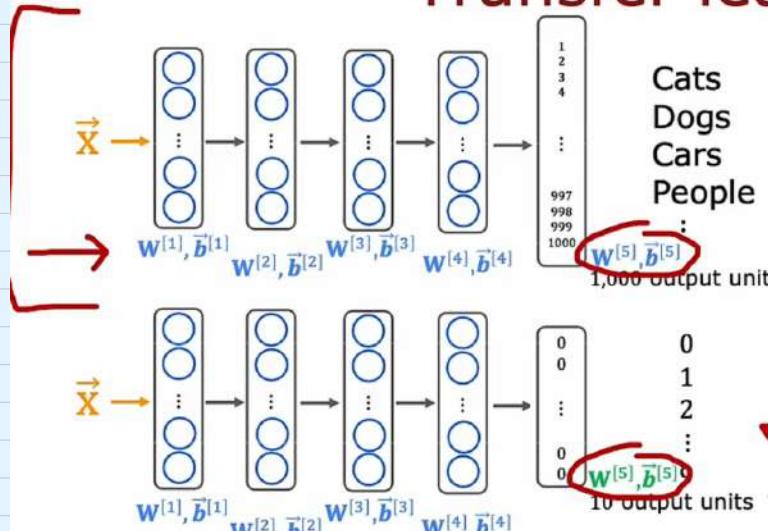
世界上好人多

就怕大师 这方面问题搞僵了

这里看

值了，原来重心一直在迁移学习啊

Transfer learning



Cats
Dogs
Cars
People

1,000
classes

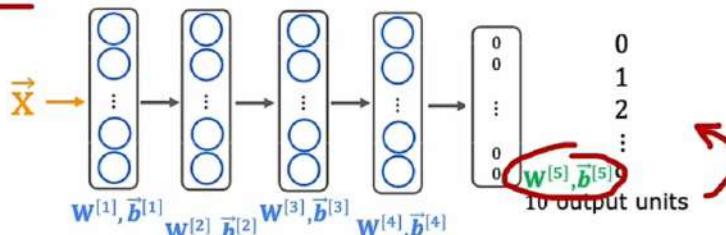
1,000 output units

0, 1, 2, ... , 9

1 million images

Supervised pretraining

Fine tuning



0
0
0
0

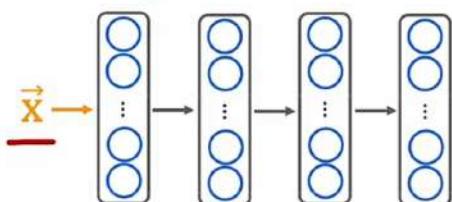
1
2
3
4

10 output units

Option 1: only train output layers parameters.

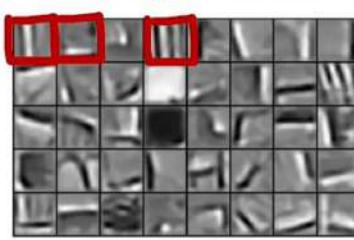
Option 2: train all parameters.

Why does transfer learning work?

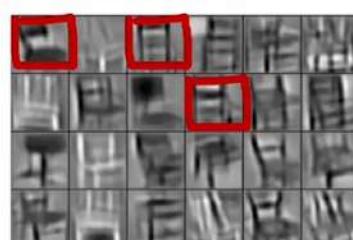


use the same input type

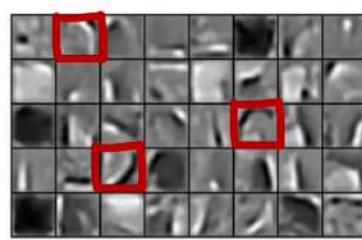
detects edges
detects corners
detects curves/basic shapes



Edges



Corners



Curves / basic shapes

Transfer learning summary

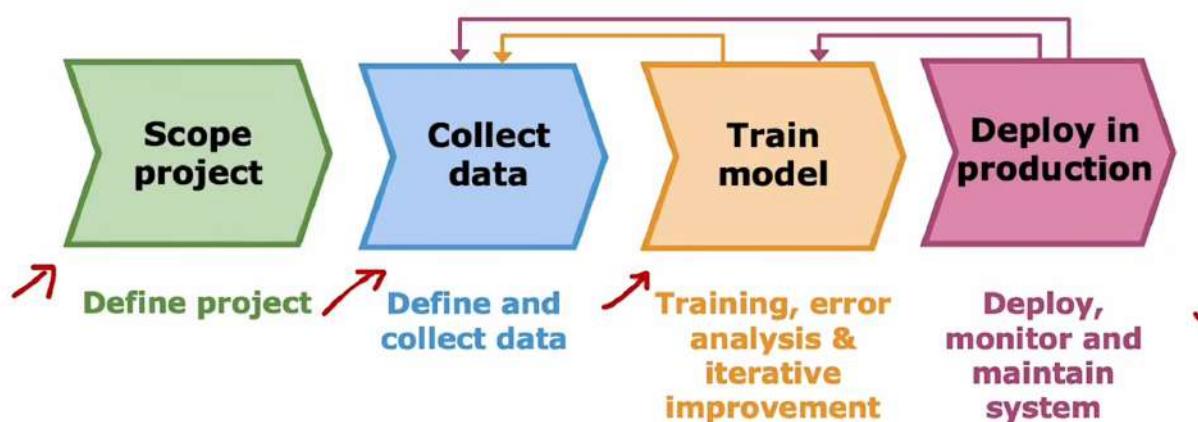
→ 1. Download neural network parameters pretrained on a large dataset with same input type (e.g., images, audio, text) as your application (or train your own). 1M

→ 2. Further train (fine tune) the network on your own data.

