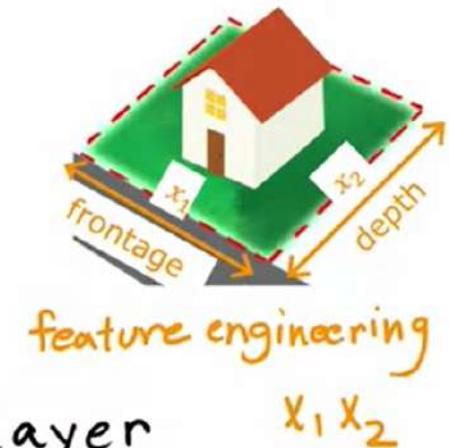


# Demand Prediction

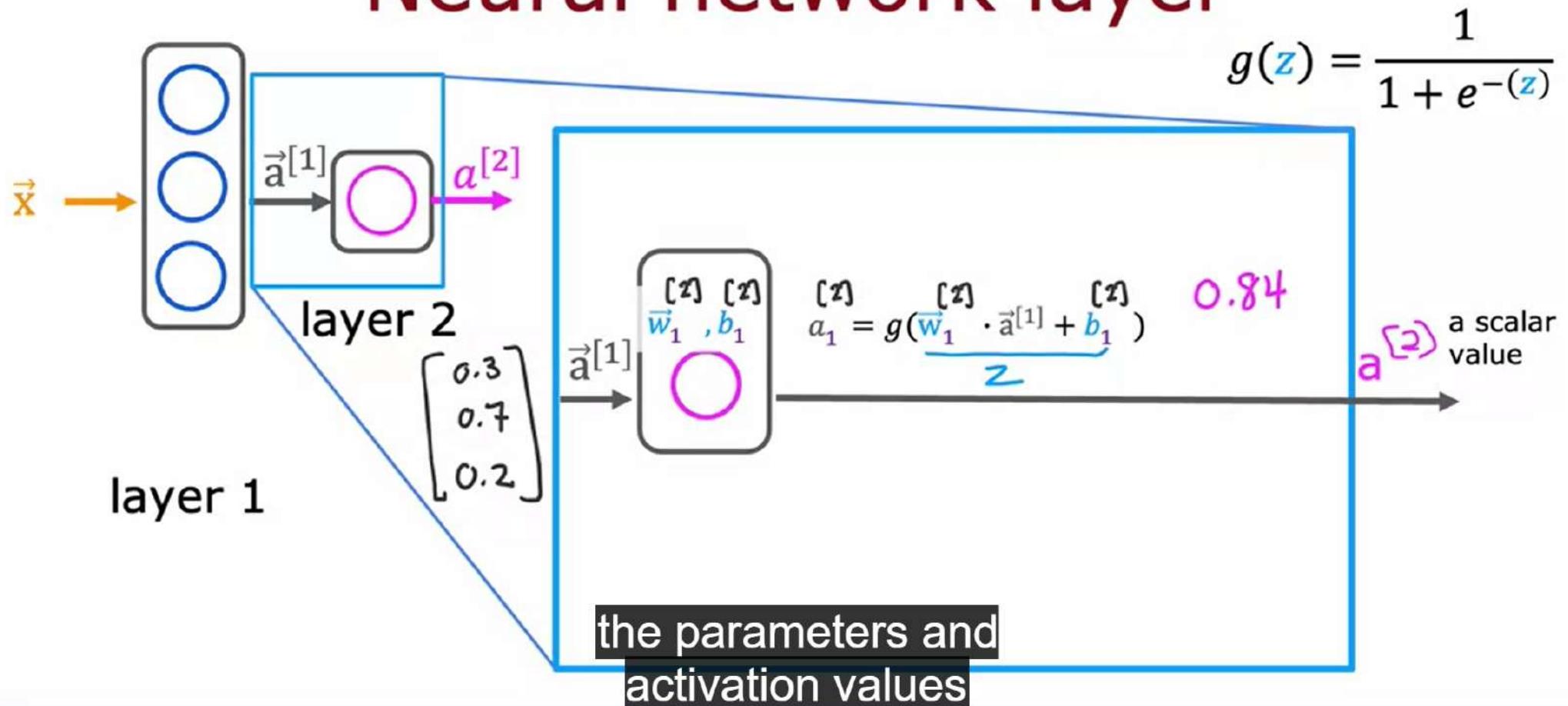


# Demand Prediction

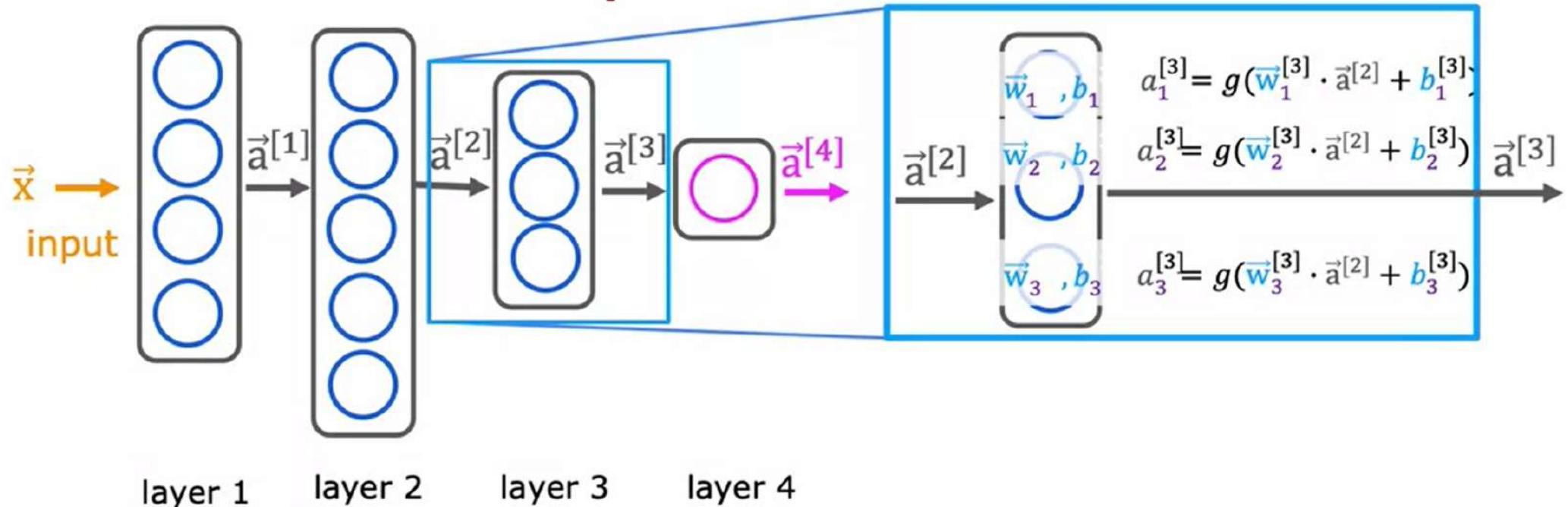
$\vec{x} \rightarrow (x, y)$  "hidden layer"  
input layer  
layer can have multiple neurons



# Neural network layer

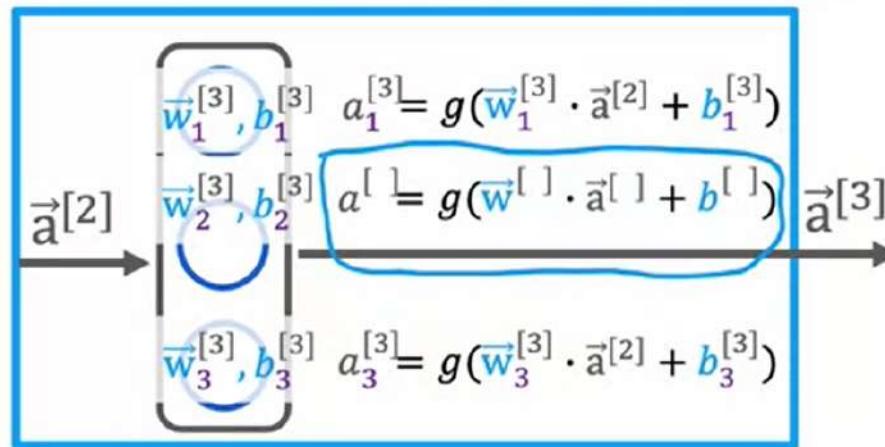


# More complex neural network



which is another vector.

# [In Video Quiz]



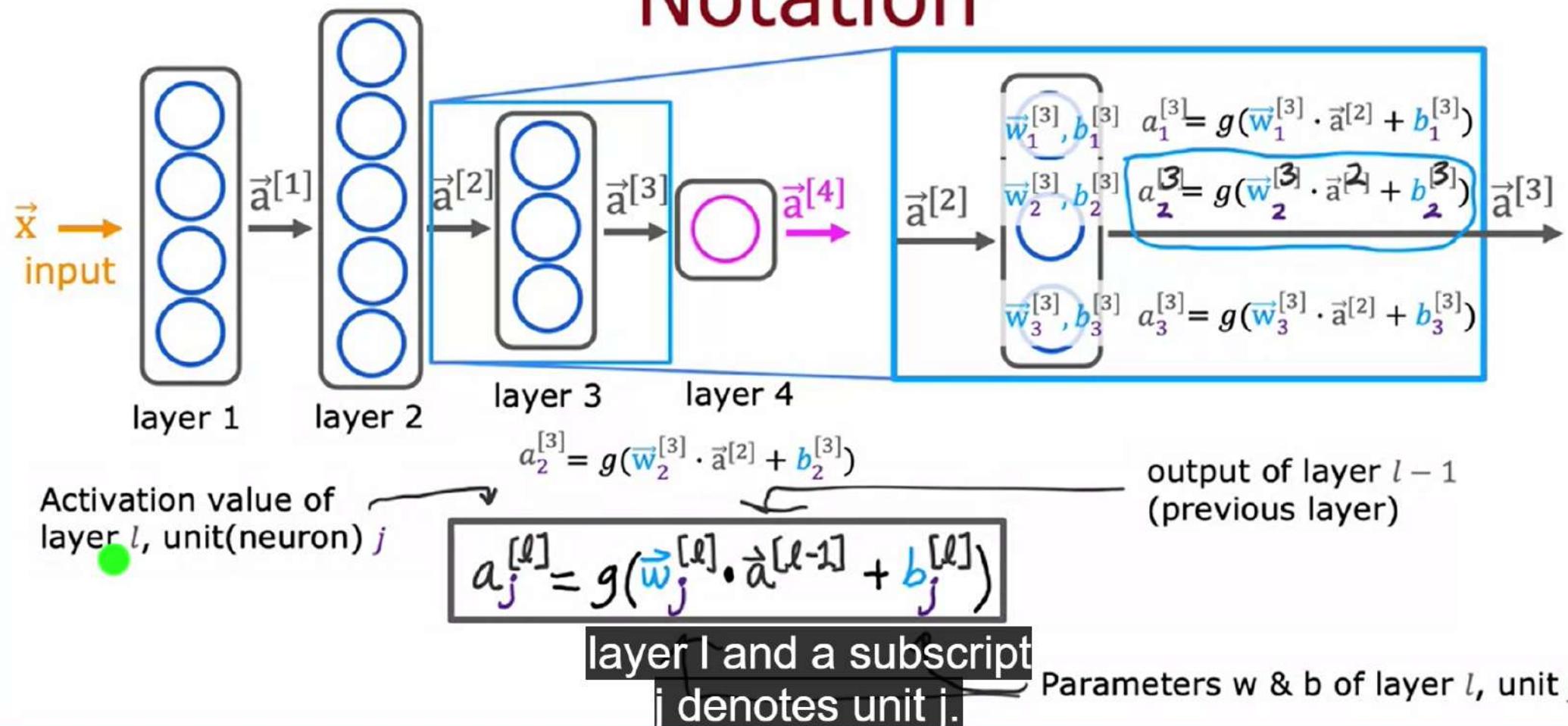
Can you fill in the superscripts and subscripts for the second neuron?

