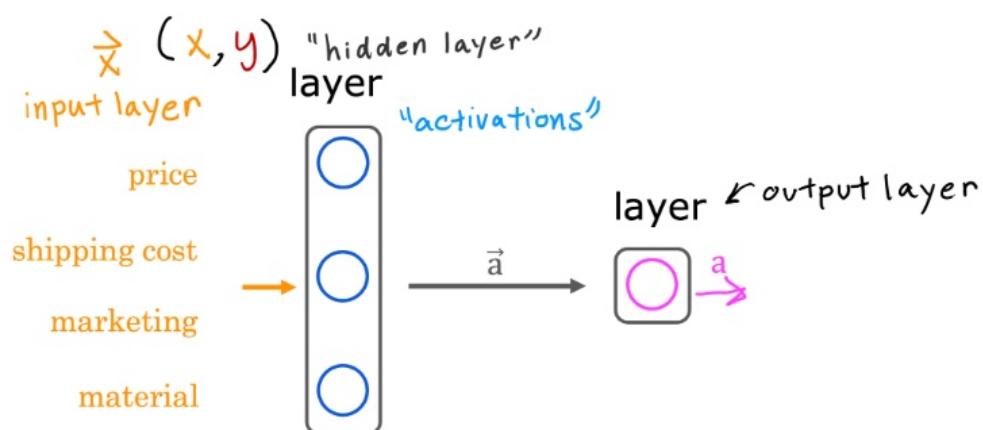


1.

1/1分



Which of these are terms used to refer to components of an artificial neural network?

activation function

正确

Yes, an activation is the number calculated by a neuron (and an activation is a vector that is output by a layer that contains multiple neurons).

layers

正确

Yes, a layer is a grouping of neurons in a neural network

neurons

正确

Yes, a neuron is a part of a neural network

axon

1.

1/1分

For a neural network, here is the formula for calculating the activation of the third neuron in layer 2, given the activation vector from layer 1: $a_3^{[2]} = g(\vec{w}_3^{[2]} \cdot \vec{d}^{[1]} + b_3^2)$. Which of the following are correct statements?

- The activation of layer 2 is determined using the activations from the previous layer.

正确

Correct. The previous layer's activations can be considered as the input features into the current layer.

- The activation of unit 3 (neuron 3) of layer 2 is calculated using a parameter vector \vec{w} and b that are specific to unit 3 (neuron 3).

正确

Correct. Each neuron has its own set of parameters \vec{w} and b .

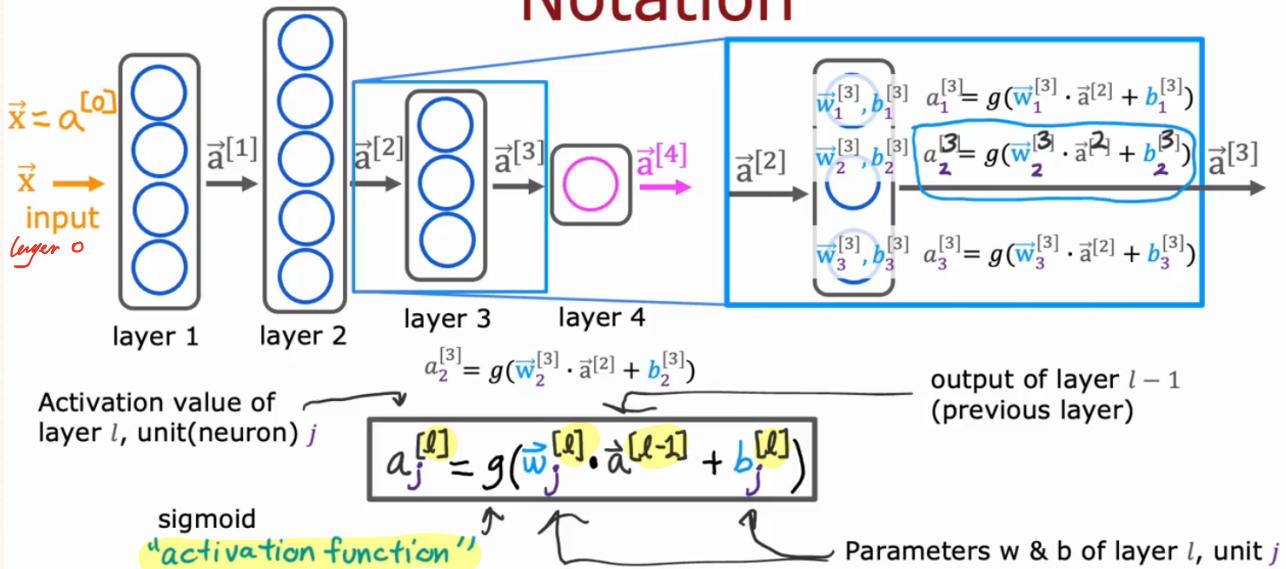
- Unit 3 (neuron 3) outputs a single number (a scalar).

正确

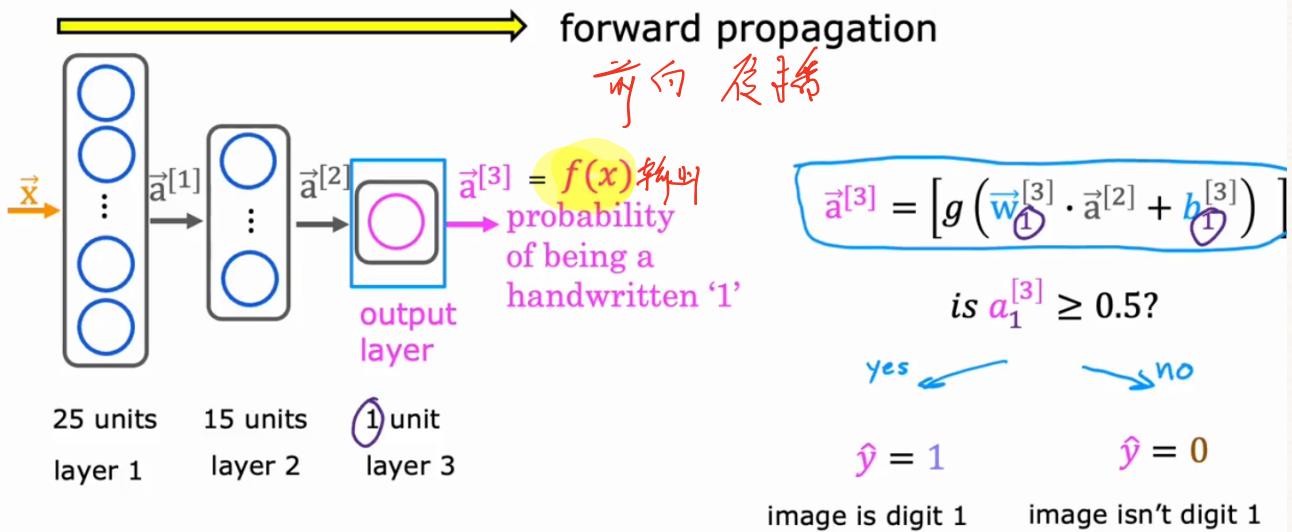
Correct, the calculation above for unit 3 (neuron 3) is a single number (and not a vector).