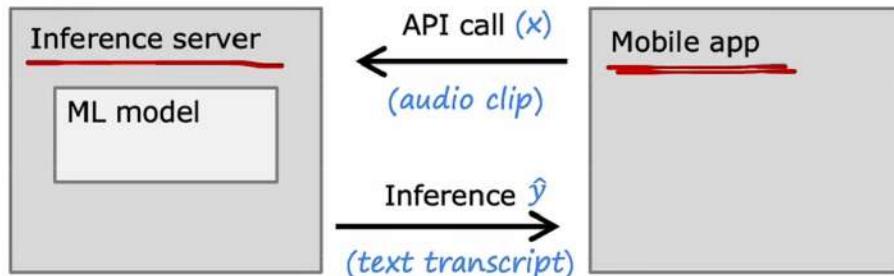
1000
50

Full cycle of machine Learning

Full cycle of a machine learning project



Deployment



→ Software engineering may be needed for:
Ensure reliable and efficient predictions
Scaling
Logging
System monitoring
Model updates

MLOps
machine learning operations

Fairness, bias and ethics (↳ if)

Bias

Hiring tool that discriminates against women.

Facial recognition system matching dark skinned individuals to criminal mugshots.

Biased bank loan approvals.

Toxic effect of reinforcing negative stereotypes.

Adverse use cases

Deepfakes

Spreading toxic/incendiary speech through optimizing for engagement.

Generating fake content for commercial or political purposes.

Using ML to build harmful products, commit fraud etc.

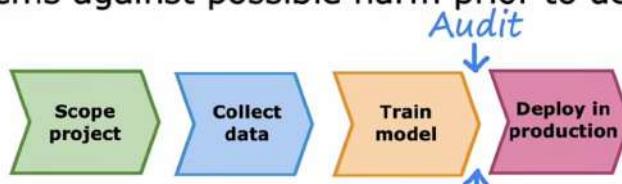
Spam vs anti-spam : fraud vs anti-fraud.

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.



Develop mitigation plan (if applicable), and after deployment, monitor for possible harm.

Error metrics for skewed datasets 10/24/2023
#後方誤差

Rare disease classification example

Train classifier $f_{\vec{w}, b}(\vec{x})$

($y = 1$ if disease present,
 $y = 0$ otherwise)

Find that you've got 1% error on test set
(99% correct diagnoses)

Only 0.5% of patients have the disease

print("y=0")

99.5% accuracy, 0.5% error ←
1% ←
1.2%

Precision/recall

$y = 1$ in presence of rare class we want to detect.

		Actual Class	
		1	0
Predict -ed Class	1	True positive 15	False positive 5
	0	False negative 10	True negative 70
		↓	↓
		25	75

Precision:

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

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

Recall:

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

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

Trading off precision and Recall

在曲线开始急剧下降之前选择一个点

Trading off precision and recall

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

→ Predict 1 if $f_{\vec{w}, b}(\vec{x}) \geq 0.3$

Predict 0 if $f_{\bar{w}, b}(\vec{x}) \leq 0.3$

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

→ higher precision, lower recall.

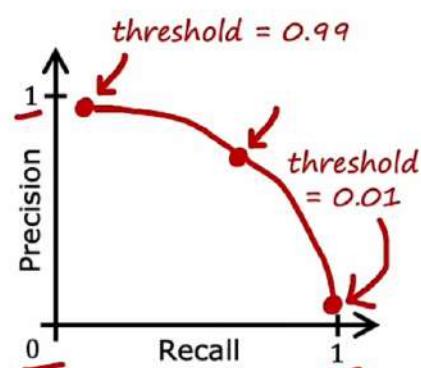
Suppose we want to avoid missing too many cases of rare disease (when in doubt predict $y = 1$)

→ lower precision, higher recall.

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

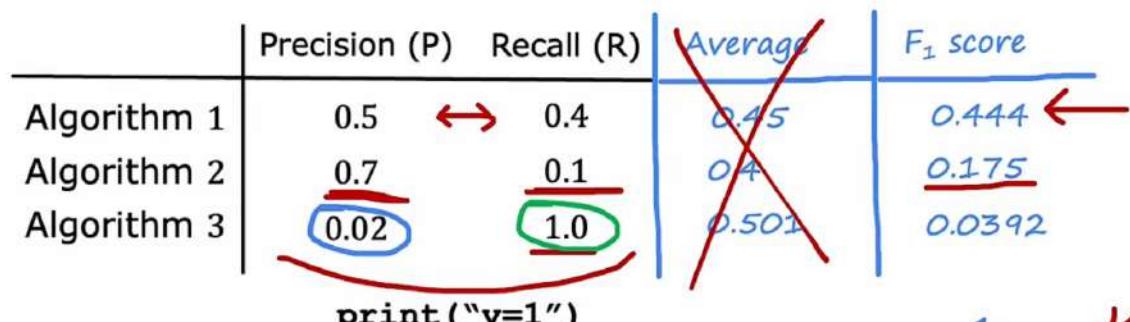
$$\text{precision} = \frac{\text{true positives}}{\text{total predicted positive}}$$

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



F1 score

How to compare precision/recall numbers?