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

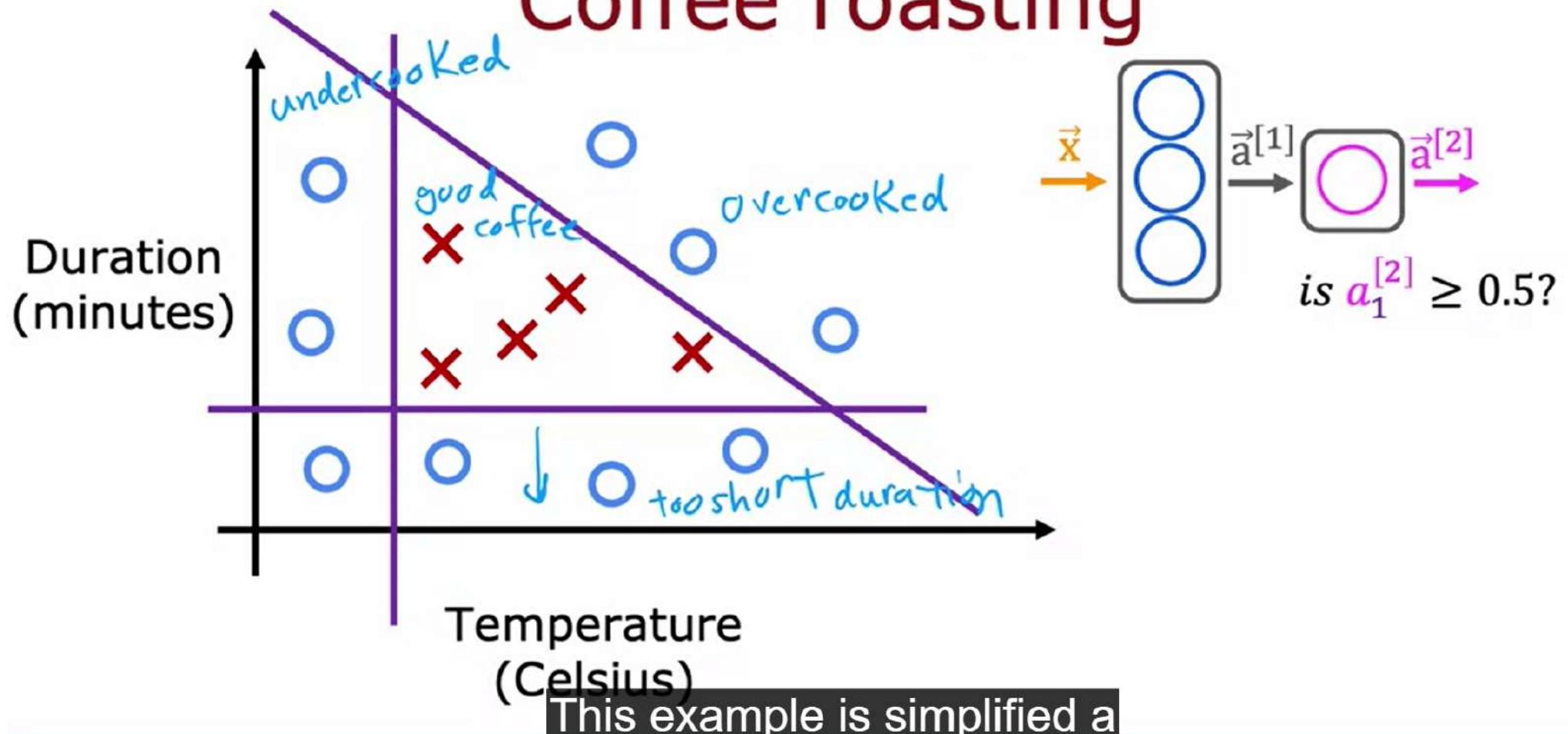
$a_2^{[3]} = g(\vec{w}_2^{[3]} \cdot \vec{a}^{[2]} + b_2^{[3]})$   
 $a_2^{[3]} = g(\vec{w}_2^{[3]} \cdot \vec{a}^{[2]} + b_2^{[3]})$

The activation of the 2nd neuron at layer 3 is denoted

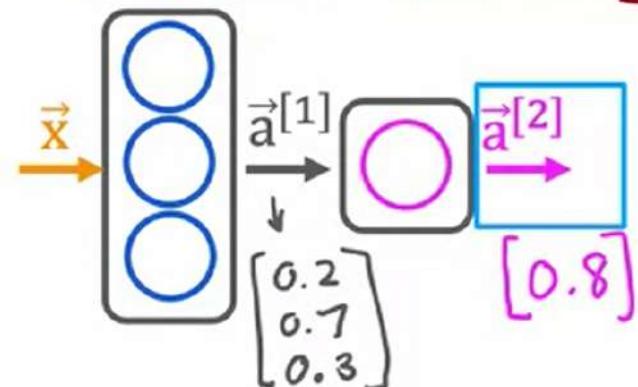
# Notation



# Coffee roasting



# Build the model using TensorFlow



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

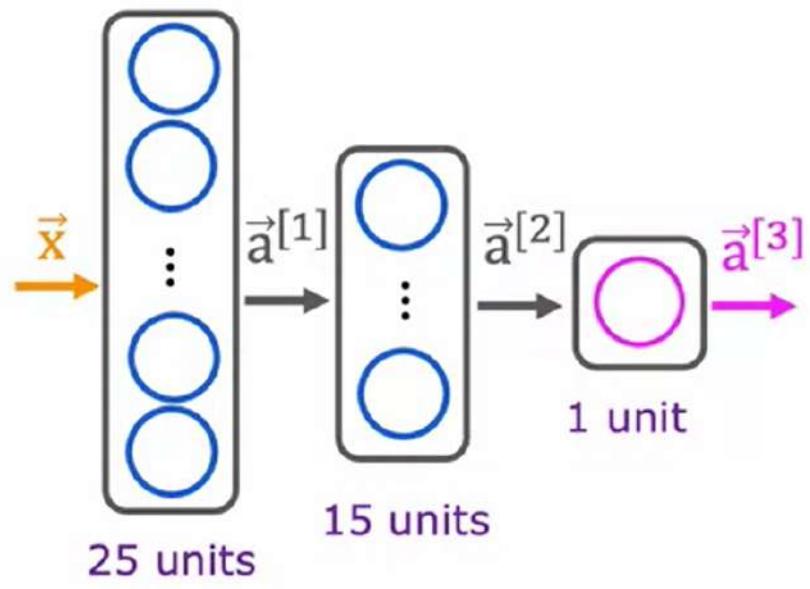
```
layer_2 = Dense(units=1, activation='sigmoid')
a2 = layer_2(a1)
```

But these are the key  
steps for forward propagation in

is  $a_1^{[2]} \geq 0.5$ ?  
yes  $\hat{y} = 1$  no  $\hat{y} = 0$

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

# Digit classification model



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

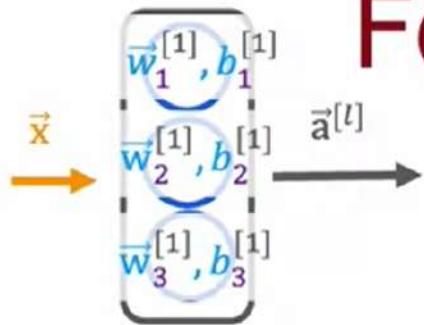
x = np.array([[0..., 245, ..., 17],
              [0..., 200, ..., 184]])
y = np.array([1,0])

model.fit(x,y)

model.predict(x_new)
```

you into a new network and same as before.

# Forward prop in NumPy



$$\vec{w}_1^{[1]} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \vec{w}_2^{[1]} = \begin{bmatrix} -3 \\ 4 \end{bmatrix} \quad \vec{w}_3^{[1]} = \begin{bmatrix} 5 \\ -6 \end{bmatrix}$$

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

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

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

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

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

and so will pull out  $w_{1,1}$ .

# For loops vs. vectorization

```
x = np.array([200, 17])  
W = np.array([[1, -3, 5],  
             [-2, 4, -6]])  
b = np.array([-1, 1, 2])  
  
def dense(a_in,W,b):  
    units = W.shape[1]  
    a_out = np.zeros(units)  
    for j in range(units):  
        w = W[:,j]  
        z = np.dot(w, a_in) + b[j]  
        a_out[j] = g(z)  
    return a_out
```

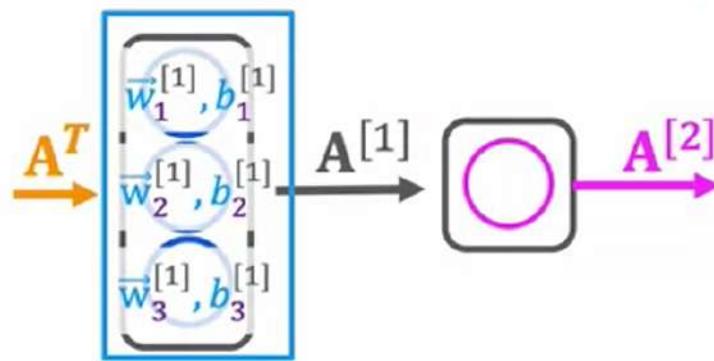
[1,0,1] ↘

vectorized

```
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]]
```

one step of forward propagation

# Dense layer vectorized



$$A^T = [200 \ 17]$$
$$W = \begin{bmatrix} 1 & -3 & 5 \\ -2 & 4 & -6 \end{bmatrix}$$
$$B = [-1 \ 1 \ 2]$$

$\text{1 } \times \text{ 2}$   
 $\text{2 } \times \text{ 3}$   
 $\text{1 } \times \text{ 3}$

$$Z = A^T W + B$$
$$\begin{bmatrix} 165 \\ z_1^{[1]} \end{bmatrix} \begin{bmatrix} -531 \\ z_2^{[1]} \end{bmatrix} \begin{bmatrix} 900 \\ z_3^{[1]} \end{bmatrix}$$

$$A = g(Z)$$
$$\begin{bmatrix} 1 & 0 & 1 \end{bmatrix}$$

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

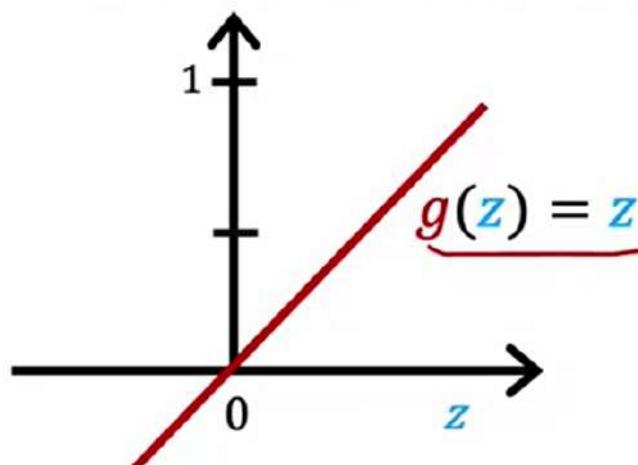
```
def dense(AT,W,b):
    z = np.matmul(AT,W) + b
    a_out = g(z)
```

the activation function applied

# Examples of Activation Functions

"No activation function"

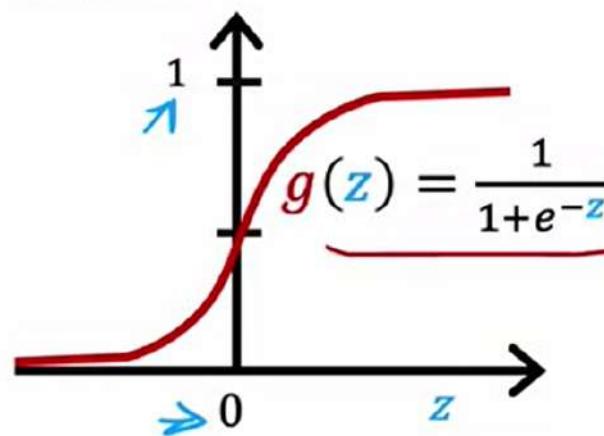
Linear activation function



$$a = g(z) = \underbrace{\vec{w} \cdot \vec{x} + b}_{z}$$

$$a_2^{[1]} = g(\underbrace{\vec{w}_2^{[1]} \cdot \vec{x}}_z + b_2^{[1]})$$

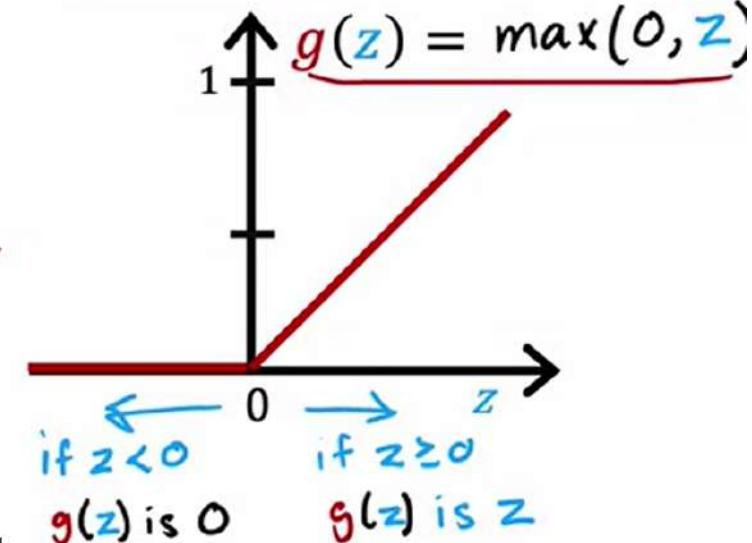
Sigmoid



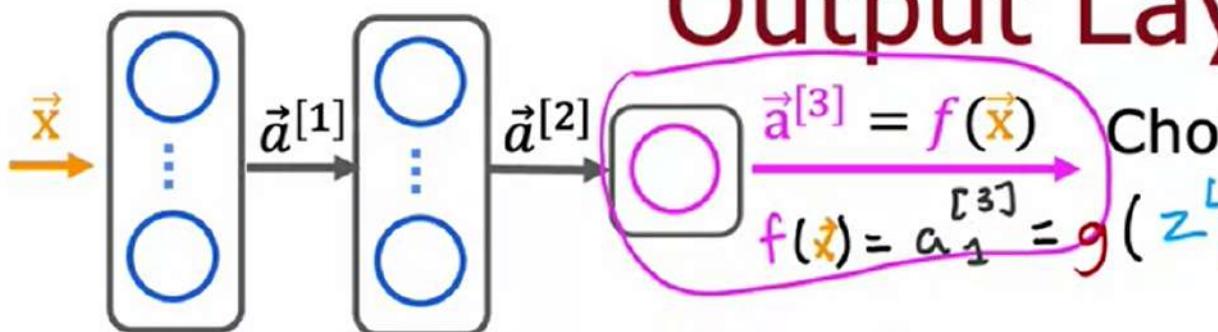
$0 < g(z) < 1$   
It just refers to the linear activation function.

ReLU

Rectified Linear Unit

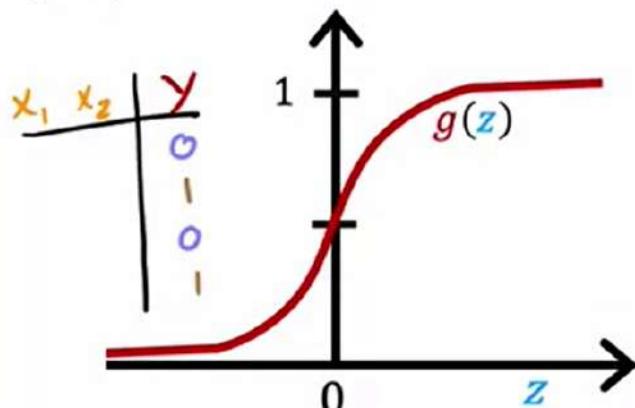


# Output Layer



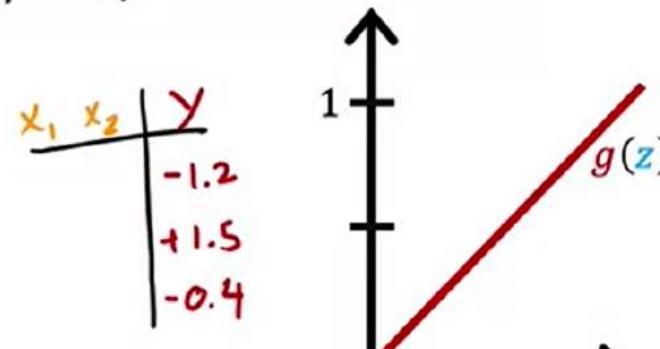
Binary classification

Sigmoid  
 $y=0/1$



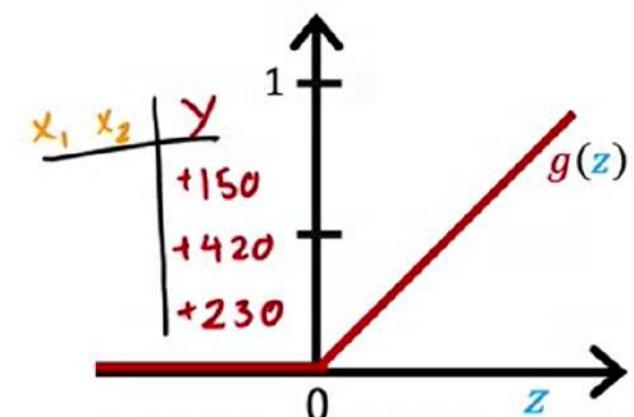
Regression

Linear activation function  
 $y = +/-$



Regression

ReLU



the ReLU activation function  
because as you see here,

Logistic regression  
(2 possible output values)

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

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

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  $w_1, w_2, \dots, w_N$   
 $b_1, b_2, \dots, b_N$

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

△  $z_4 = \vec{w}_4 \cdot \vec{x} + b_4$   
another variable k to index  
the summation because

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

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

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

X O □ △

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

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

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$
$$= P(y=2|\vec{x}) \quad 0.20$$

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

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

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

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

# Cost

## Logistic regression

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

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

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

$$\text{loss} = -y \underbrace{\log a_1}_{\text{if } y=1} - (1-y) \underbrace{\log(1-a_1)}_{\text{if } y=0}$$

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

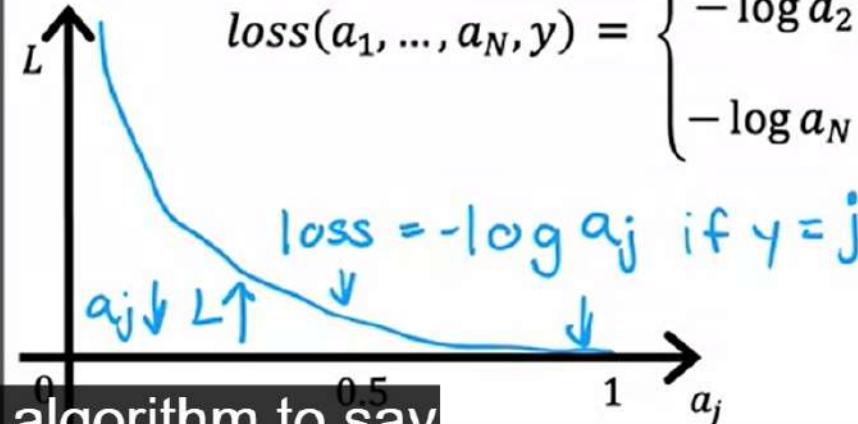
## Softmax regression

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

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

Crossentropy loss

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



you want the algorithm to say

# MNIST with softmax

- 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  $J(\vec{w}, b)$