- If you are calculating the activation for layer 1, then the previous layer's activations would be denoted by \vec{d}^{-1}

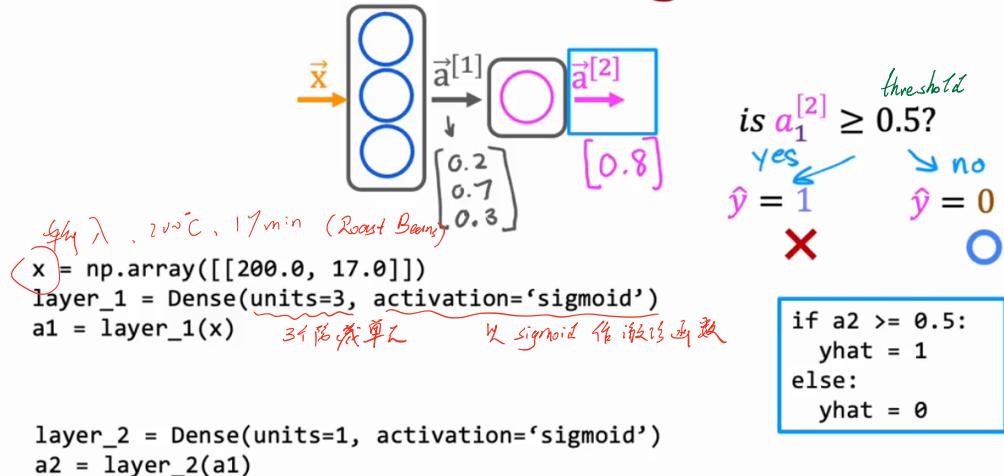
Notation



Handwritten digit recognition



Build the model using TensorFlow

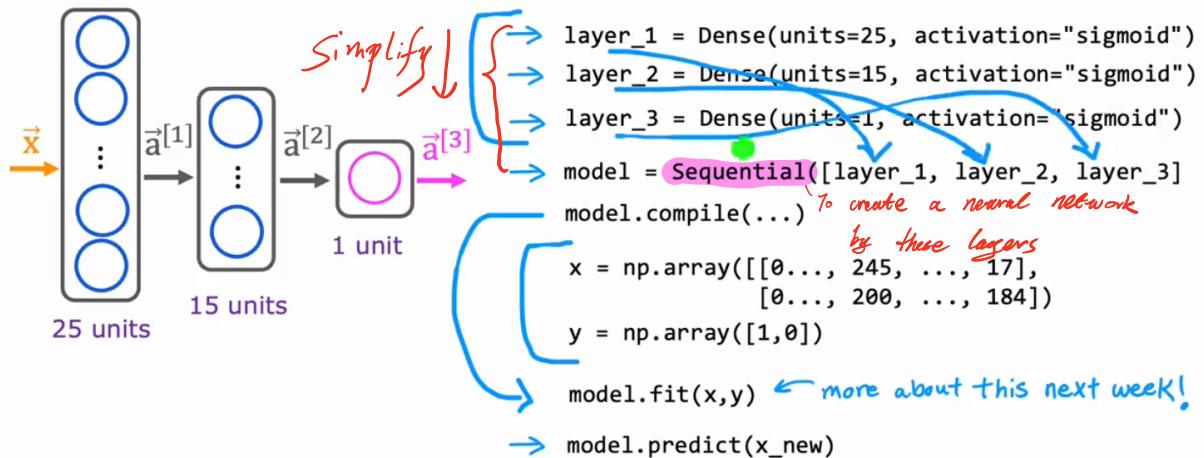


Note about numpy arrays

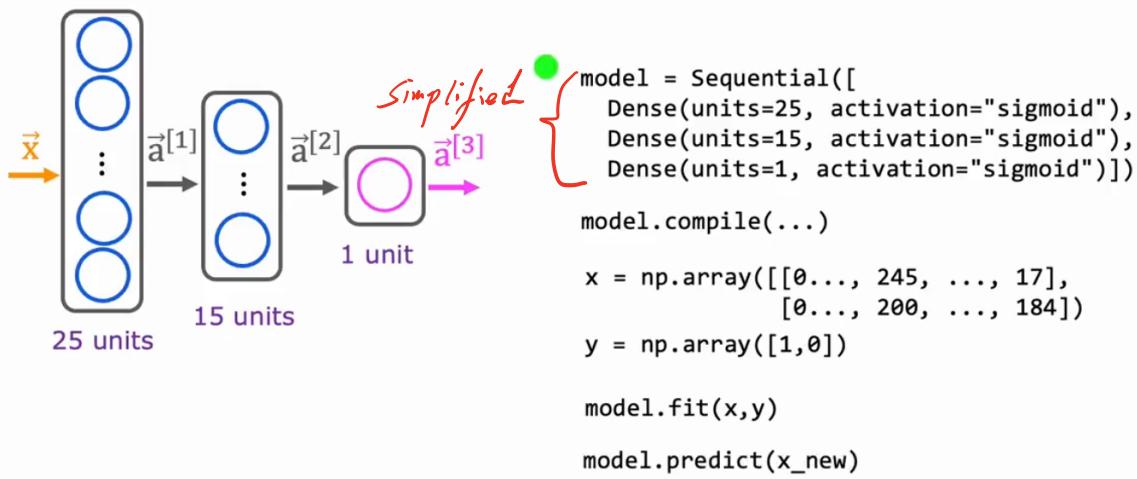
$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$	$x = np.array([1, 2, 3, 4, 5, 6])$	2D array
2 rows 3 columns 2 x 3 matrix	$\begin{bmatrix} 1, 2, 3 \\ 4, 5, 6 \end{bmatrix}$	2 x 3
$\begin{bmatrix} 0.1 & 0.2 \\ -3 & -4 \\ -0.5 & -0.6 \\ 7 & 8 \end{bmatrix}$	$x = np.array([[0.1, 0.2], [-3.0, -4.0], [-0.5, -0.6], [7.0, 8.0]])$	4 x 2
4 rows 2 columns 4 x 2 matrix	$\begin{bmatrix} [0.1, 0.2] \\ [-3.0, -4.0] \\ [-0.5, -0.6] \\ [7.0, 8.0] \end{bmatrix}$	1 x 2
		2 x 1

Fix the tensorflow of v2 Tensor (2nd Martin) fix it

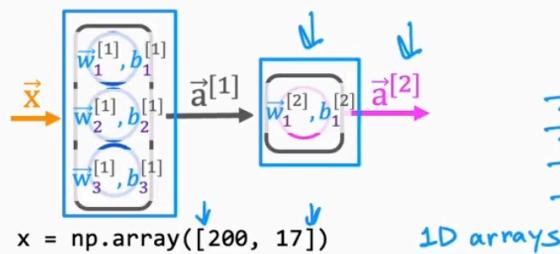
Digit classification model



Digit classification model



forward prop (coffee roasting model)



Complicated

$$a_1^{[2]} = g(\vec{w}_1^{[2]} \cdot \vec{a}_1^{[1]} + b_1^{[2]})$$

$$\rightarrow w2_1 = \text{np.array}([-7, 8])$$

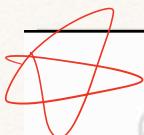
$$\rightarrow b2_1 = \text{np.array}([3])$$

$$\rightarrow z2_1 = \text{np.dot}(w2_1, a1) + b2_1$$

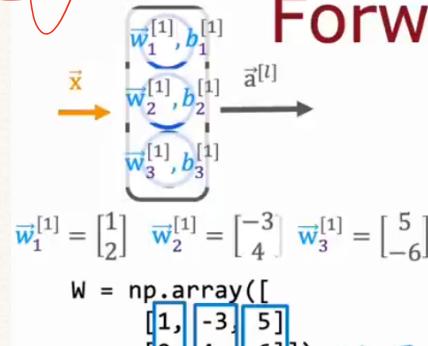
$$\rightarrow a2_1 = \text{sigmoid}(z1_1)$$

$w_1^{[2]}$
 $w_1 \rightarrow w2_1$

```
w1_1 = np.array([1, 2])      w1_2 = np.array([-3, 4])      w1_3 = np.array([5, -6])
b1_1 = np.array([-1])        b1_2 = np.array([1])          b1_3 = np.array([2])
z1_1 = np.dot(w1_1, x) + b  z1_2 = np.dot(w1_2, x) + b  z1_3 = np.dot(w1_3, x) + b
a1_1 = sigmoid(z1_1)        a1_2 = sigmoid(z1_2)        a1_3 = sigmoid(z1_3)
a1    = np.array([a1_1, a1_2, a1_3])
```