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

$$F1 \text{ score} = \frac{1}{\frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)} = 2 \frac{PR}{P+R}$$

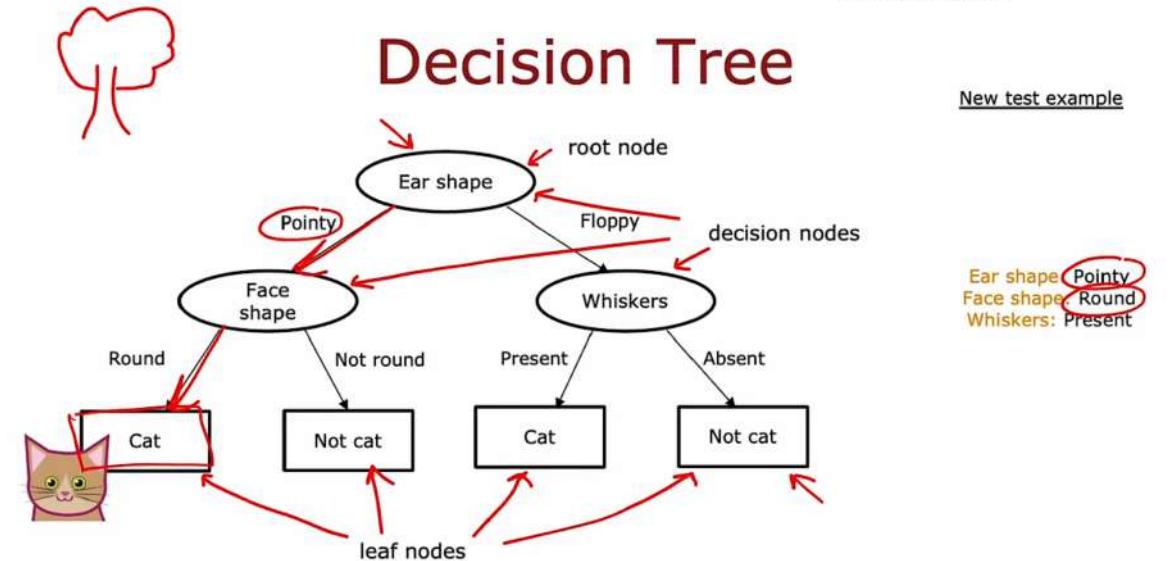
Faith

Decision Tree Model

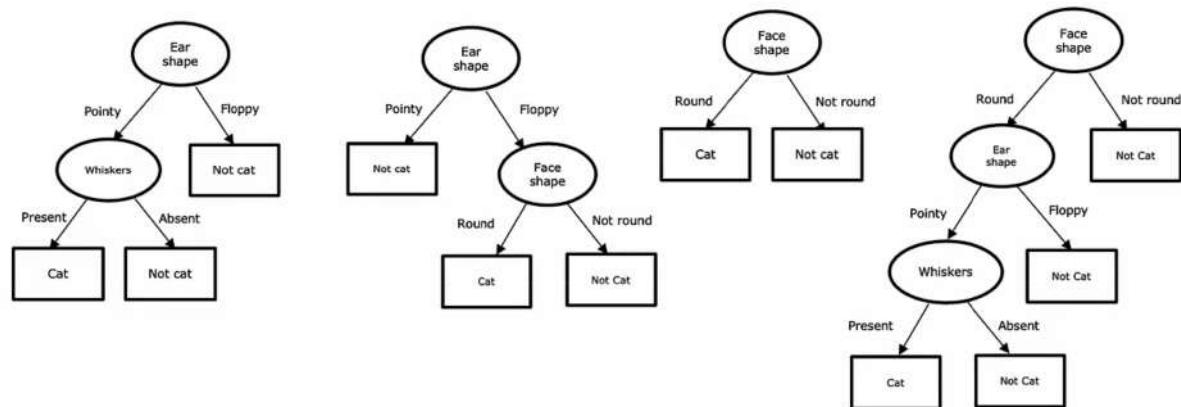
Cat classification example

	Ear shape (x_1)	Face shape (x_2)	Whiskers	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
	Floppy	Not round	Absent	0
	Pointy	Round	Absent	1
	Floppy	Round	Absent	0
	Floppy	Round	Absent	0

Categorical (discrete values)