```
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)
```

Just one important note, if you use this code exactly as I've written here, it will

# Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

Logistic regression:

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

Original loss

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

More accurate loss (in code)

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

either of these implementations  
actually works okay,

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

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

# More numerically accurate implementation of softmax

## Softmax regression

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

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

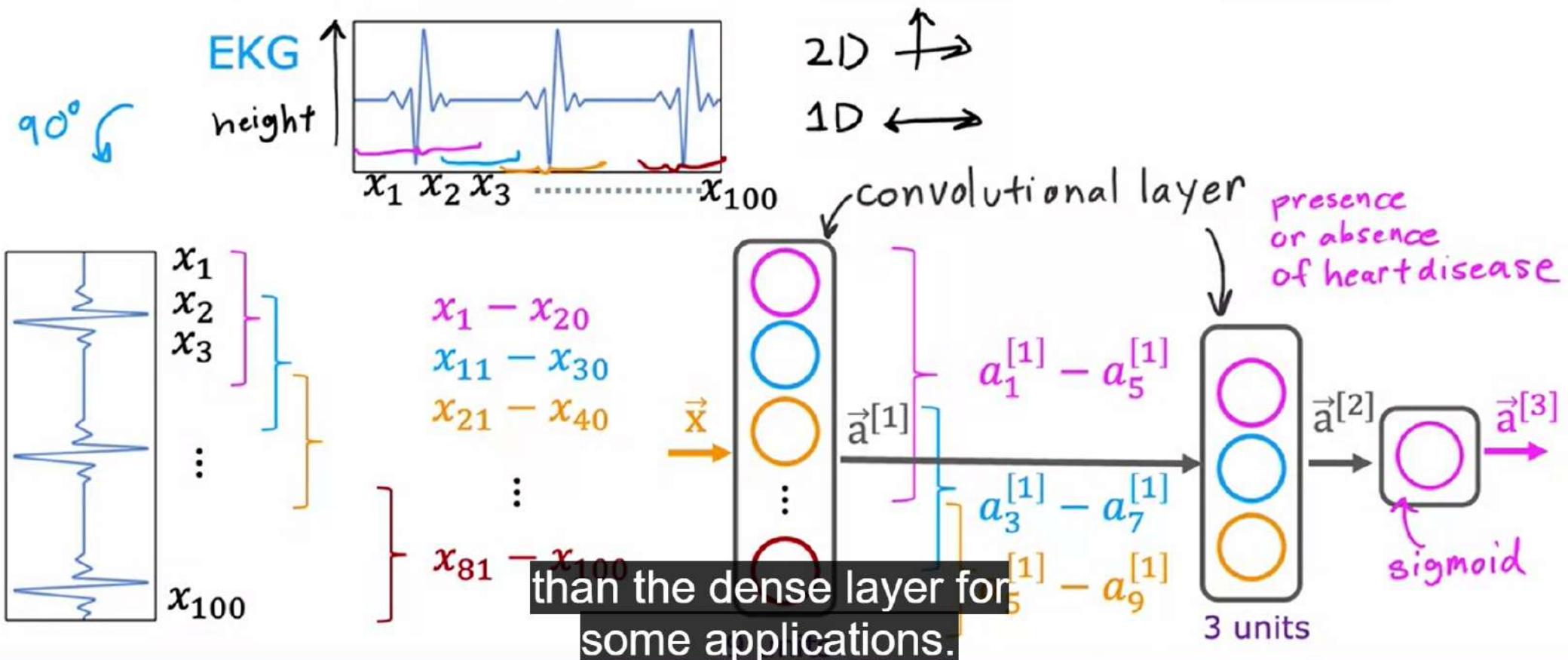
```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
'linear'  
model.compile(loss=SparseCategoricalCrossEntropy())
```

## 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=SparseCategoricalCrossEntropy(from_logits=True))  
and this whole computation of
```

# Convolutional Neural Network



# Derivative Example

Cost function

$$J(w) = w^2$$

Say  $w = 3$

$$J(w) = 3^2 = 9$$

$$\epsilon = 0.002$$

If we increase  $w$  by a tiny amount  $\epsilon = 0.001$  how does  $J(w)$  change?

$$w = 3 + \cancel{0.001} \quad 0.002$$

$$J(w) = w^2 = \cancel{9.006001}$$

$$\begin{array}{r} 9.012004 \\ \hline 9.012 \end{array}$$

$$\text{If } \underline{w} \uparrow \cancel{0.001} \quad \overset{0.002}{\epsilon} \leftarrow$$

$$\begin{aligned} J(w) &\uparrow 6 \times \cancel{0.001} \quad 6 \times \epsilon \quad 6 \times 0.002 = 0.012 \\ \frac{\partial}{\partial w} J(w) &= 6 \end{aligned}$$

That's why the  
derivative of  $J$  of  $w$

# More Derivative Examples

$w = 3$

$J(w) = w^2 = 9$

$w \uparrow 0.001$

$J(w) = J(3.001) = 9.006001$

$\frac{\partial}{\partial w} J(w) = 6$

$J(w) \uparrow 6 \times 0.001$

$w = 2$

$J(w) = w^2 = 4$

$w \uparrow 0.001$

$J(w) = J(2.001) = 4.004001$

$\frac{\partial}{\partial w} J(w) = 4$

$J(w) \uparrow 4 \times 0.001$

$w = -3$

$J(w) = w^2 = 9$

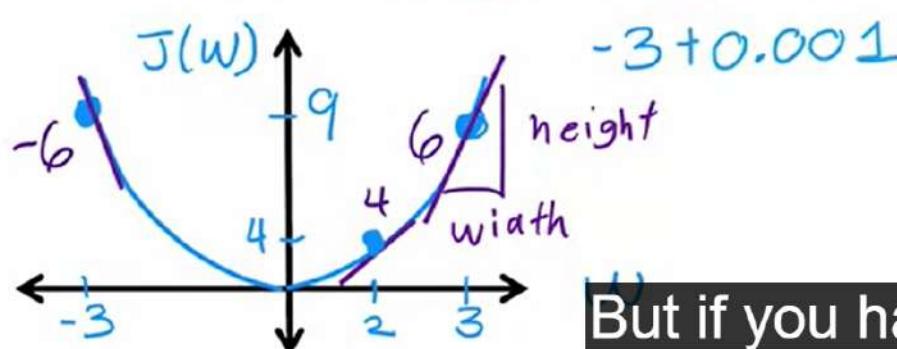
$w \uparrow 0.001$

$J(w) = J(-2.999) = 8.994001$

$\frac{\partial}{\partial w} J(w) = -6$

$J(w) \downarrow 6 \times 0.001$

$J(w) \uparrow -6 \times 0.001$



But if you haven't taken a calculus class before and

# jupyter Using code to get derivatives (unsaved changes)



File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [1]: `import sympy`

In [2]: `J, w = sympy.symbols('J,w')`

In [33]: `J = 1/w`  
J

Out[33]:  $\frac{1}{w}$

In [34]: `dJ_dw = sympy.diff(J,w)`  
`dJ_dw`

Out[34]:  $-\frac{1}{w^2}$

In [35]: `dJ_dw.subs([(w,2)])`

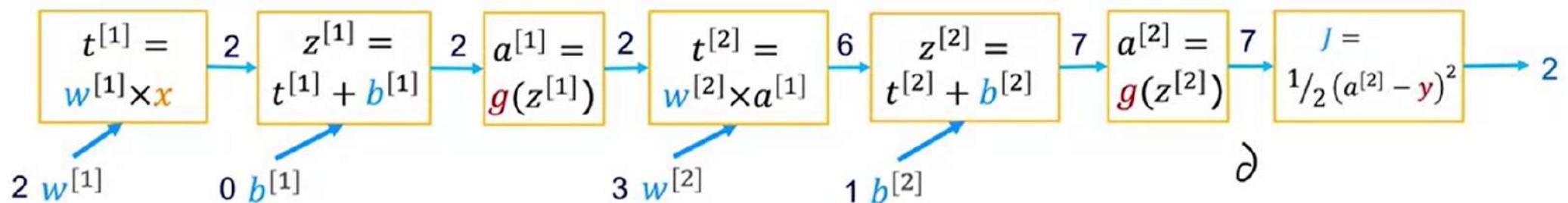
Out[35]:  $-\frac{1}{4}$

Let's copy this back  
to our other slide.