Forward prop in NumPy



$$b_1^{[l]} = -1 \quad b_2^{[l]} = 1 \quad b_3^{[l]} = 2$$

$$b = \text{np.array}([-1, 1, 2])$$

$$\vec{a}^{[0]} = \vec{x}$$

$$a_{in} = \text{np.array}([-2, 4])$$

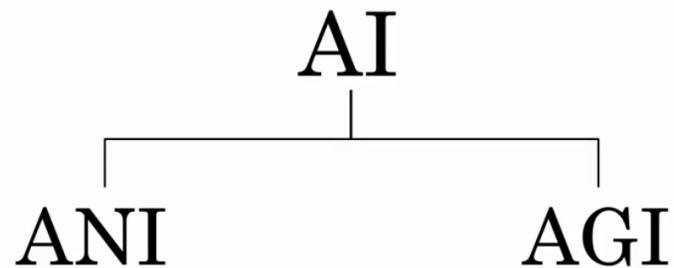
定义层

```
def dense(a_in, W, b, g):
    units = W.shape[1] [0,0,0]
    a_out = np.zeros(units)
    for j in range(units): 0,1,2
        w = W[:,j] for i in range(2)
        z = np.dot(w, a_in) + b[j]
        a_out[j] = g(z) sigmoid
```

定义网络

```
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
```

capital W refers to a matrix



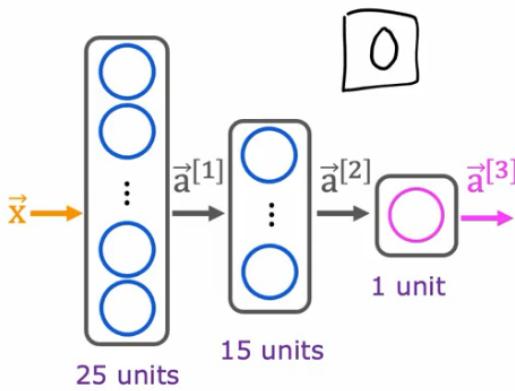
ANI
(artificial narrow intelligence)

E.g., smart speaker,
self-driving car, web search,
AI in farming and factories

AGI
(artificial general intelligence)

Do anything a human can do

Train a Neural Network in TensorFlow



Given set of (x, y) examples
How to build and train this in code?

```
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'),
])
from tensorflow.keras.losses import BinaryCrossentropy
model.compile(loss=BinaryCrossentropy())
model.fit(X, Y, epochs=100) epochs: number of steps in gradient descent
```

① ② ③

Model Training Steps

	TensorFlow
① specify how to compute output given input x and parameters w, b (define model)	<p>logistic regression</p> $z = \mathbf{np}.dot(w, x) + b$ $f_x = 1/(1+\mathbf{np}.exp(-z))$
② specify loss and cost	<p>logistic loss</p> $\text{loss} = -y * \mathbf{np}.log(f_x) - (1-y) * \mathbf{np}.log(1-f_x)$
③ Train on data to minimize $J(\vec{w}, b)$	<p>neural network</p> $\text{model} = \text{Sequential}([\text{Dense}(...), \text{Dense}(...), \text{Dense}(...)])$ <p><i>y=0/1</i></p> <p>binary cross entropy</p> $\text{model.compile(loss=BinaryCrossentropy())}$ $\text{model.fit}(X, y, \text{epochs}=100)$

Logistic regression
(2 possible output values)

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

0.71

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

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

0.29

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

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

parameters w_1, w_2, \dots, w_N
 b_1, b_2, \dots, b_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) $y=1, 2, 3, 4$

$$\text{X } z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

= $P(y=1|\vec{x})$ 0.30

$$\text{O } z_2 = \vec{w}_2 \cdot \vec{x} + b_2$$

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

= $P(y=2|\vec{x})$ 0.20

$$\square z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

= $P(y=3|\vec{x})$ 0.15

$$\Delta z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

= $P(y=4|\vec{x})$ 0.35

DeepLearning.AI Stanford | ONLINE

但它最终会变成逻辑回归模型。

Andrew Ng

Cost

Logistic regression

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

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

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

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

if $y=1$ if $y=0$

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

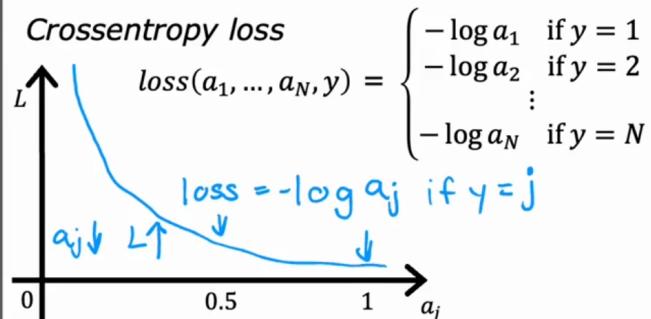
Softmax regression

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

$$\vdots$$

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

Crossentropy loss

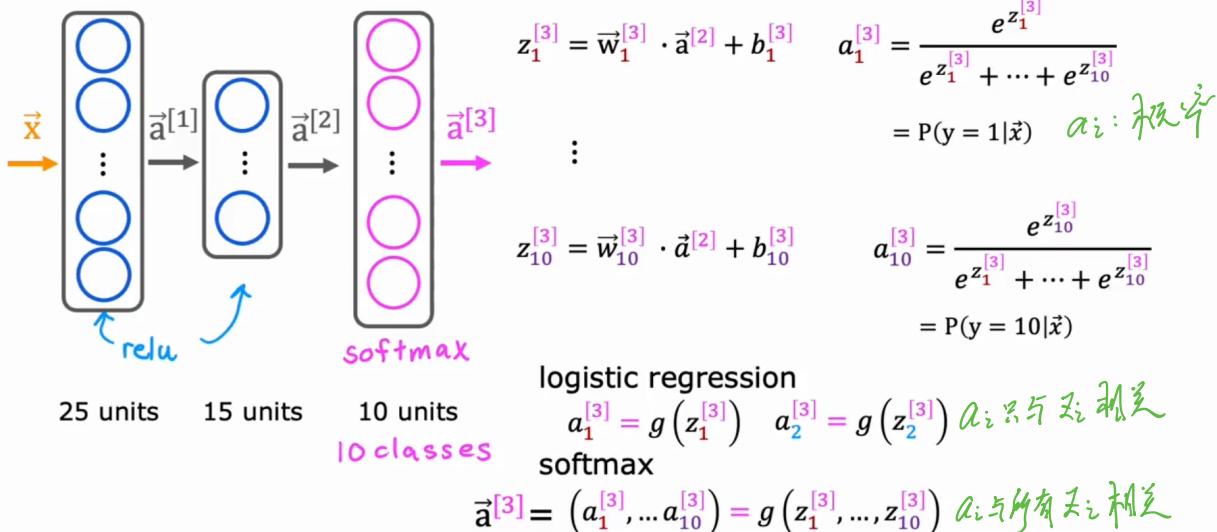


DeepLearning.AI Stanford | ONLINE

例如，如果 $y=2$ ，

Andrew Ng

Neural Network with Softmax output



MNIST with softmax

① specify the model

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

```
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')
])
from tensorflow.keras.losses import
    SparseCategoricalCrossentropy,
model.compile(loss= SparseCategoricalCrossentropy() )
model.fit(X, Y, epochs=100)
Note: better (recommended) version later.
```

Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

Logistic regression:

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

Original loss

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

$| + \frac{1}{10,000}$ $| - \frac{1}{10,000}$

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='sigmoid')])
```

~~model.compile(loss=BinaryCrossEntropy())~~

More accurate loss (in code)

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

logits: z

More numerically accurate implementation of softmax

Softmax regression

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

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

~~model = Sequential([~~

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

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

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

'linear'

~~model.compile(loss=SparseCategoricalCrossEntropy())~~

More Accurate

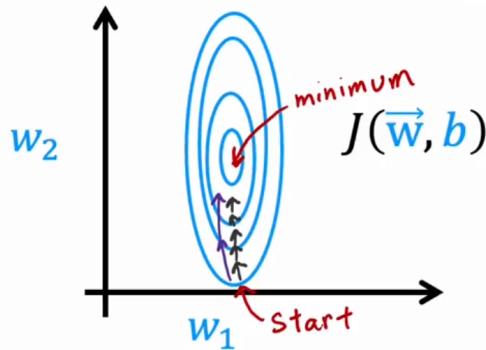
$$L(\vec{a}, y) = \begin{cases} -\log \frac{e^{z_1}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 1 \\ \vdots \\ -\log \frac{e^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 10 \end{cases}$$

~~model.compile(loss=SparseCrossEntropy(from_logits=True))~~

除了这个推荐的版本，在数字上更准确，

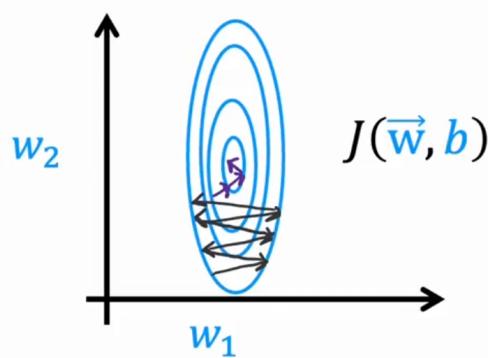
except that the version that is recommended is more numerically accurate,

Adam Algorithm Intuition



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

*adjust α
automatically*

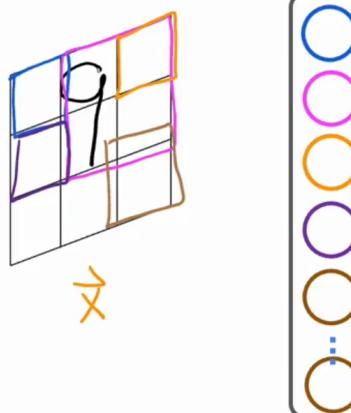


If w_j (or b) keeps oscillating,
reduce α_j .

Convolutional Layer

卷积

ReLU

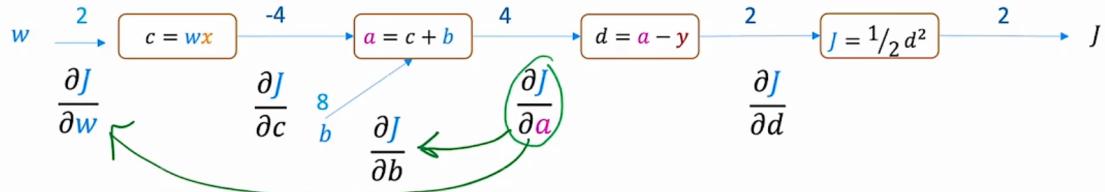


Each Neuron only looks at part of the previous layer's inputs.

Why?

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

Backprop is an efficient way to compute derivatives



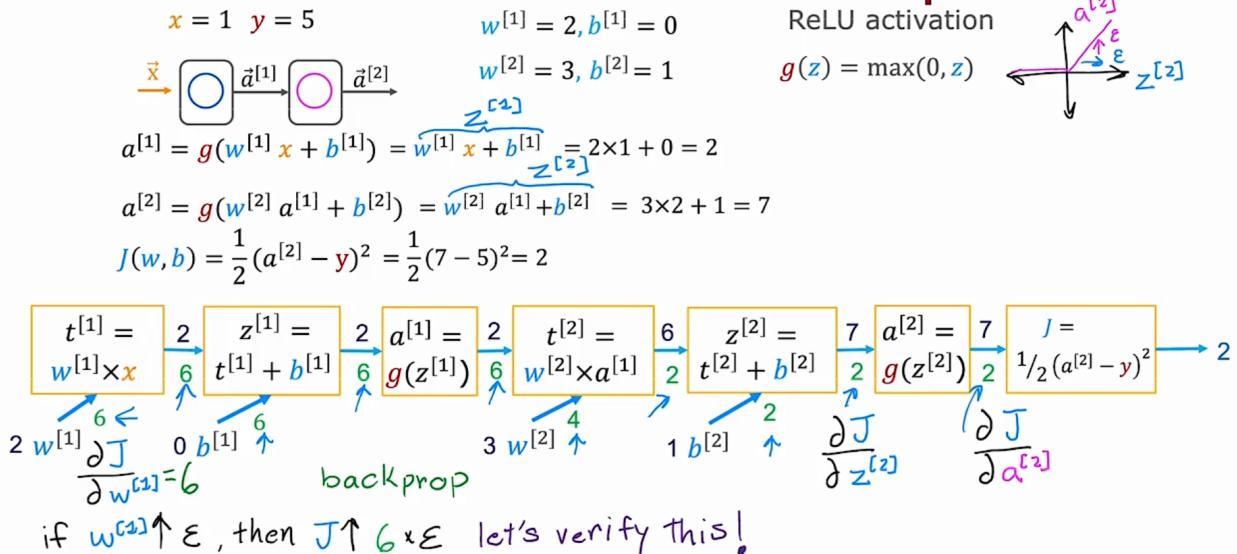
Compute $\frac{\partial J}{\partial a}$ once and use it to compute both $\frac{\partial J}{\partial w}$ and $\frac{\partial J}{\partial b}$.

If N nodes and P parameters, compute derivatives
in roughly $N + P$ steps rather than $N \times P$ steps.

N	P	N + P	N × P
10,000	100,000	1.1×10^5	10^9

DeepLearning.AI Stanford ONLINE Andrew Ng
the computation graph gives you a very efficient way to compute all the derivatives, and that's

Neural Network Example