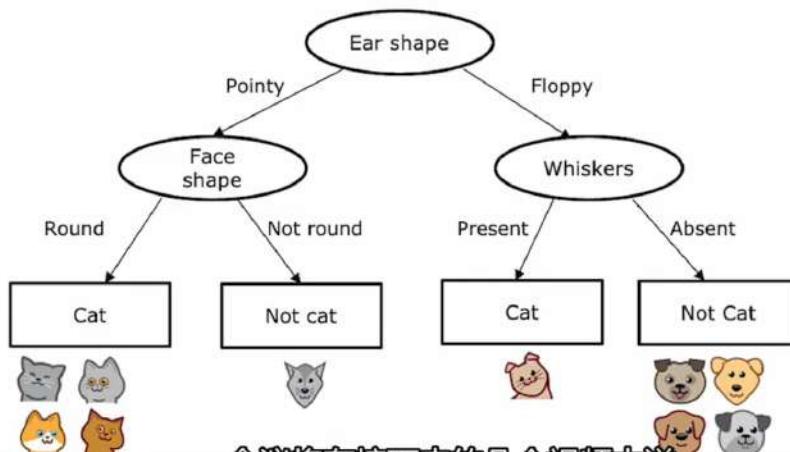
Decision Tree



Learning Process

思考：学习了，做了什么？

Decision Tree Learning

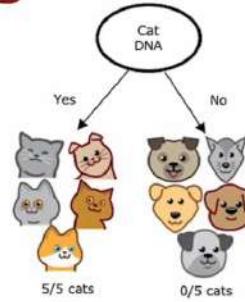
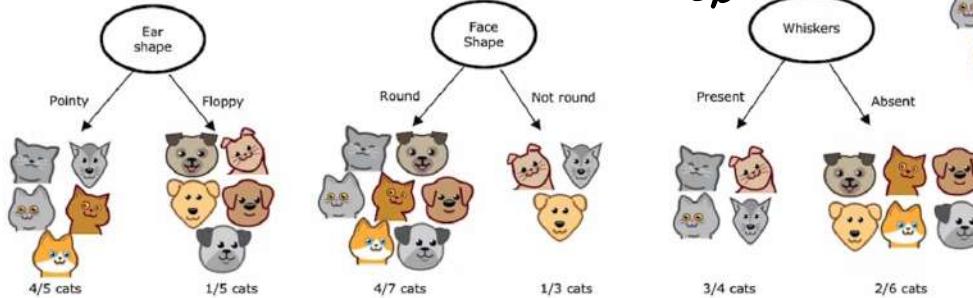


Decision Tree Learning

Decision 1: How to choose what feature to split on at each node?

Maximize purity (or minimize impurity)

→ Cat Information Gain

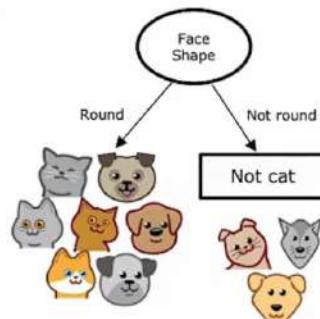


Decision Tree Learning

Decision 2: When do you stop splitting?

→ 停止

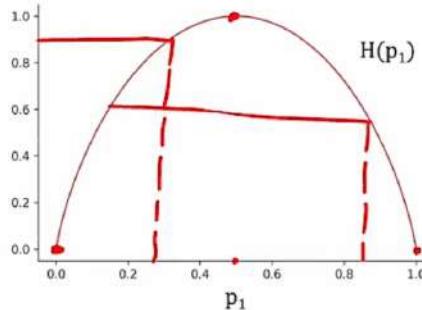
- When a node is 100% one class
- When splitting a node will result in the tree exceeding a maximum depth
- When improvements in purity score are below a threshold
- When number of examples in a node is below a threshold



Masuring purity

Entropy as a measure of impurity

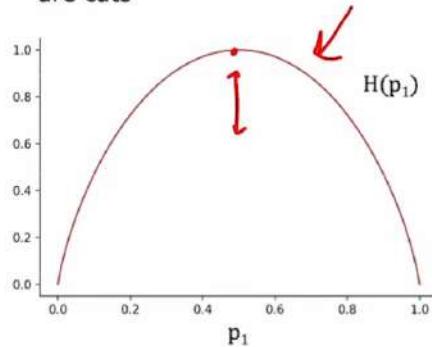
p_1 = fraction of examples that are cats



		$p_1 = 0 \quad H(p_1) = 0$
		$p_1 = 2/6 \quad H(p_1) = 0.92$
		$p_1 = 3/6 \quad H(p_1) = 1$
		$p_1 = 5/6 \quad H(p_1) = 0.65$
		$p_1 = 6/6 \quad H(p_1) = 0$

Entropy as a measure of impurity

p_1 = fraction of examples that are cats



$$p_0 = 1 - p_1$$

$$\begin{aligned} H(p_1) &= -p_1 \log_2(p_1) - p_0 \log_2(p_0) \\ &= -p_1 \log_2(p_1) - (1 - p_1) \log_2(1 - p_1) \end{aligned}$$

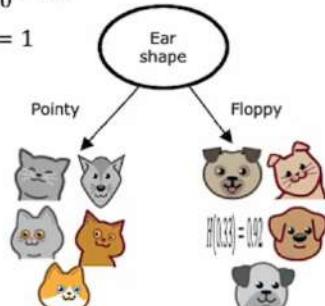
Note: " $0 \log(0) = 0$ "