In [ ]:

# Neural Network Example

$$\begin{array}{l} \textcolor{brown}{x} = 1 \quad \textcolor{red}{y} = 5 \\ \textcolor{brown}{x} \rightarrow \text{ReLU} \rightarrow \text{ReLU} \rightarrow \text{ReLU} \\ \textcolor{blue}{w}^{[1]} = 2, \textcolor{blue}{b}^{[1]} = 0 \quad \text{ReLU activation} \\ \textcolor{blue}{w}^{[2]} = 3, \textcolor{blue}{b}^{[2]} = 1 \quad g(z) = \max(0, z) \\ a^{[1]} = g(w^{[1]} x + b^{[1]}) = \underbrace{w^{[1]} x}_{z^{[1]}} + \underbrace{b^{[1]}}_{z^{[2]}} = 2 \times 1 + 0 = 2 \\ a^{[2]} = g(w^{[2]} a^{[1]} + b^{[2]}) = \underbrace{w^{[2]} a^{[1]}}_{z^{[1]}} + \underbrace{b^{[2]}}_{z^{[2]}} = 3 \times 2 + 1 = 7 \\ J(w, b) = \frac{1}{2}(a^{[2]} - y)^2 = \frac{1}{2}(7 - 5)^2 = 2 \end{array}$$



But if you were to carry out backprop,  
the first thing you do is ask,

# Model selection (choosing a model)

$d=1$

$$1. f_{\vec{w}, b}(\vec{x}) = w_1 x + b$$

$$\rightarrow w^{<1>} , b^{<1>} \rightarrow J_{test}(w^{<1>} , b^{<1>})$$

$d=2$

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

$$\rightarrow w^{<2>} , b^{<2>} \rightarrow J_{test}(w^{<2>} , b^{<2>})$$

$d=3$

$$3. f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + w_3 x^3 + b \rightarrow w^{<3>} , b^{<3>} \rightarrow J_{test}(w^{<3>} , b^{<3>})$$

$\vdots$

$d=10$

$$10. f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + \cdots + w_{10} x^{10} + b \rightarrow J_{test}(w^{<10>} , b^{<10>})$$

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

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

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

$w, b$  are overly optimistic estimate of generalization error on training data.  
of the generalization error.

# Model selection

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

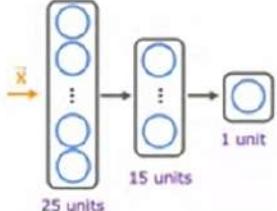
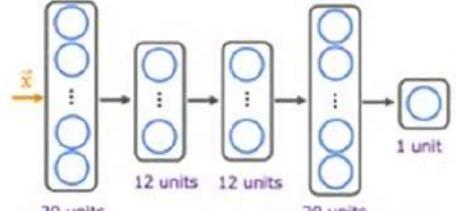
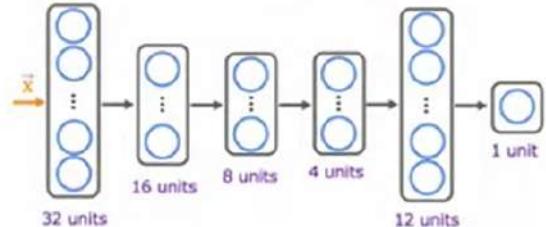
→ Pick  $w_1 x + \dots + w_4 x^4 + b$   $(J_{cv}(w^{(4)}, b^{(4)}))$

Estimate generalization error using test set:  
automatically make a  
decision like what order

$J_{test}(w^{(4)}, b^{(4)})$



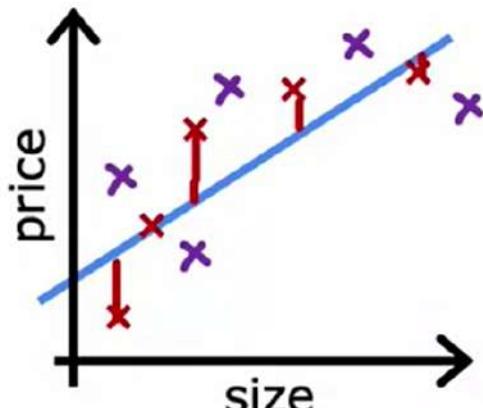
# Model selection – choosing a neural network architecture

1.  $w^{(1)}, b^{(1)}$   $J_{cv}(w^{(1)}, b^{(1)})$
2.  $w^{(2)}, b^{(2)}$   $J_{cv}(w^{(2)}, b^{(2)})$
3.  $w^{(3)}, b^{(3)}$   $J_{cv}(w^{(3)}, b^{(3)})$

Pick  $w^{(2)}, b^{(2)}$

Estimate generalization error using the test set:  $J_{test}(w^{(2)}, b^{(2)})$   
You would compute this using

# Bias/variance

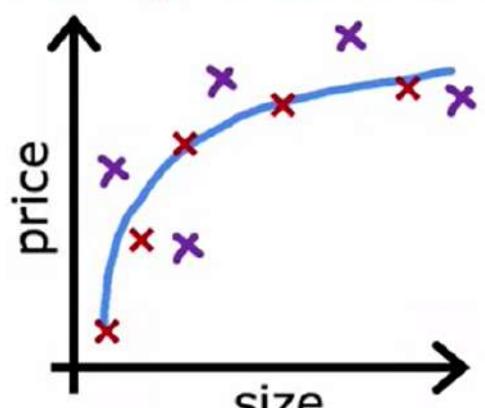


$$f_{\vec{w},b}(x) = w_1x + b$$



High bias  
(underfit)

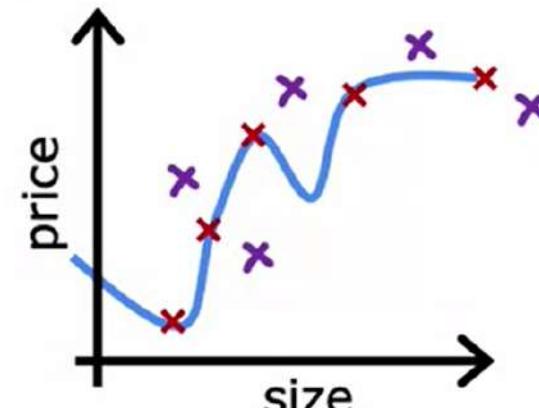
$d=1$



$$f_{\vec{w},b}(x) = w_1x + w_2x^2 + b$$

"Just right"

but  $J_{cv}$  is high.

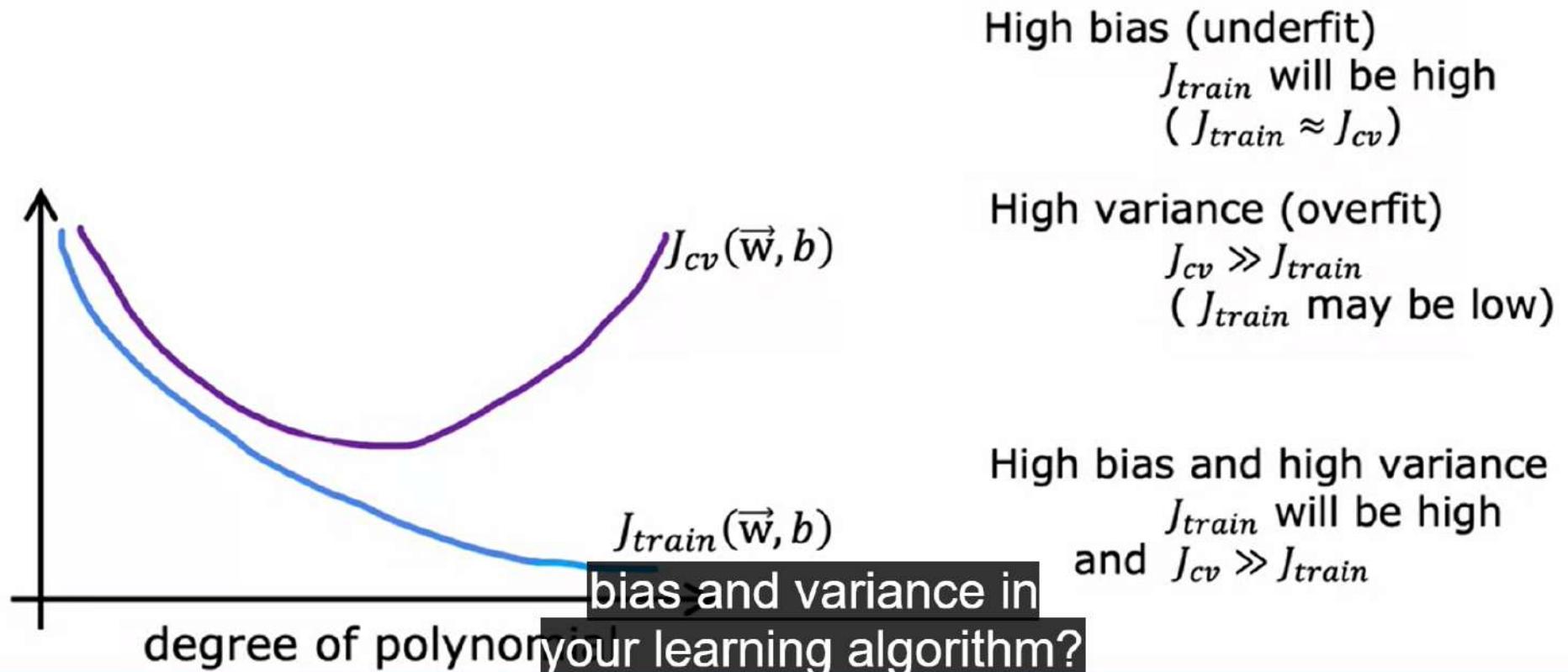


$$f_{\vec{w},b}(x) = w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$$

High variance  
(overfit)

# Diagnosing bias and variance

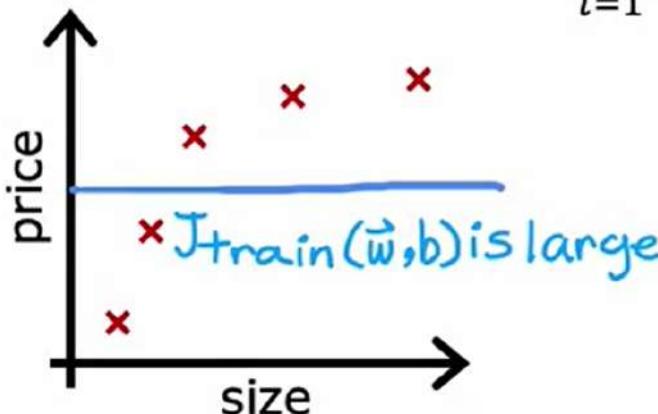
How do you tell if your algorithm has a bias or variance problem?



# Linear regression with regularization

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

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

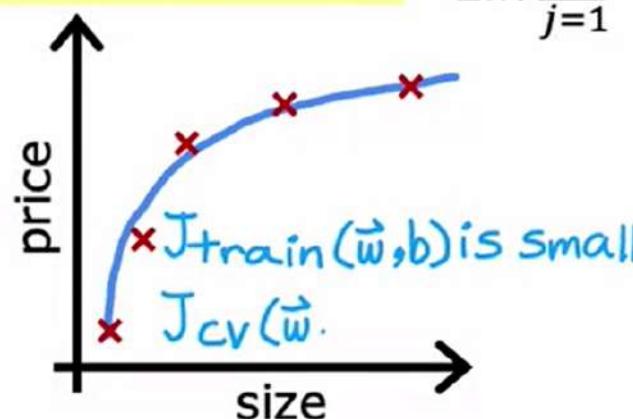


Large  $\lambda$

High bias (underfit)

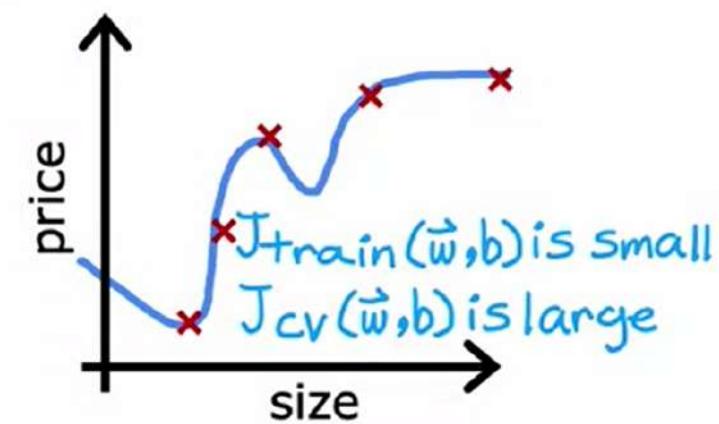
$$\lambda = 10,000 \quad w_1 \approx 0, w_2 \approx 0$$

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



Intermediate  $\lambda$

with small  $J_{\text{train}}$   
and small  $J_{\text{cv}}$ .