阶段

确定模型 (d) 确定参数 (w, b) 测试误差

Training/cross validation/test set

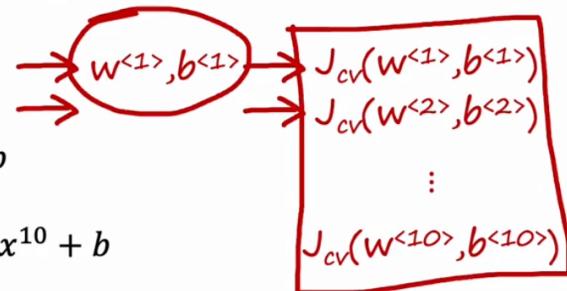
Training error: $J_{train}(\vec{w}, b) = \frac{1}{2m_{train}} \left[\sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 \right]$

Cross validation error: $J_{cv}(\vec{w}, b) = \frac{1}{2m_{cv}} \left[\sum_{i=1}^{m_{cv}} (f_{\vec{w}, b}(\vec{x}_{cv}^{(i)}) - y_{cv}^{(i)})^2 \right]$ (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. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + w_3 x^3 + b$
- : :
- $d=10$ 10. $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x^2 + \dots + w_{10} x^{10} + b$



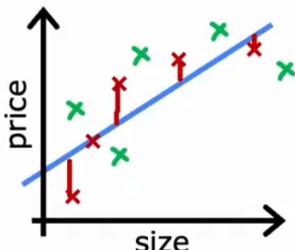
→ Pick $w_1 x_1 + \dots + w_4 x^4 + b$ ($J_{cv}(w^{<4>}, b^{<4>})$)

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

所以到目前为止，你还没有把任何参数 w b d 放到测试集中

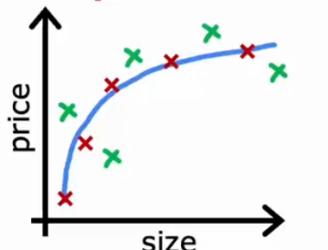
and so up until this point, you've not fit any parameters, either w or b or d to the test set

Bias/variance



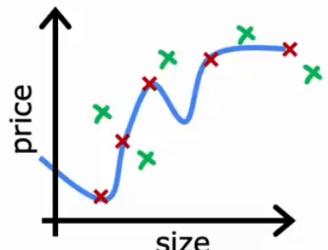
→ High bias
(underfit)

J_{train} is high
 J_{cv} is high



"Just right"

J_{train} is low
 J_{cv} is low

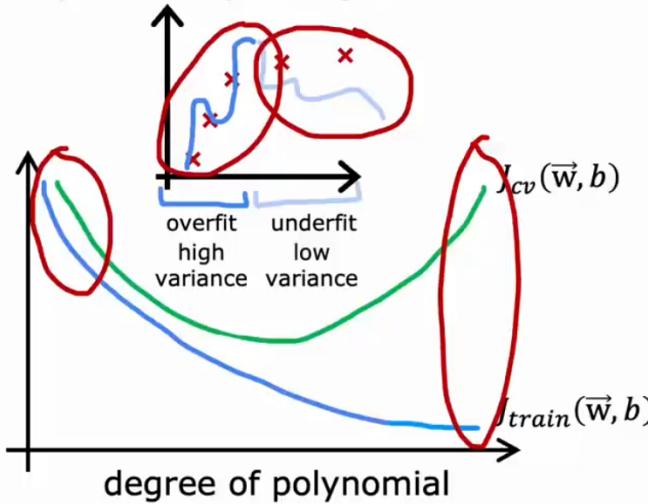


High variance
(overfit)

J_{train} is low
 J_{cv} is high

Diagnosing bias and variance

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



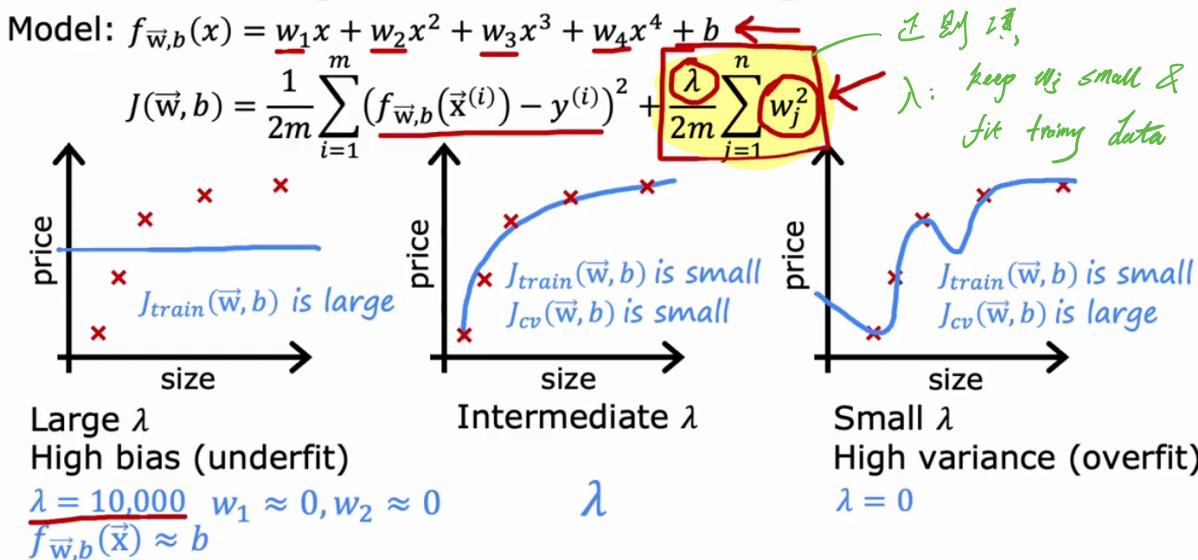
- High bias (underfit)
 J_{train} will be high
(J_{train} \approx J_{cv})
- High variance (overfit)
 J_{cv} \gg J_{train}
(J_{train} may be low)
- High bias and high variance
and J_{cv} \gg J_{train}

High Bias : not even well in Training Set

High Variance : not good in CVs, but may
全连接层 (dense layer) 类型，这一层中的每个神经元都从前一层得到所有的激活。

good in TS

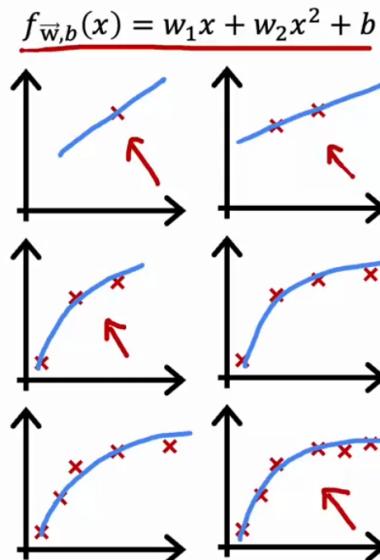
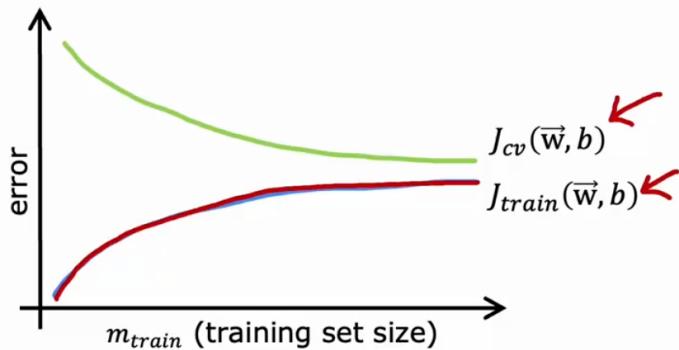
Linear regression with regularization



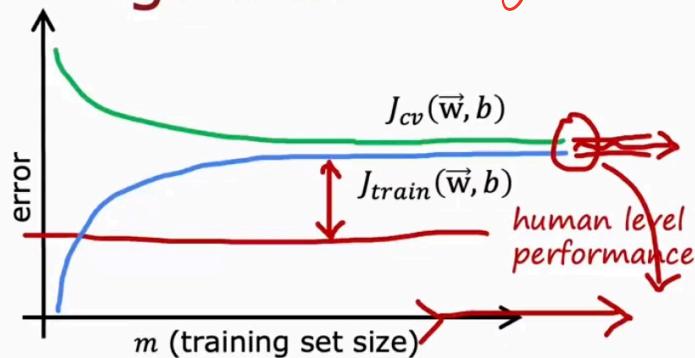
Learning curves

J_{train} = training error

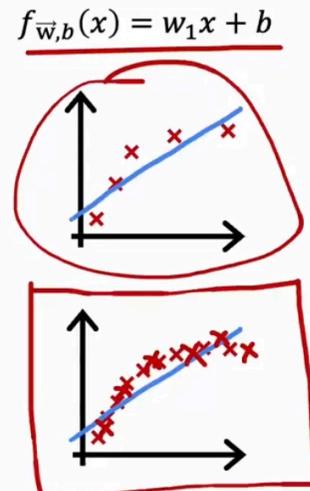
J_{cv} = cross validation error



High bias *underfit*



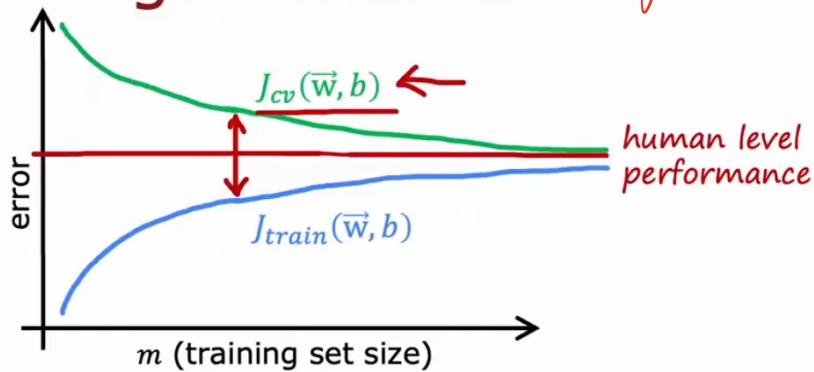
if a learning algorithm suffers from high bias, getting more training data will not (by itself) help much.



如果一个学习算法有很高的偏差，获得更多的训练数据本身不会有
多大帮助。

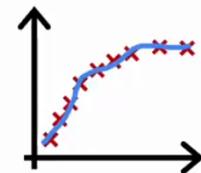