Choosing a split: information Gain

Choosing a split

$$p_1 = \frac{5}{10} = 0.5$$

$$H(0.5) = 1$$

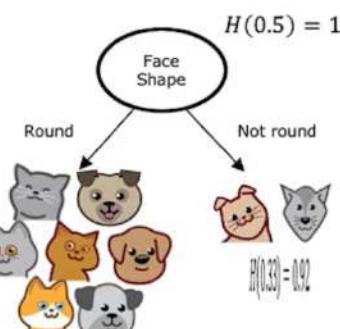


$$p_1 = \frac{4}{5} = 0.8 \quad p_1 = \frac{1}{5} = 0.2$$

$$H(0.8) = 0.72 \quad H(0.2) = 0.72$$

$$H(0.5) - \left(\frac{5}{10} H(0.8) + \frac{5}{10} H(0.2) \right)$$

$$= 0.28$$

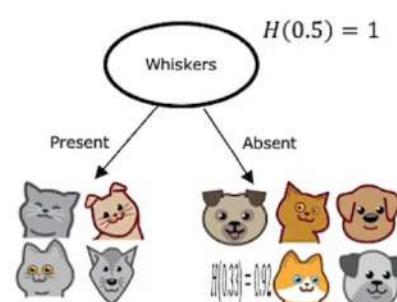


$$p_1 = \frac{4}{7} = 0.57 \quad p_1 = \frac{1}{3} = 0.33$$

$$H(0.57) = 0.99 \quad H(0.33) = 0.92$$

$$H(0.5) - \left(\frac{7}{10} H(0.57) + \frac{3}{10} H(0.33) \right)$$

$$= 0.03$$



$$p_1 = \frac{3}{4} = 0.75 \quad p_1 = \frac{2}{6} = 0.33$$

$$H(0.75) = 0.81 \quad H(0.33) = 0.92$$

$$H(0.5) - \left(\frac{4}{10} H(0.75) + \frac{6}{10} H(0.33) \right)$$

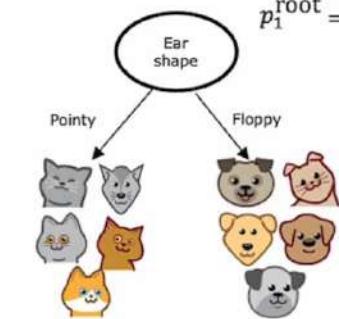
$$= 0.12$$

Information Gain

↑ 10¹⁶²



$$p_1^{\text{root}} = 5/10 = 0.5$$



$$p_1^{\text{left}} = 4/5$$

$$p_1^{\text{right}} = 1/5$$

$$w^{\text{left}} = 5/10$$

$$w^{\text{right}} = 5/10$$

Information gain

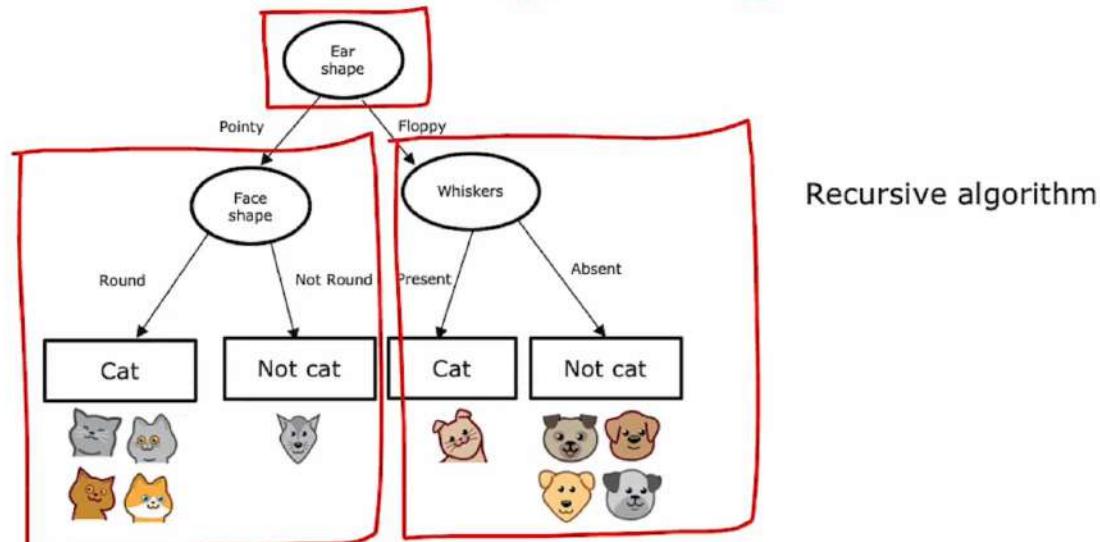
$$= H(p_1^{\text{root}}) - \left(w^{\text{left}} H(p_1^{\text{left}}) + w^{\text{right}} H(p_1^{\text{right}}) \right)$$

Putting it together

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

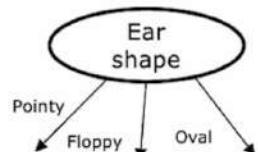
Recursive splitting



Using one-hot encoding of categorical features

Features with three possible values

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



3 possible values

Ear-shape	Pointy ears	Floppy ears	Oval 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	0
	0	0	1	Round	Absent	1
	0	1	0	Round	Absent	0
	0	1	0	Round	Absent	0