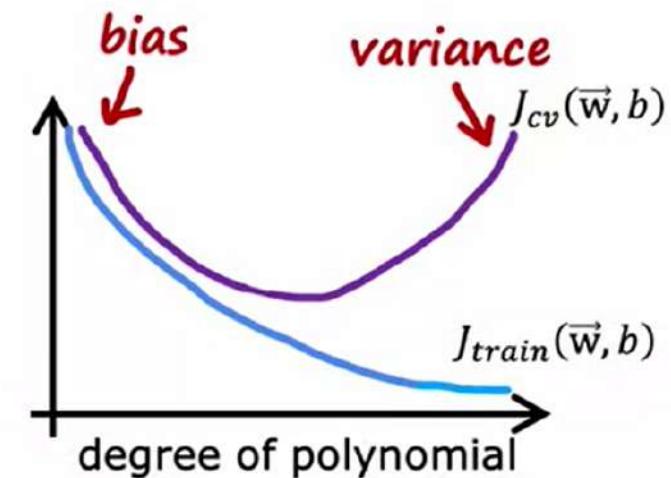
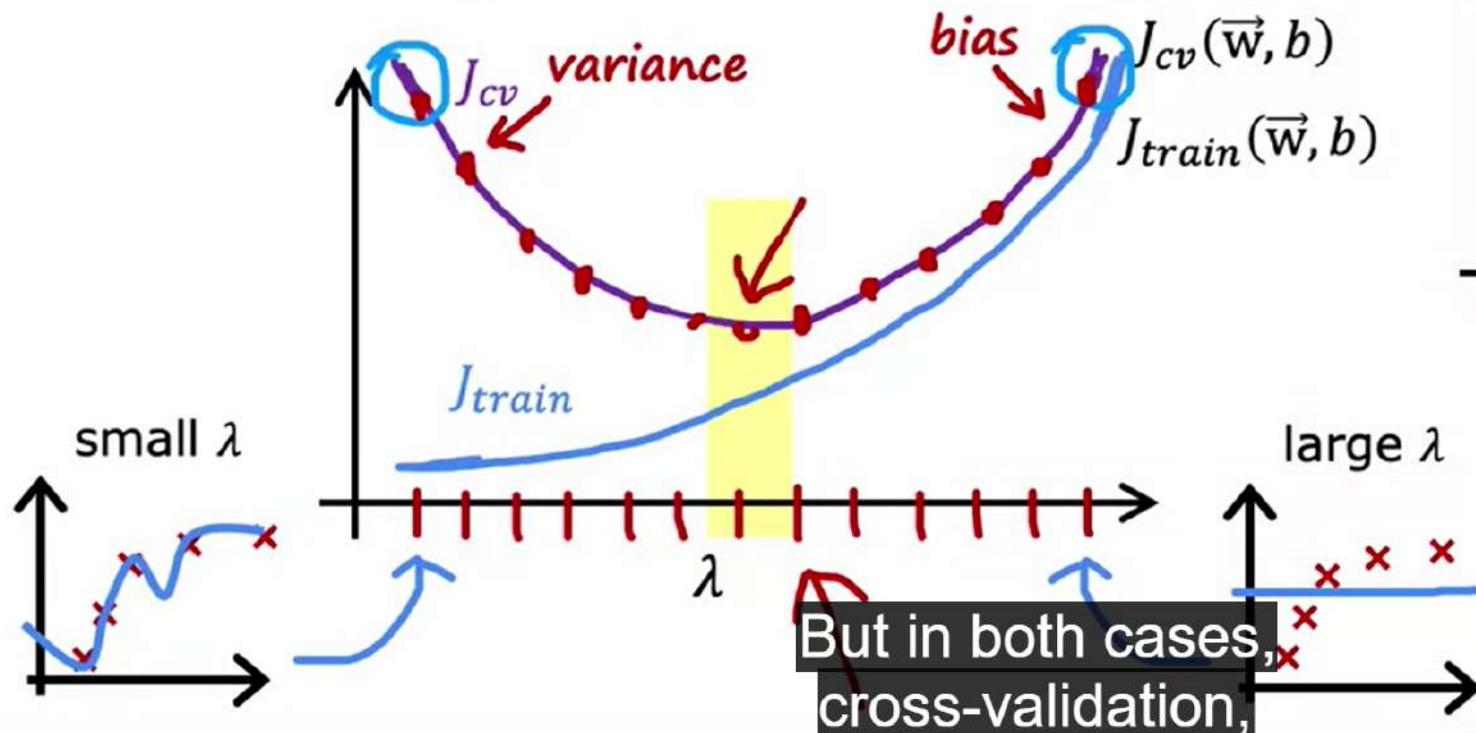
Small  $\lambda$

High variance (overfit)

$$\lambda = 0$$

## Bias and variance as a function of regularization parameter $\lambda$

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



# Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

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

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

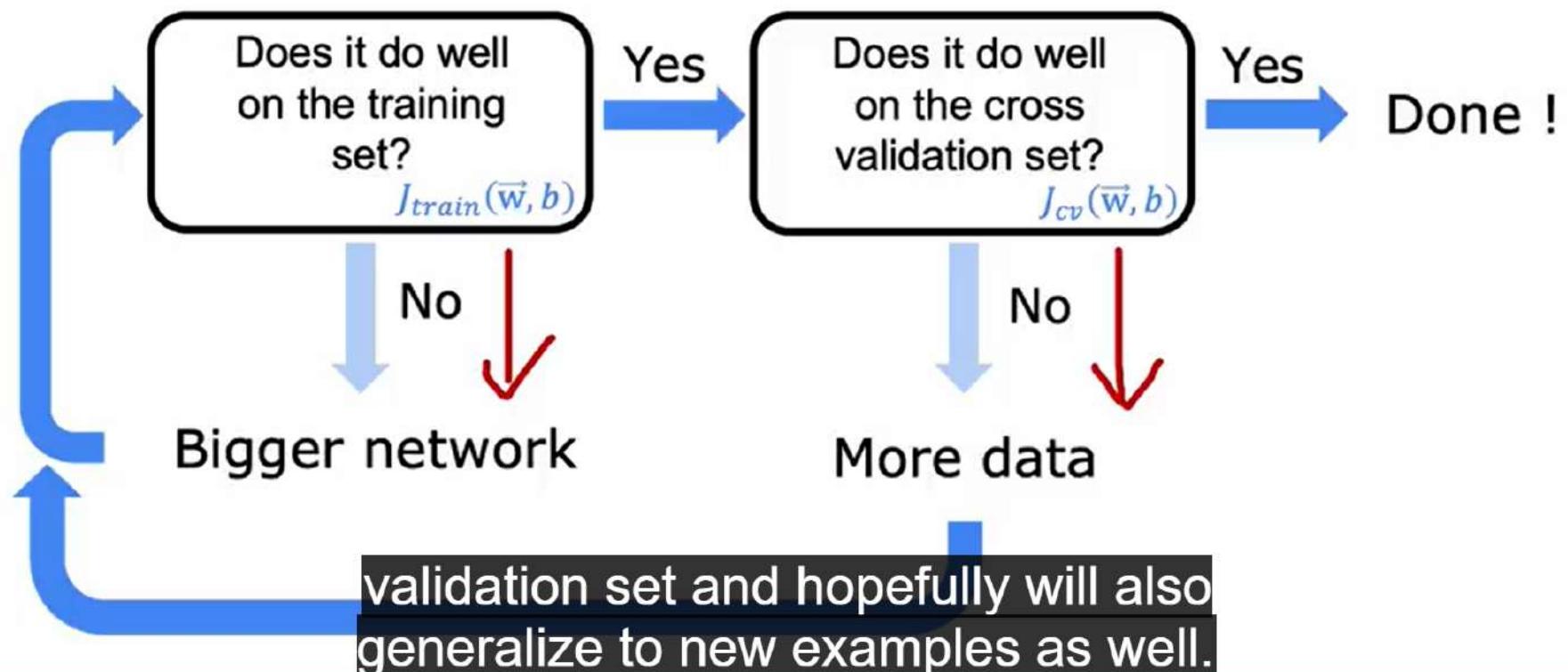
- Get more training examples
- Try smaller sets of features  $x, x^2, \cancel{x^3}, \cancel{x^4}, \cancel{x^5}, \cancel{y}$ ...
- Try getting additional features
- Try adding polynomial features ( $x_1^2, x_2^2, x_1x_2, \text{etc}$ )
- Try decreasing  $\lambda$
- Try increasing  $\lambda$

I realized that this was a  
lot of stuff on this slide.

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 network regularization

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

b

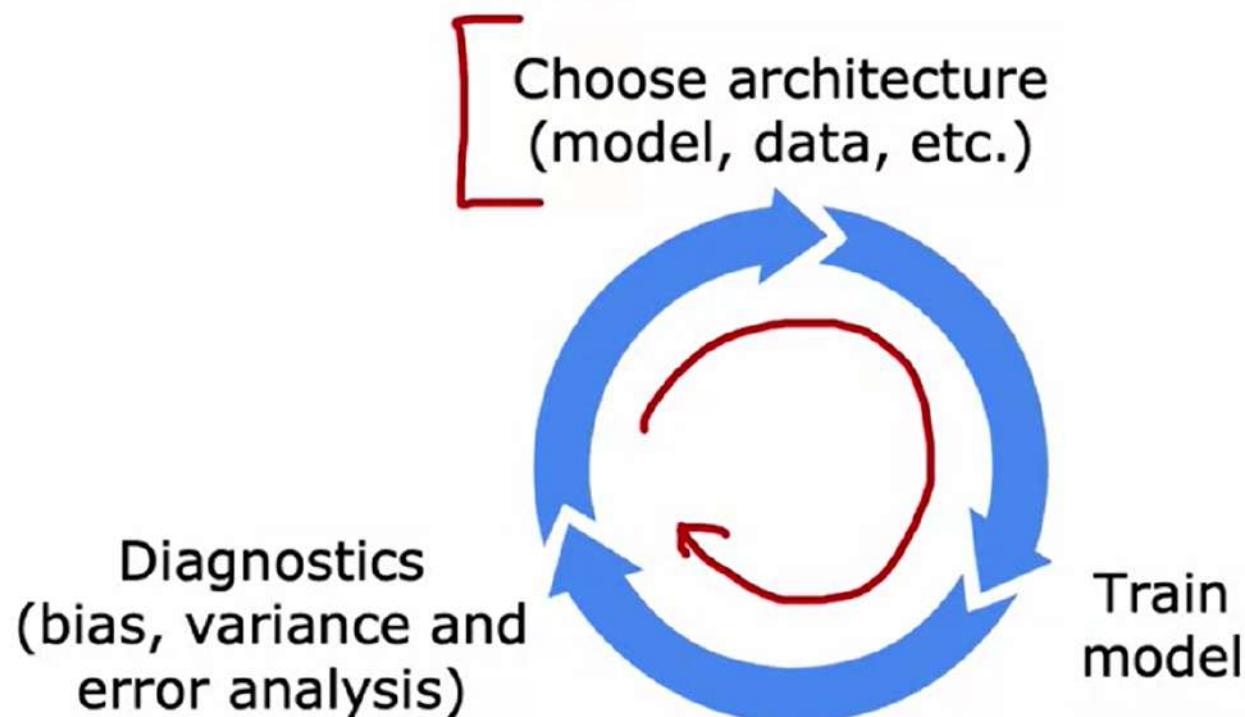
## Unregularized MNIST model

```
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])
```

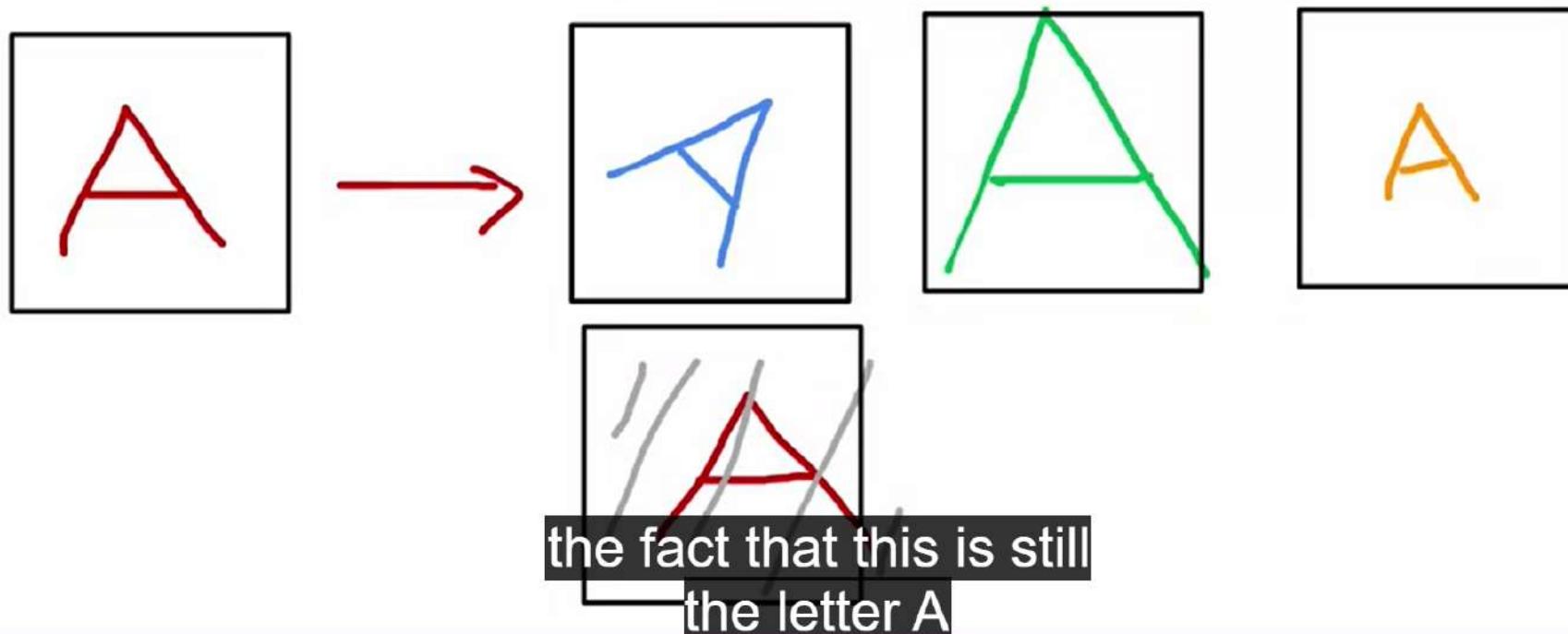
# Iterative loop of ML development



variance diagnostic as well as

# Data augmentation

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



# 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
  - audio sound like you're recording it on  
a bad cell phone connection like this.