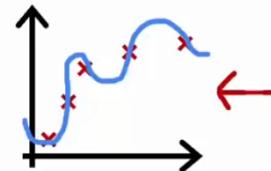
that if a learning algorithm has high bias, getting more training data will not
by itself help that much.

High variance *overfit*



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

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



Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

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

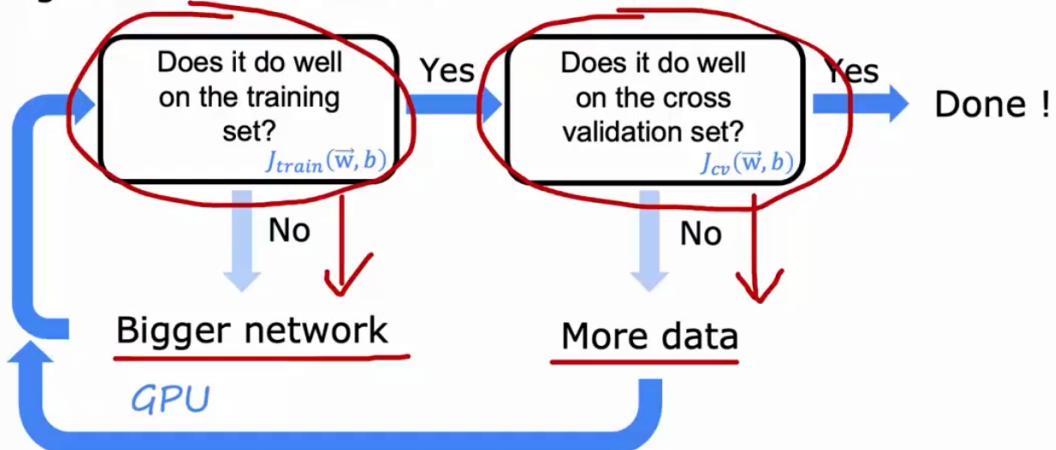
But it makes unacceptably large errors in predictions. What do you try next?

- Get more training examples
- Try smaller sets of features $x, x^2, \cancel{x}, \cancel{x^2}, \cancel{...}$
- Try getting additional features \leftarrow
- Try adding polynomial features $(x_1^2, x_2^2, x_1 x_2, \text{etc})$
- Try decreasing $\lambda \leftarrow$
- Try increasing $\lambda \leftarrow$

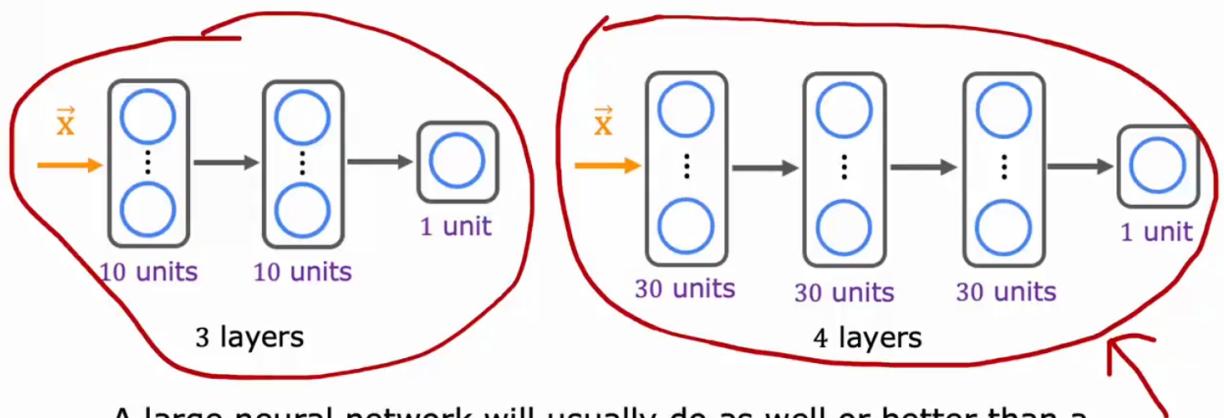
- fixes high variance
- fixes high variance
- fixes high bias
- fixes high bias
- fixes high bias
- fixes high variance

Neural networks and bias variance

Large neural networks are low bias machines



Neural networks and regularization



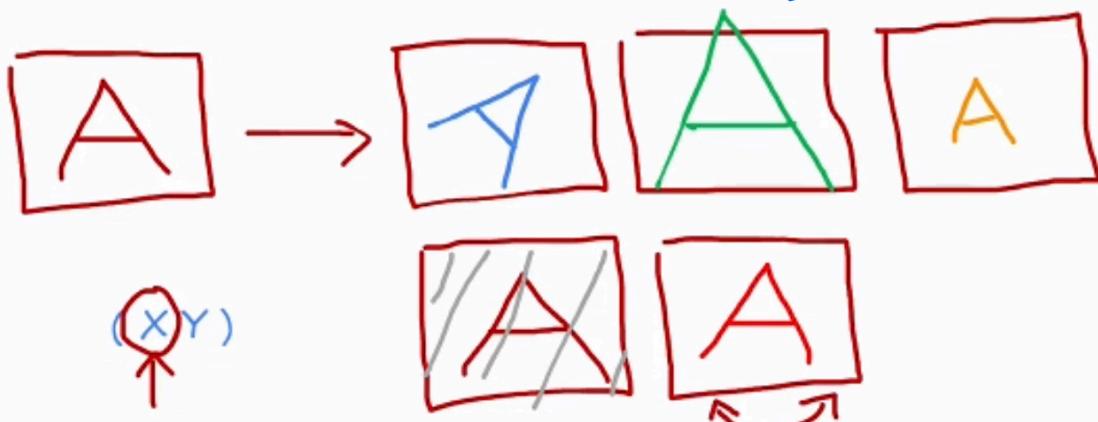
A large neural network will usually do as well or better than a smaller one so long as regularization is chosen appropriately.

数据增强

Data augmentation *To modify examples*

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

Increases dataset manually



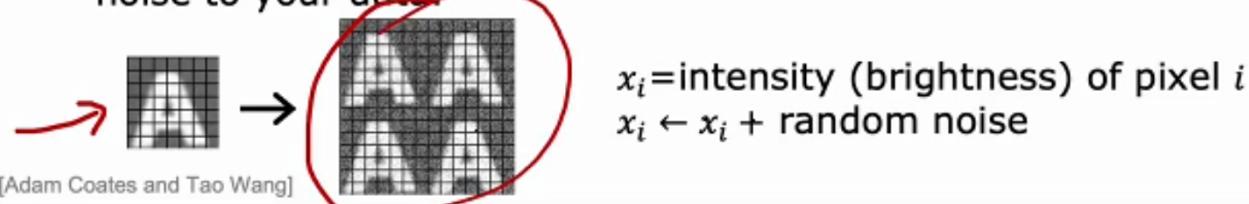
Data augmentation by introducing distortions

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



Audio:
Background noise,
bad cellphone connection

Usually does not help to add purely random/meaningless noise to your data.

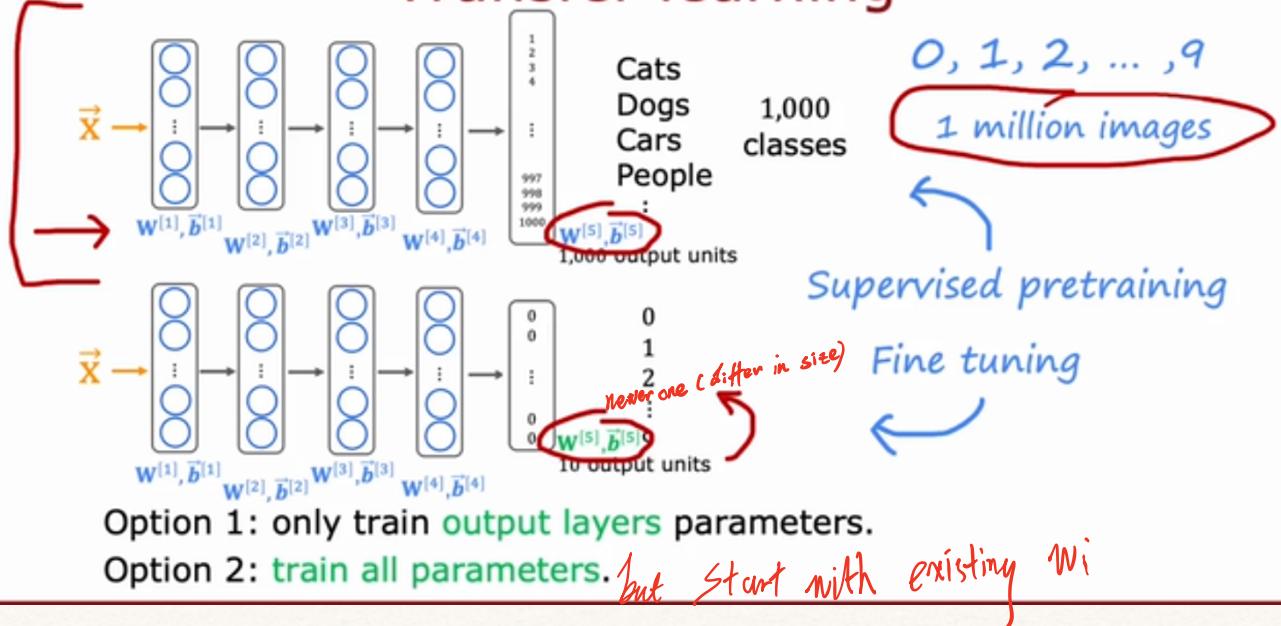


x_i = intensity (brightness) of pixel i
 $x_i \leftarrow x_i + \text{random noise}$

数据合成 To create examples Data synthesis

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

Transfer learning



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). IM
- 2. Further train (fine tune) the network on your own data. Pre-training

1000

50

Pre-training

Precision/recall

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

		Actual Class	
		1	0
Predicted 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:

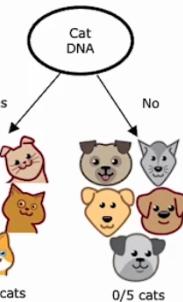
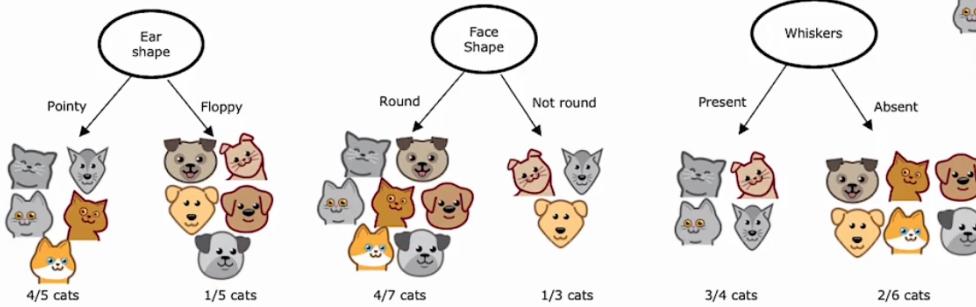
(of all patients that actually have the rare disease, what fraction did we correctly detect as having it?)

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

Decision Tree Learning

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

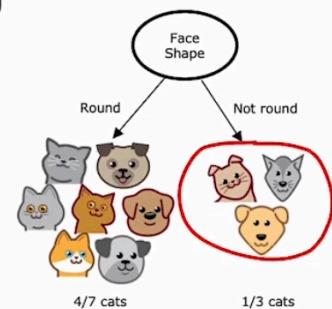
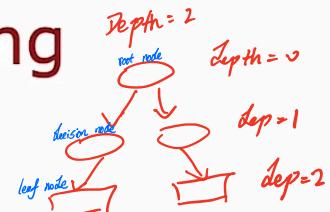
Maximize purity (or minimize impurity)



Decision Tree Learning

Decision 2: When do you stop splitting?

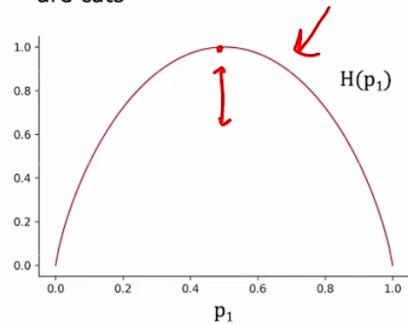
- When a node is 100% one class
- When splitting a node will result in the tree exceeding a maximum depth
- When improvements in purity score are below a threshold
- When number of examples in a node is below a threshold



Entropy as a measure of impurity

impurity $\uparrow \rightarrow$ Entropy \uparrow

p_1 = fraction of examples that are cats



$$p_0 = 1 - p_1$$

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

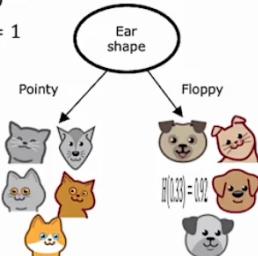
Note: "0 log(0)" = 0

Goal : Reduce entropy \rightarrow Raise purity

Choosing a split

$$p_1 = 5/10 = 0.5$$

$$H(0.5) = 1$$



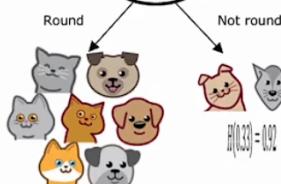
$$p_1 = 4/5 = 0.8 \quad p_1 = 1/5 = 0.2$$

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

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

$$= 0.28$$

$$H(0.5) = 1$$

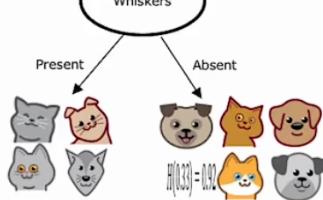


$$p_1 = 4/7 = 0.57 \quad p_1 = 1/3 = 0.33$$

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

$$= 0.03$$

$$H(0.5) = 1$$



$$p_1 = 3/4 = 0.75 \quad p_1 = 2/6 = 0.33$$

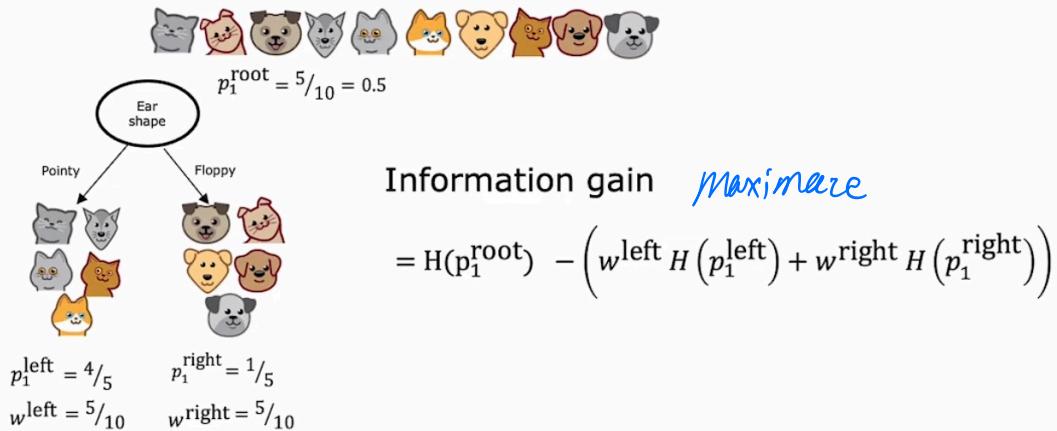
$$H(0.75) = 0.81 \quad H(0.33) = 0.92$$

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

$$= 0.12$$

Information gain : 信息熵减少量

Information Gain



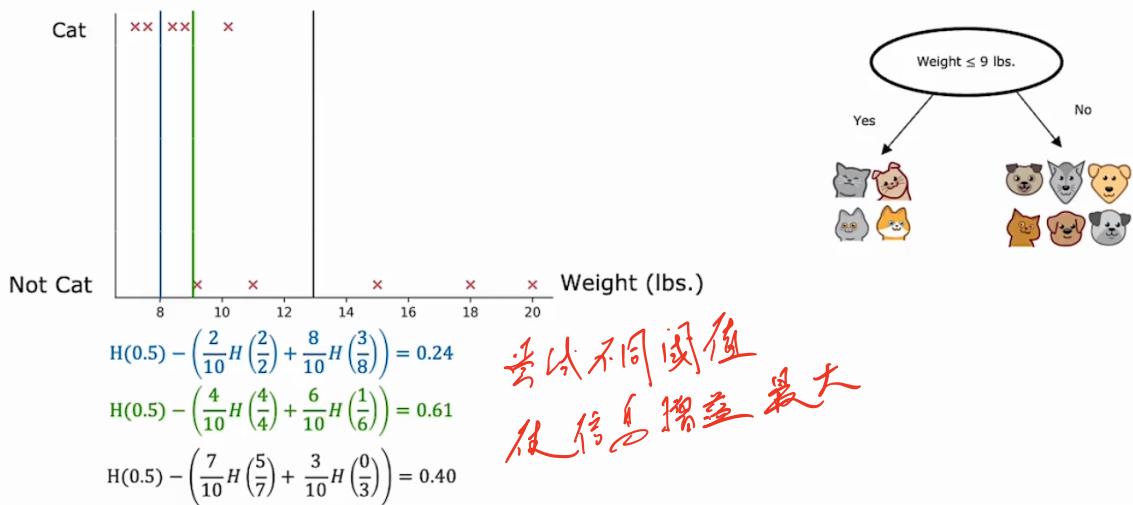
One hot encoding

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

And because one of these features will always take on the value 1
that's

因为这些特征中总有一个取1，也就是hot (热) 特征

Splitting on a continuous variable



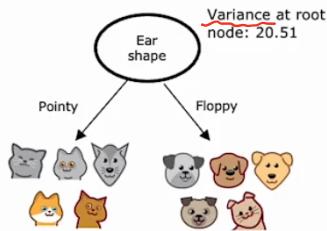
↑
Decision tree for classification

for regression

↓

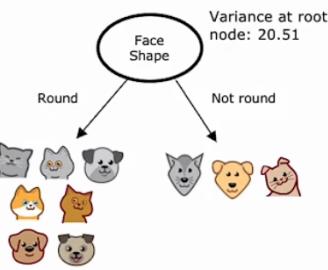
Choosing a split

predict the weight



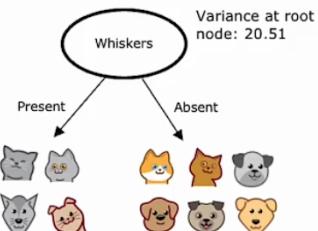