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

One hot encoding and neural networks

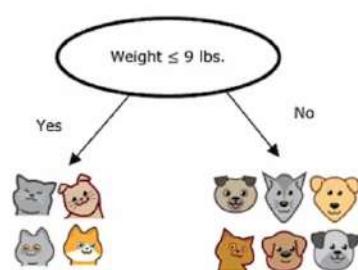
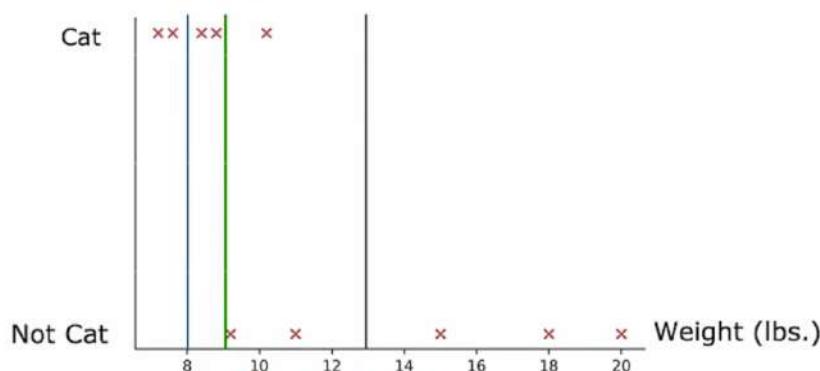
	Pointy ears	Floppy ears	Round ears	Face shape	Whiskers	Cat
	1	0	0	Round 1	Present 1	1
	0	0	1	Not round 0	Present 1	1
	0	0	1	Round 1	Absent 0	0
	1	0	0	Not round 0	Present 1	0
	0	0	1	Round 1	Present 1	1
	1	0	0	Round 1	Absent 0	1
	0	1	0	Not round 0	Absent 0	1
	0	0	1	Round 1	Absent 0	1
	0	1	0	Round 1	Absent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

Continuous valued feature

Continuous features ↗

Ear shape	Face shape	Whiskers	Weight (lbs.)	Cat
	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

Splitting on a continuous variable



$$H(0.5) - \left(\frac{2}{10} H\left(\frac{2}{2}\right) + \frac{8}{10} H\left(\frac{3}{8}\right) \right) = 0.24$$

$$H(0.5) - \left(\frac{4}{10} H\left(\frac{4}{4}\right) + \frac{6}{10} H\left(\frac{1}{6}\right) \right) = 0.61$$

$$H(0.5) - \left(\frac{7}{10} H\left(\frac{5}{7}\right) + \frac{3}{10} H\left(\frac{0}{3}\right) \right) = 0.40$$