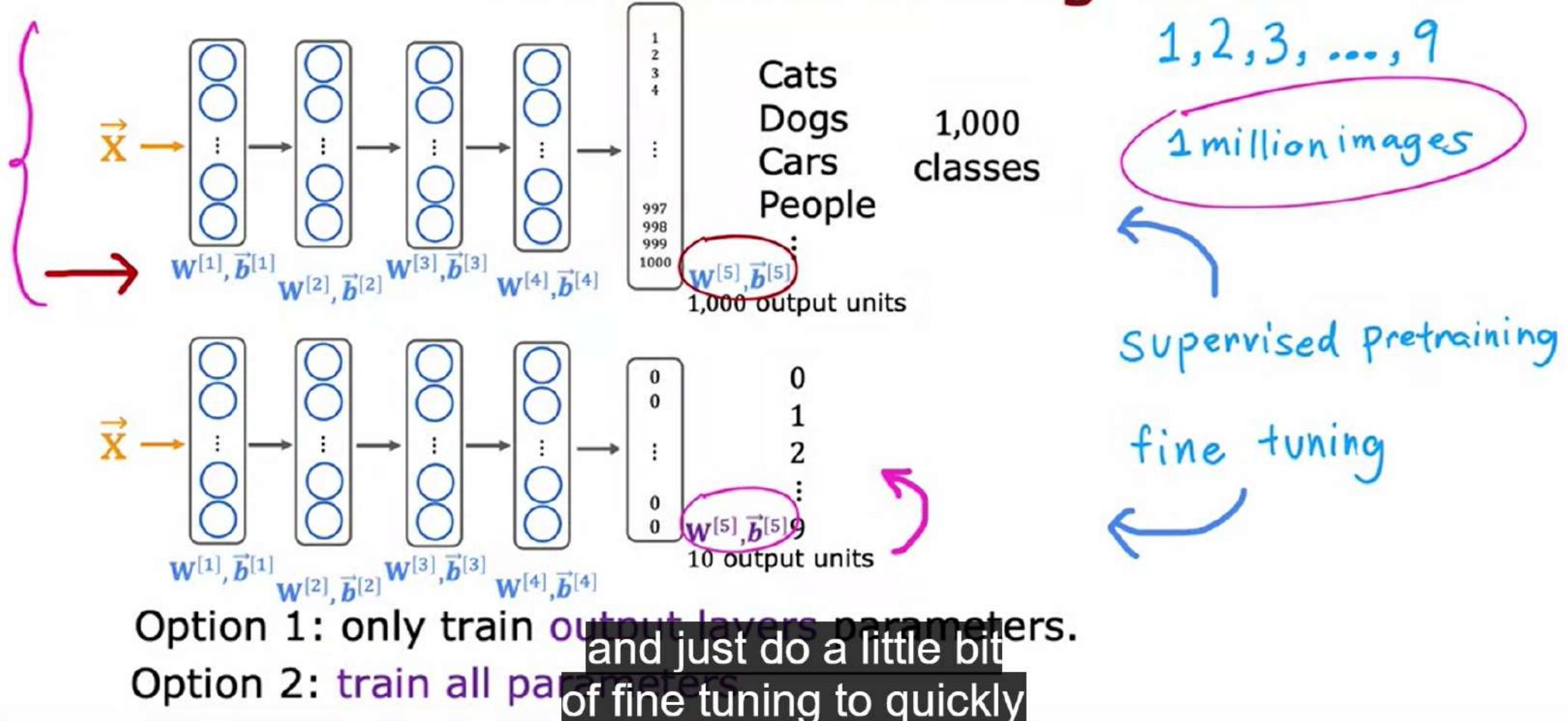
# Artificial data synthesis for photo OCR



Real data  
So with synthetic data generation  
like this you can generate a very

[Adam Coates and Tao Wang]

# Transfer learning

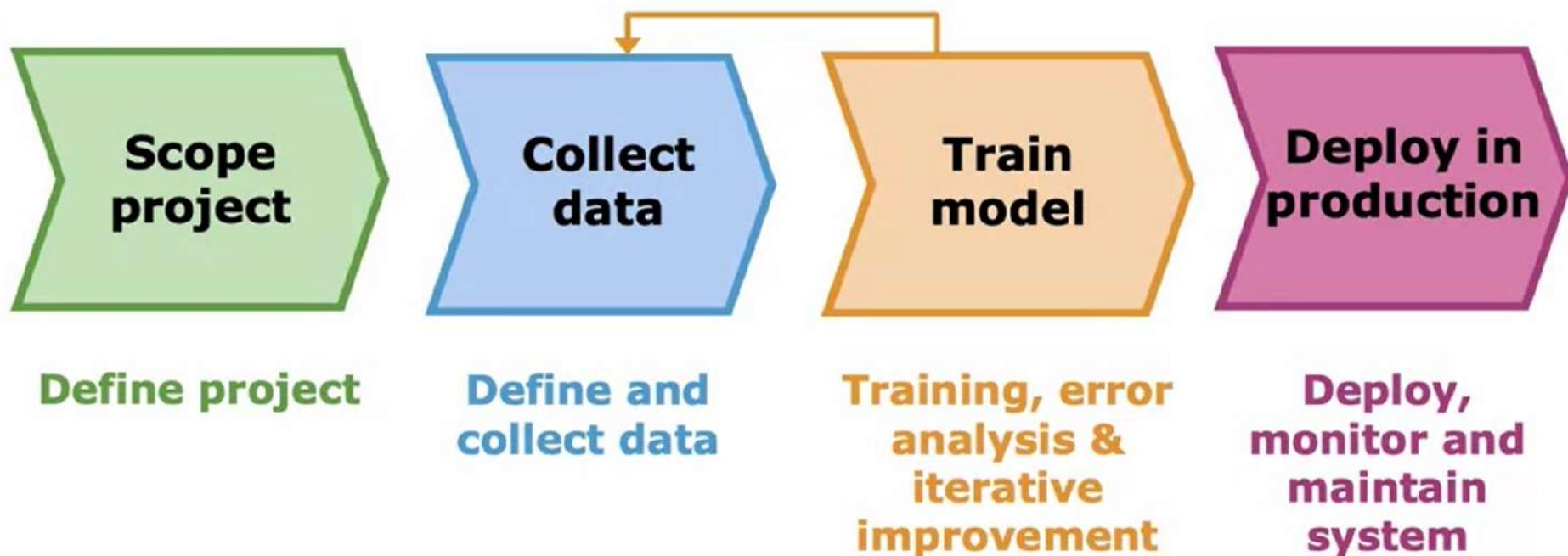


# Transfer learning summary

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).  
*1 million images*
2. Further train (fine tune) the network on your own data.  
*1000 images*  
*50 images*

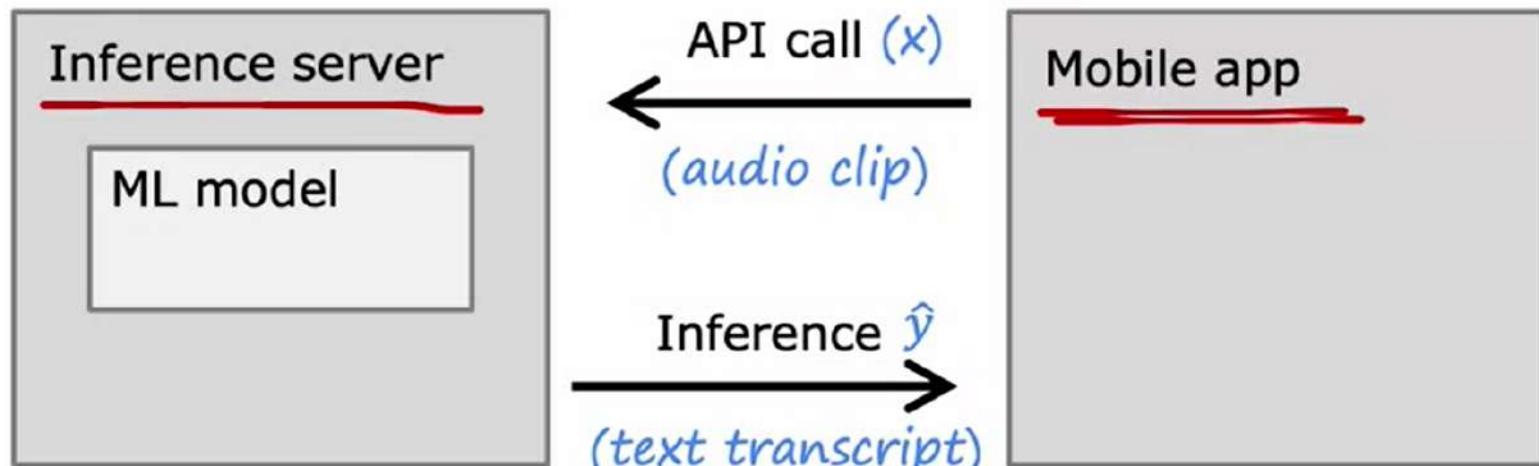
potentially much smaller dataset.

# Full cycle of a machine learning project



What that means is you make it

# Deployment



→ Software engineering may be needed for:

Ensure reliable and efficient predictions

Scaling

Logging

System monitoring

Model update **has good laws, is monitored,**

MLOps  
machine learning  
operations

# 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?)

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

which is 0.6 or 60 percent.

# Trading off precision and recall

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

- Predict 1 if  $f_{\vec{w}, b}(\vec{x}) \geq 0.5$  0.7 0.7 0.3
- Predict 0 if  $f_{\vec{w}, b}(\vec{x}) < 0.5$  0.7 0.1 0.3