Weights: 7.2,
9.2, 8.4, 7.6, 10.2
Variance: 1.47
 $w^{\text{left}} = \frac{5}{10}$

$$20.51 - \left(\frac{5}{10} * 1.47 + \frac{5}{10} * 21.87 \right) \\ = 8.84 \quad \leftarrow$$



Weights: 8.8, 15,
11, 18, 20
Variance: 21.87
 $w^{\text{left}} = \frac{7}{10}$

$$20.51 - \left(\frac{7}{10} * 27.80 + \frac{3}{10} * 1.37 \right) \\ = 0.64$$



Weights: 7.2, 8.8,
9.2, 8.4
Variance: 0.75
 $w^{\text{left}} = \frac{4}{10}$

$$20.51 - \left(\frac{4}{10} * 0.75 + \frac{6}{10} * 23.32 \right) \\ = 6.22$$

Tree Ensemble

Trees are highly sensitive to small changes of the data



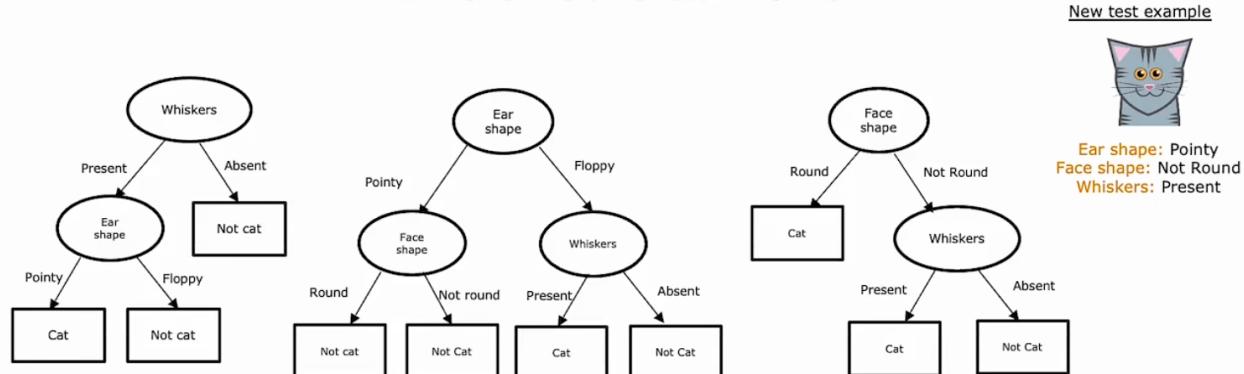
DeepLearning.AI Stanford | ONLINE

Andrew Ng

我们仅改变了单个训练样本，具有最高信息增益的节点特征就变成了胡须，而不是耳朵形状

With just changing a single training example, the highest information gain feature at the split on becomes the whiskers feature instead of the ear shape feature.

Tree ensemble



DeepLearning.AI Stanford | ONLINE

Andrew Ng

这就是我们所说的集成树，意思是多棵树的集合

And this is what we call a tree ensemble, which just means a collection of multiple trees.

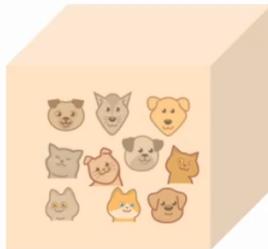
如果你想给一个新的测试样本分类，那么你要做的就是在新样本上运行这三棵树，并让它们投票做出最终预测

If you had a new test example that you wanted to classify, then what you would do is run all three of these trees on your new example and get them to vote on what is the final prediction.

For Building ensemble tree ↓

有放回 抽样

Sampling with replacement



	Ear shape	Face shape	Whiskers	Cat
1	Pointy	Round	Present	1
2	Floppy	Not round	Absent	0
3	Pointy	Round	Absent	1
4	Pointy	Not round	Present	0
5	Floppy	Not round	Absent	0
6	Pointy	Round	Absent	1
7	Pointy	Round	Present	1
8	Floppy	Not round	Present	1
9	Floppy	Round	Absent	0
10	Pointy	Round	Absent	1

其中一些结果是重复的，并且你还注意到此训练集并不包含所有的原始训练样本

some of which are repeats and you notice also that this training set does not contain all ten of the original training examples.

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

Boosted trees intuition

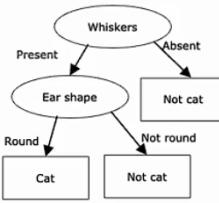
Given training set of size m

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	Absent	No
Pointy	Not Round	Present	Yes



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)
```

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