用不等式解决

连续决策

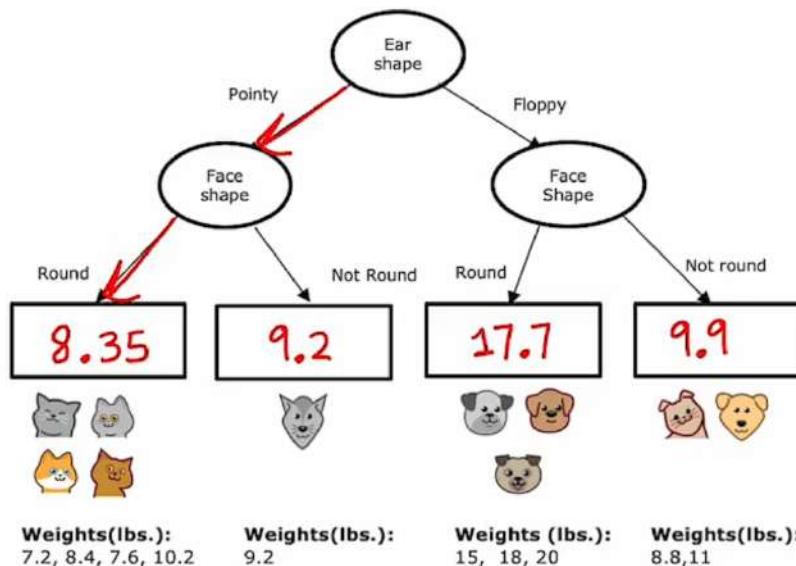
Regression Trees

Regression with Decision Trees: Predicting a number

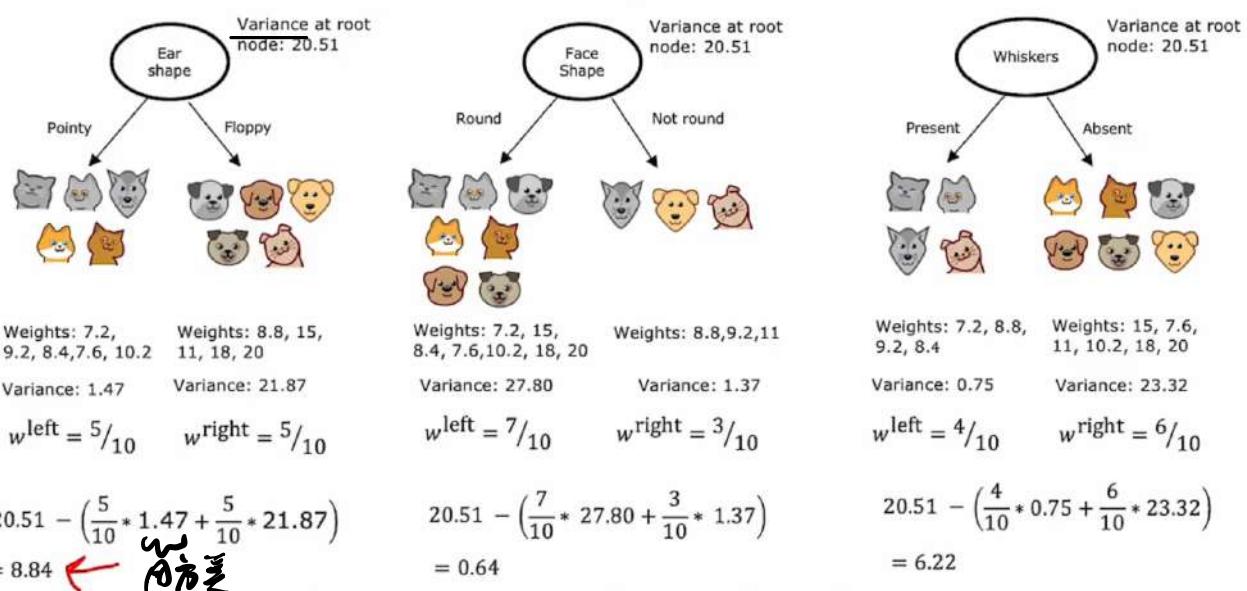
Ear shape	Face shape	Whiskers	Weight (lbs.)
Pointy	Round	Present	7.2
Floppy	Not round	Present	8.8
Floppy	Round	Absent	15
Pointy	Not round	Present	9.2
Pointy	Round	Present	8.4
Pointy	Round	Absent	7.6
Floppy	Not round	Absent	11
Pointy	Round	Absent	10.2
Floppy	Round	Absent	18
Floppy	Round	Absent	20

X Y

Regression with Decision Trees



Choosing a split

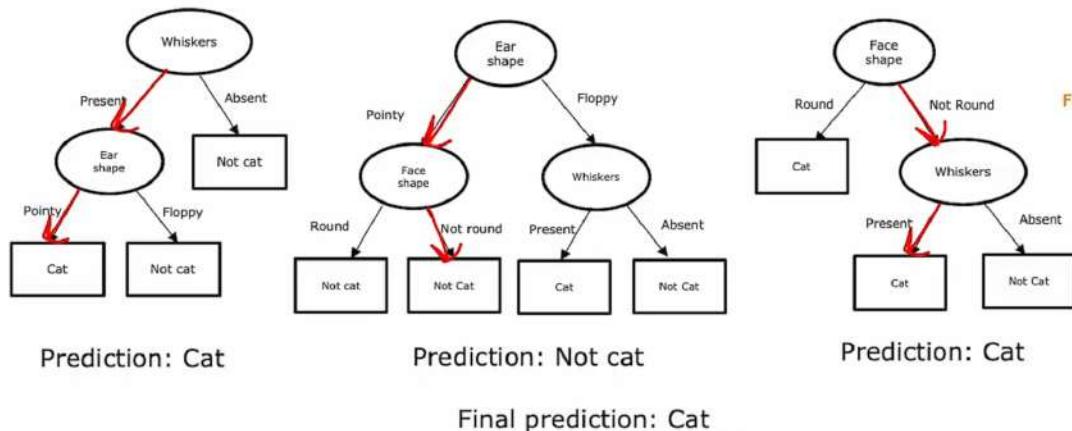


Using multiple decision trees

Trees are highly sensitive to small changes of the data



Tree ensemble



New test example



Ear shape: Pointy
Face shape: Not Round
Whiskers: Present

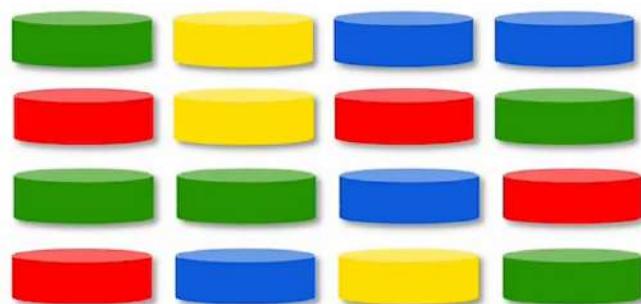
Sampling with replacement 有放回抽样

Sampling with replacement

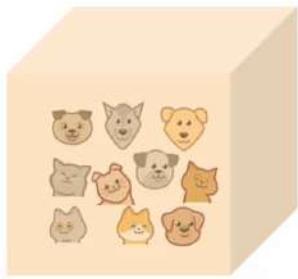
Tokens



Sampling with replacement:



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

Random forest algorithm 隨机森林

Generating a tree sample

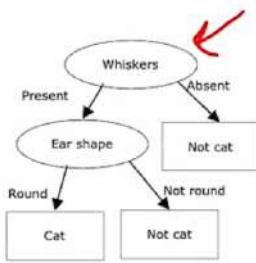
Given training set of size m

For $b = 1$ to B

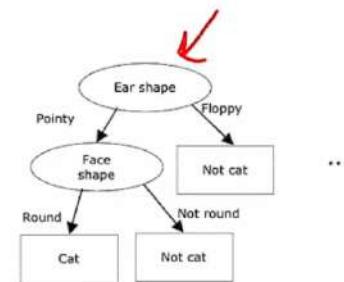
Use sampling with replacement to create a new training set of size m