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

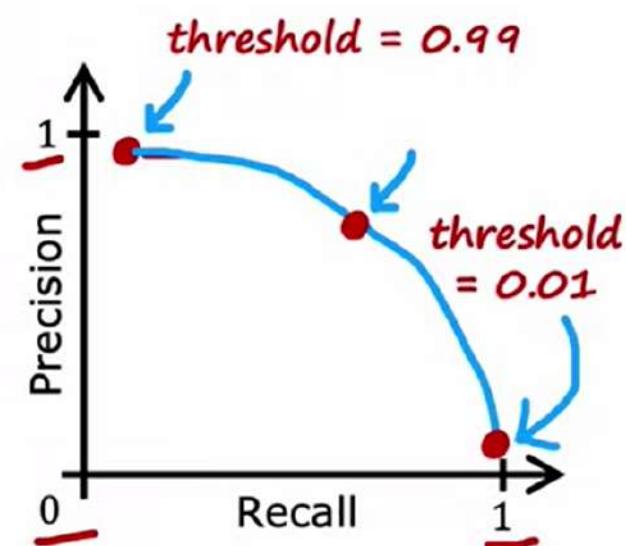
higher precision, lower recall

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

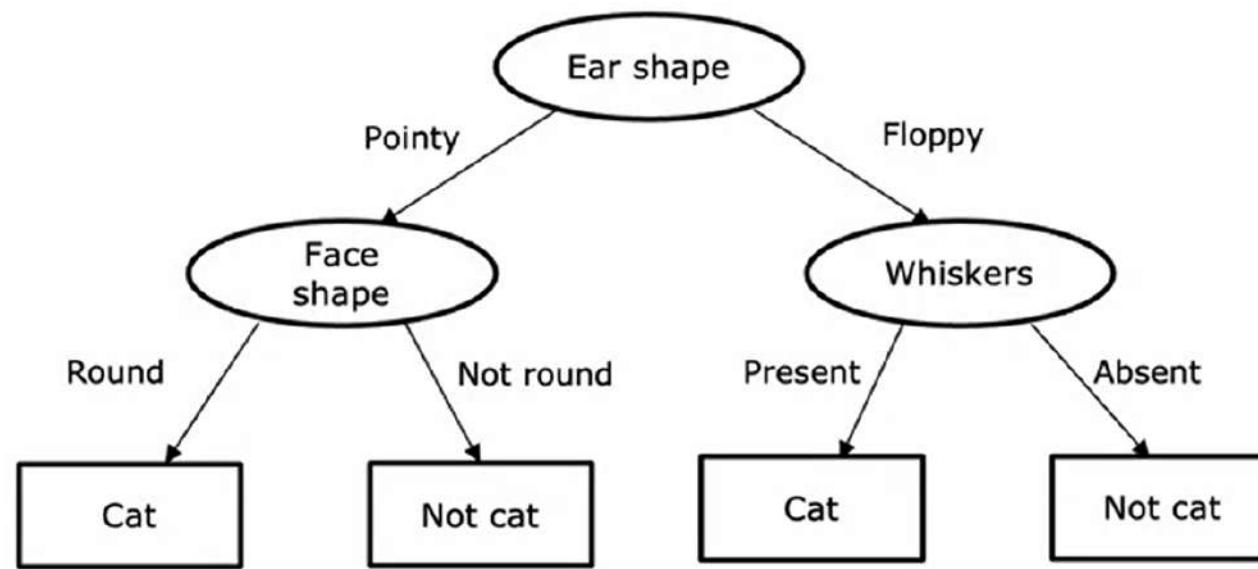
lower precision, higher recall

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

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

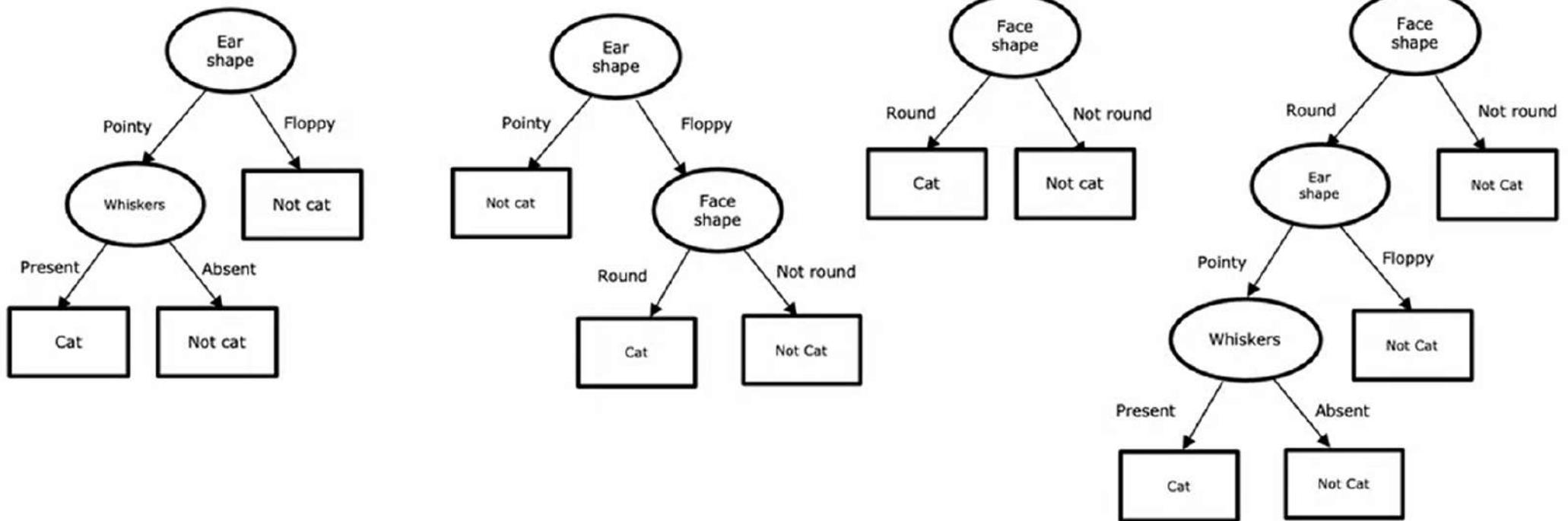


# Decision Tree



the learning algorithm  
looks like a tree,

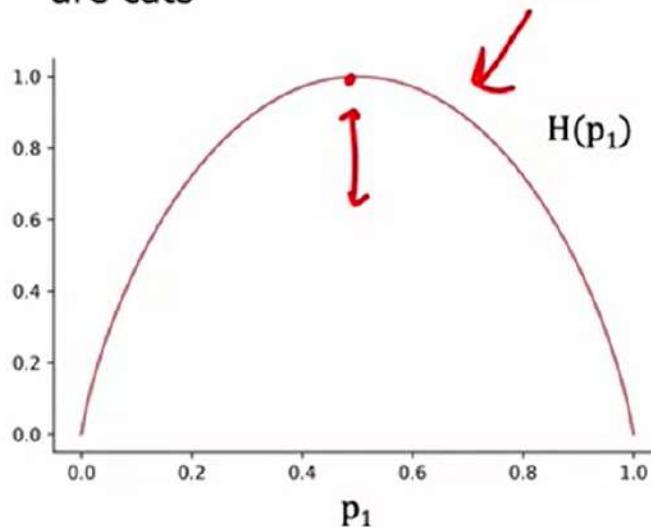
# Decision Tree



to try to pick one that

# Entropy as a measure of impurity

$p_1$  = fraction of examples that are cats



$$p_0 = 1 - p_1$$

$$H(p_1) = -p_1 \log_2(p_1) - p_0 \log_2(p_0)$$

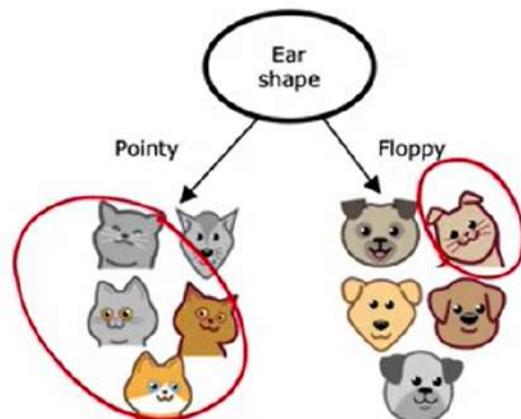
$$= -p_1 \log_2(p_1) - (1 - p_1) \log_2(1 - p_1)$$



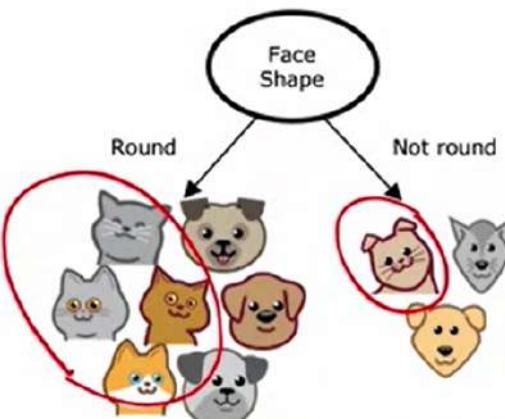
Note: “ $0 \log(0)$ ” = 0

To summarize, the entropy function is

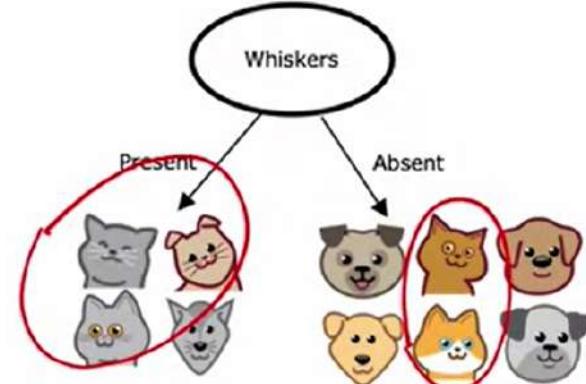
# Choosing a split



$$p_1 = \frac{4}{5} = 0.8 \quad p_1 = \frac{1}{5} = 0.2$$
$$H(0.8) = 0.72 \quad H(0.2) = 0.72$$



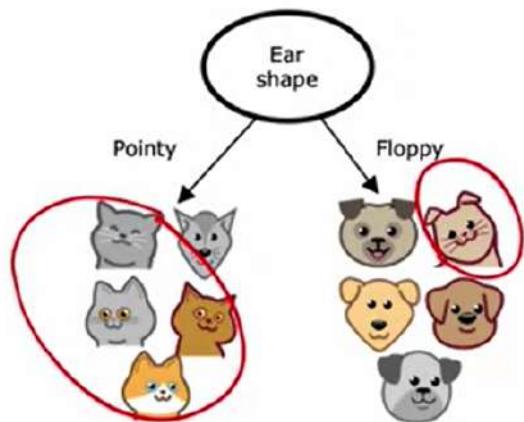
$$p_1 = \frac{4}{7} = 0.57 \quad p_1 = \frac{1}{3} = 0.33$$
$$H(0.57) = 0.99 \quad H(0.33) = 0.92$$



$$p_1 = \frac{3}{4} = 0.75 \quad p_1 = \frac{2}{6} = 0.33$$
$$H(0.75) = 0.81 \quad H(0.33) = 0.92$$

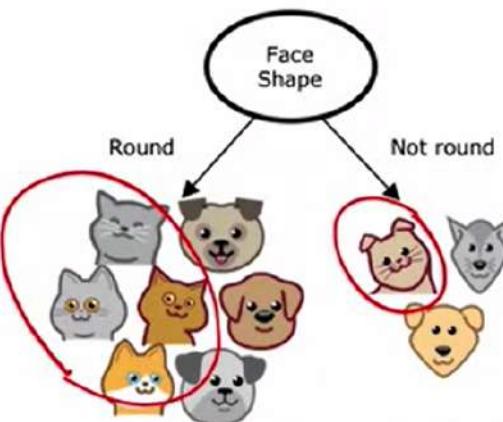
given these three options

# Choosing a split



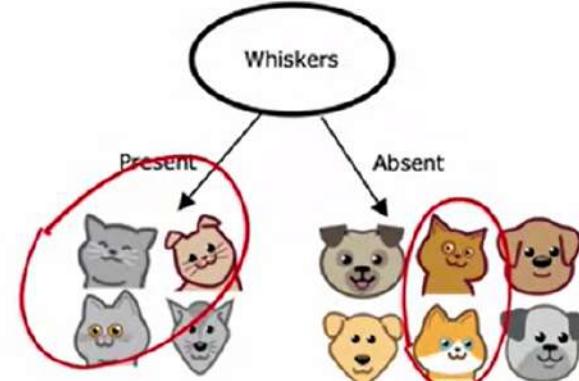
$$p_1 = \frac{4}{5} = 0.8 \quad p_2 = \frac{1}{5} = 0.2$$
$$H(0.8) = 0.72 \quad H(0.2) = 0.72$$

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



$$p_1 = \frac{4}{7} = 0.57 \quad p_2 = \frac{1}{3} = 0.33$$
$$H(0.57) = 0.99 \quad H(0.33) = 0.92$$

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



$$p_1 = \frac{3}{4} = 0.75 \quad p_2 = \frac{2}{6} = 0.33$$
$$H(0.75) = 0.81 \quad H(0.33) = 0.92$$

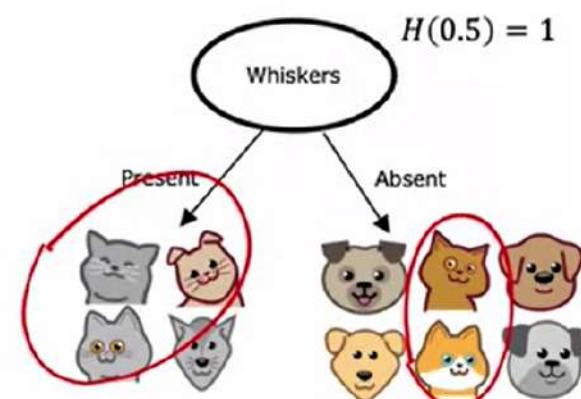
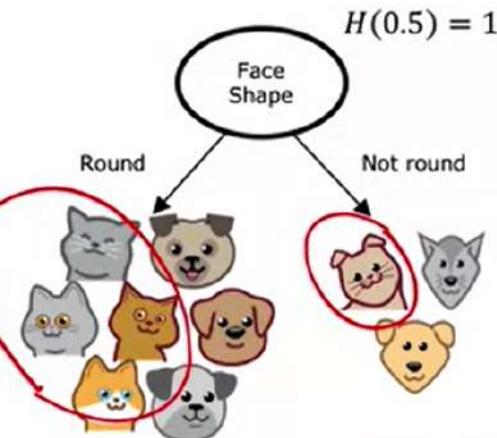
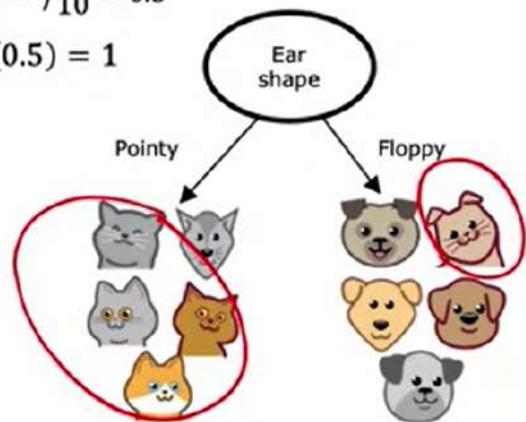
$$\left( \frac{4}{10}H(0.75) + \frac{6}{10}H(0.33) \right)$$

If we go to the root node,

# Choosing a split

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

$$H(0.5) = 1$$



$$p_1 = \frac{4}{5} = 0.8 \quad p_1 = \frac{1}{5} = 0.2$$
$$H(0.8) = 0.72 \quad H(0.2) = 0.72$$

$$H(0.5) - \left( \frac{5}{10} H(0.8) + \frac{5}{10} H(0.2) \right)$$
$$= 0.28$$

$$p_1 = \frac{4}{7} = 0.57 \quad p_1 = \frac{1}{3} = 0.33$$
$$H(0.57) = 0.99 \quad H(0.33) = 0.92$$

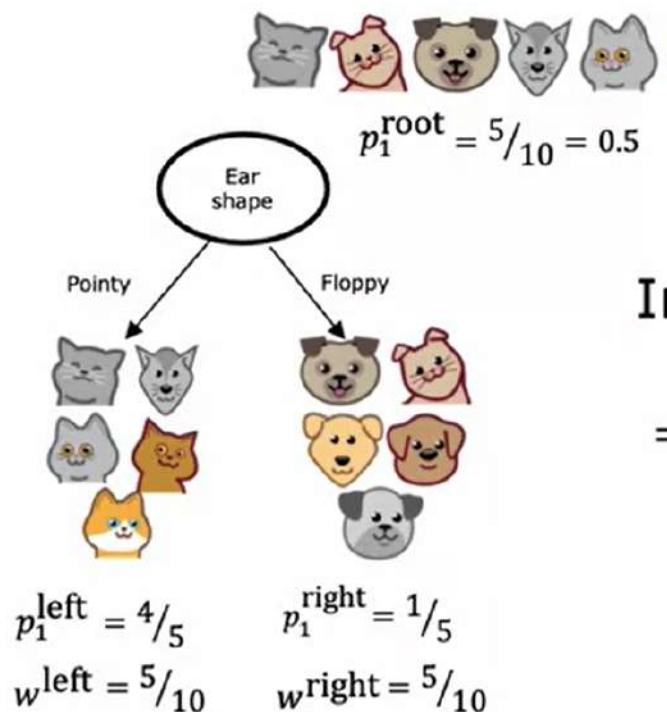
$$H(0.5) - \left( \frac{7}{10} H(0.57) + \frac{3}{10} H(0.33) \right)$$
$$= 0.03$$

$$p_1 = \frac{3}{4} = 0.75 \quad p_1 = \frac{2}{6} = 0.33$$
$$H(0.75) = 0.81 \quad H(0.33) = 0.92$$

$$H(0.5) - \left( \frac{4}{10} H(0.75) + \frac{6}{10} H(0.33) \right)$$
$$= 0.12$$

Information gain

# Information Gain

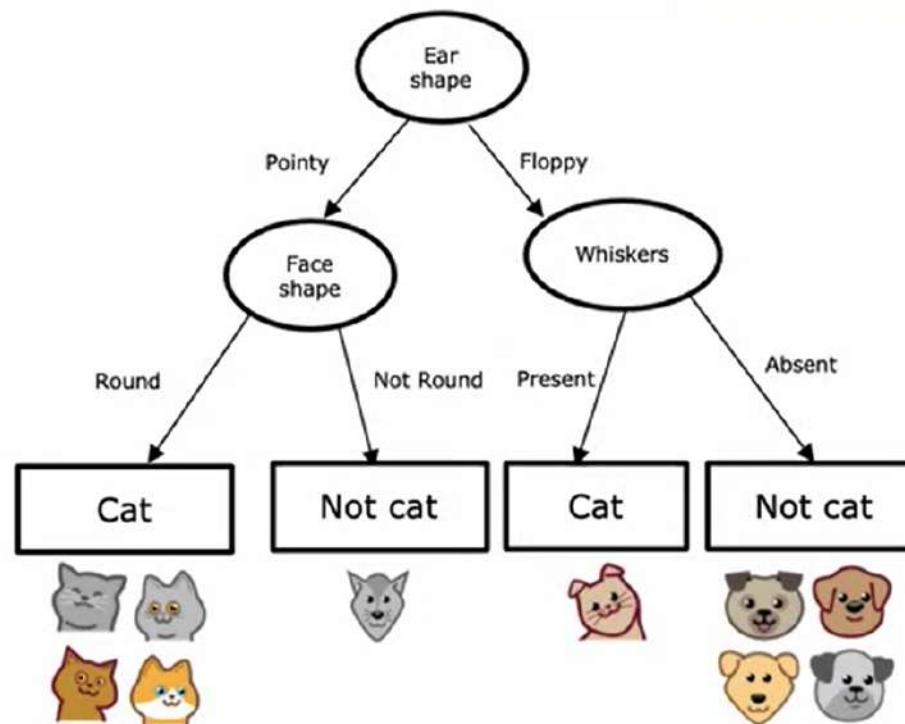


Information gain

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

That will result in, hopefully,

# Recursive splitting



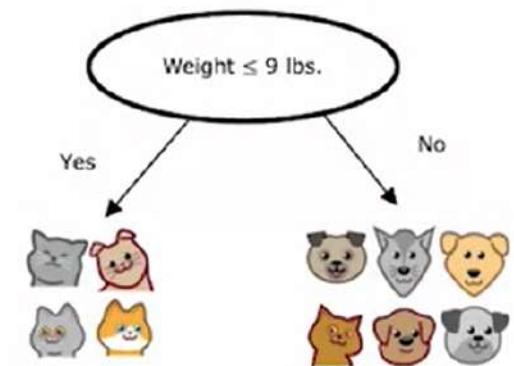
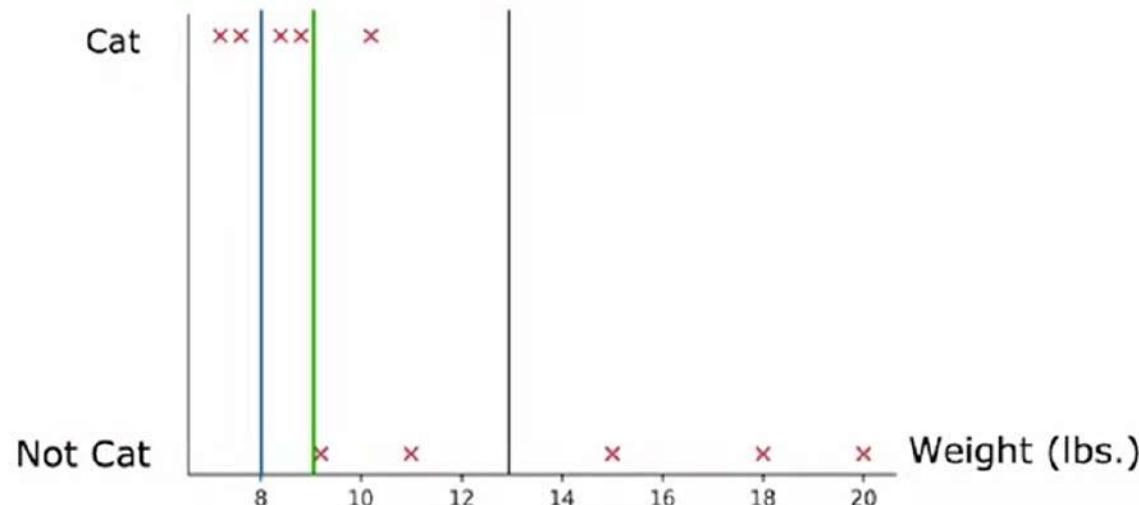
Notice that there's

# 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 so instead of one feature  
taking on three possible values,

# Splitting on a continuous variable



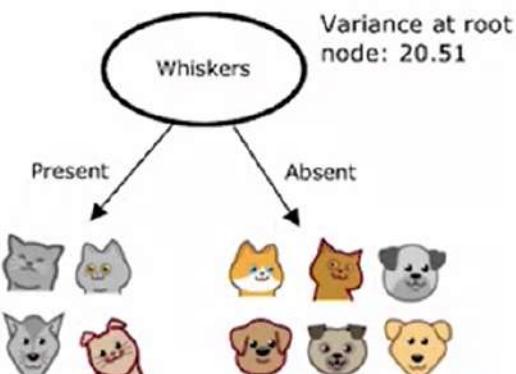
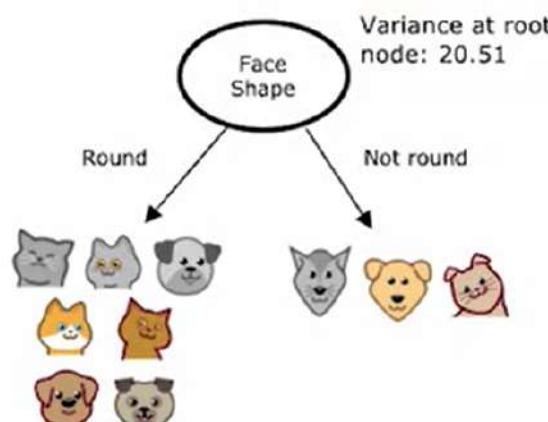
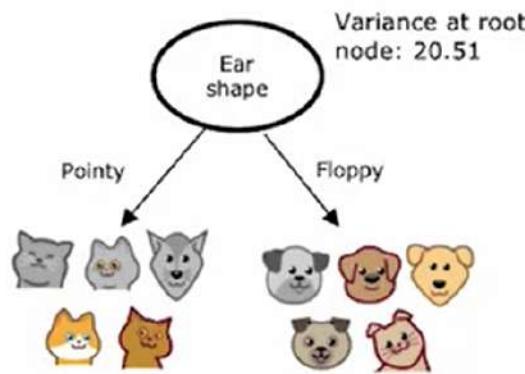
$$H(0.5) - \left( \frac{2}{10} H\left(\frac{2}{2}\right) + \frac{8}{10} H\left(\frac{3}{8}\right) \right) = 0.24$$

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

$H(0.5) - \left( \frac{7}{10} H\left(\frac{5}{7}\right) + \frac{3}{10} H\left(\frac{2}{3}\right) \right) = 0.40$

you can then build recursive,  
li additional decision trees using

## Choosing a split



Weights: 7.2, 9.2, 8.4, 7.6, 10.2  
Variance: 1.47

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

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

Weights: 7.2, 15, 8.4, 7.6, 10.2, 18, 20  
Variance: 27.80

$$w^{\text{left}} = 7/10 \quad w^{\text{right}} = 3/10$$

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

Weights: 7.2, 8.8, 9.2, 8.4  
Variance: 0.75

$$w^{\text{left}} = 4/10 \quad w^{\text{right}} = 6/10$$

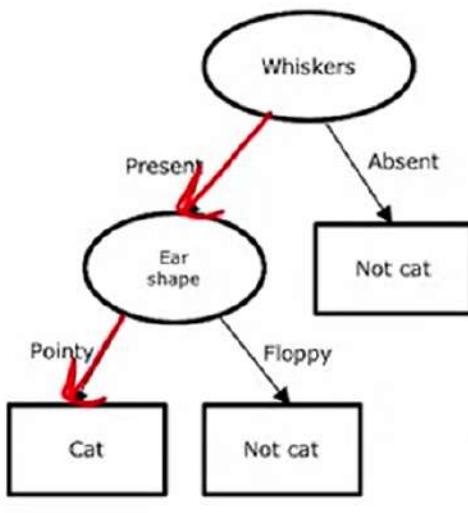
$$20.51 - \left( \frac{4}{10} * 0.75 + \frac{6}{10} * 23.32 \right) \\ = 6.22$$

# Tree ensemble

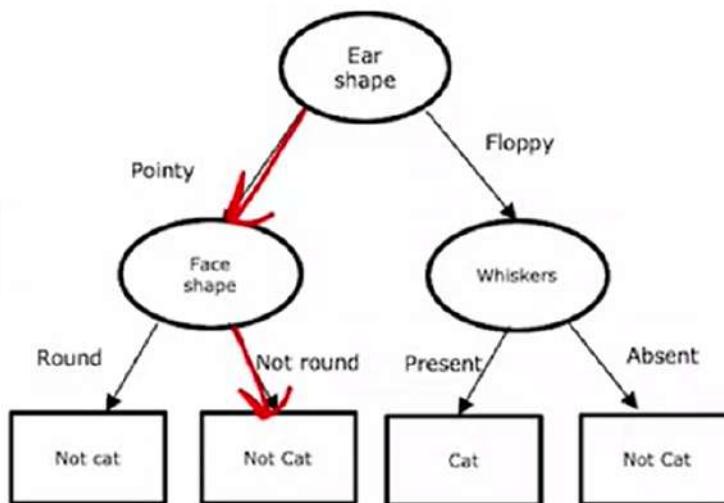
New test example



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

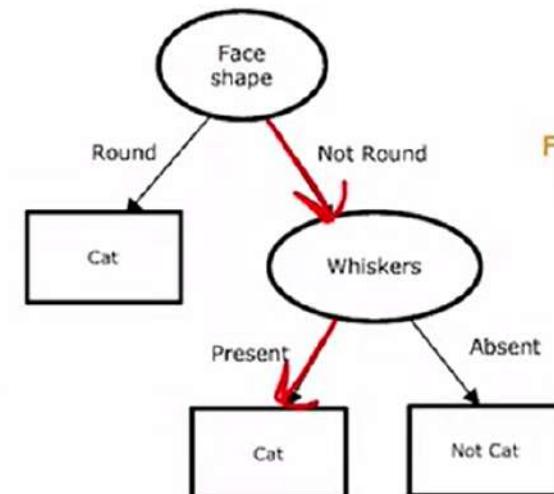


Prediction: Cat



Prediction: Not cat

which happens to be the  
Final prediction: Cat

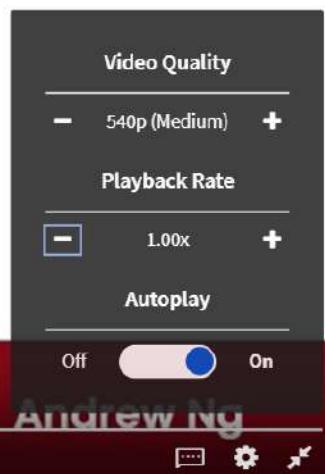


Prediction: Cat

# 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.

So in our example we had three features available rather than picking from all end



# Using XGBoost

## Classification

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

practitioners will use the open source  
libraries that implement XGBoost.

# 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)
```

classification, then this line  
here just becomes XGBRegressor and

# 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

critical to getting  
competitive performance.