Train a decision tree on the new dataset

Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	Yes
Floppy	Round	Absent	No
Floppy	Round	Absent	No
Pointy	Round	Present	Yes
Pointy	Not Round	Present	Yes
Floppy	Not Round	Absent	No
Floppy	Not Round	Absent	No
Pointy	Not Round	Absent	No
Pointy	Not Round	Present	Yes



Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	Yes
Pointy	Round	Absent	Yes
Floppy	Not Round	Absent	No
Floppy	Not Round	Absent	No
Pointy	Not Round	Present	Yes
Floppy	Not Round	Absent	No
Floppy	Not Round	Absent	No
Pointy	Not Round	Absent	No
Pointy	Not Round	Present	Yes



Bagged decision tree

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.

$$k = \sqrt{n}$$

随机 100% 特征

Random forest algorithm

Boosted trees intuition

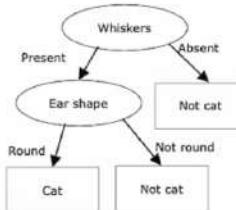
Given training set of size m

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 examples that the previously trained trees misclassify

Train a decision tree on the new dataset

Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	Yes
Floppy	Round	Absent	No
Floppy	Round	Absent	No
Pointy	Round	Present	Yes
Pointy	Not Round	Present	Yes
Floppy	Round	Absent	No
Floppy	Round	Present	Yes
Pointy	Not Round	Absent	No
Pointy	Not Round	Present	Yes
Pointy	Not Round	Present	Yes



Boosted trees intuition

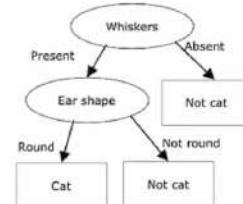
Given training set of size m

For $b = 1$ to \underline{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 examples that the previously trained trees misclassify

Train a decision tree on the new dataset

Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	Yes
Floppy	Round	Absent	No
Floppy	Round	Absent	No
Pointy	Round	Present	Yes
Pointy	Not Round	Present	Yes
Floppy	Round	Absent	No
Floppy	Round	Present	Yes
Pointy	Not Round	Absent	No
Pointy	Not Round	Absent	No
Pointy	Not Round	Present	Yes



Ear shape	Face shape	Whiskers	Prediction
Pointy	Round	Present	Cat ✓
Floppy	Not Round	Present	Not cat ✗
Floppy	Round	Absent	Not cat ✓
Pointy	Not Round	Present	Not cat ✓
Pointy	Round	Present	Cat ✓
Floppy	Not Round	Absent	Not cat ✗
Floppy	Round	Absent	Not cat ✓
Floppy	Round	Absent	Not cat ✗
Pointy	Not Round	Absent	Not cat ✓

1, 2, ..., b-1 ↗
b

XGBoost (eXtreme Gradient Boosting)

- Open source implementation of boosted trees
- Fast efficient implementation
- Good choice of default splitting criteria and criteria for when to stop splitting
- Built in regularization to prevent overfitting
- Highly competitive algorithm for machine learning competitions (eg: Kaggle competitions)

Using XGBoost

Classification



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

Regression

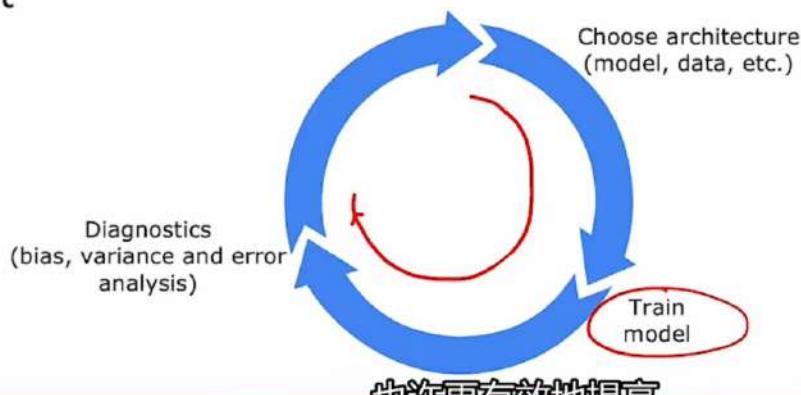
```
from xgboost import XGBRegressor  
model = XGBRegressor()  
model.fit(X_train, y_train)  
y_pred = model.predict(X_test)
```

When to use decision trees decision tree Vs Neural Networks

Decision Trees vs Neural Networks

Decision Trees and Tree ensembles

- Works well on tabular (structured) data
- Not recommended for unstructured data (images, audio, text)
- Fast



Decision Trees vs Neural Networks

Decision Trees and Tree ensembles

- Works well on tabular (structured) data
- Not recommended for unstructured data (images, audio, text)
- Fast
- Small decision trees may be human interpretable

Neural Networks

- Works well on all types of data, including tabular (structured) and unstructured data
- May be slower than a decision tree
- Works with transfer learning
- When building a system of multiple models working together, it might be easier to string together multiple neural networks