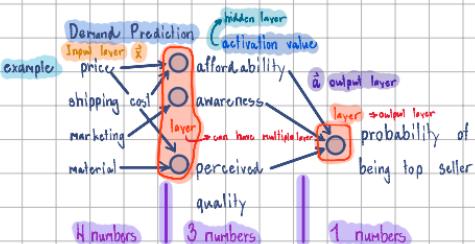
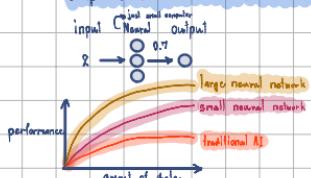


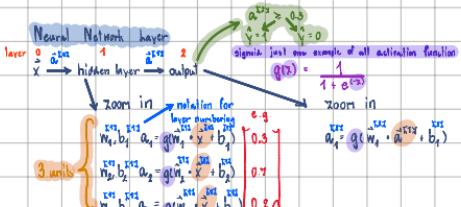
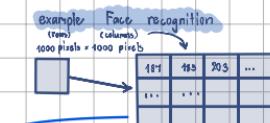
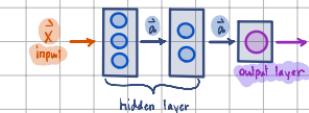
## Neural networks

Origins: Algorithms that try to mimic the brain  
 speech → images → text (Neural language program) → anything

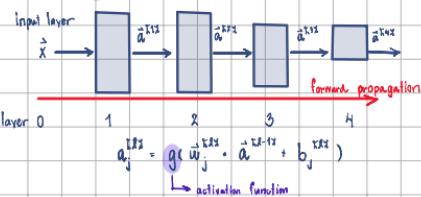
## Simplified mathematical model of a neuron



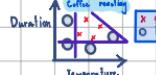
**summary** If just logis regression, they can learn its own features like "feature engineering" but Neural Networks if doing by yourself that make it easier to accurate predictions



## More complex neural networks



## Inference in Code (Tensorflow implementation)



Build the model using TensorFlow

```
x = np.array([1, 2, 3, 4, 5])
layer_1 = tf.layers.dense(inputs=x, activation='sigmoid')
# layer_1 = tf.sigmoid(x)

# inference (1) - return from tensor to array
y = np.array(layer_1)
print(y)

# inference (2) - pass back to tensorflow
y2 = np.array([0.5, 0.9, 0.1])
x2 = tf.layers.dense(inputs=y2, activation='sigmoid')
z2 = tf.sigmoid(x2)

# print(z2)
```

## Data in Tensorflow

```
x = np.array([[1, 2, 3], [4, 5, 6]])
x = np.array([1, 2, 3]) # 1D array
x = np.array([1, 2, 3, 4]) # 2D row vector
```

## Building a neural network

### Building a neural network architecture

```
x = np.array([1, 2, 3, 4])
x = Dense(2, activation='softmax')
x = Dense(3, activation='sigmoid')
x = Dense(4, activation='softmax')

y = np.array([1, 2, 3, 4])
y = Dense(2, activation='softmax')
y = Dense(3, activation='sigmoid')
y = Dense(4, activation='softmax')

model.compile(...)

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

## Forward prop in a single layer

### forward prop (coffee roasting model)

```
# 1D arrays
w1_1 = np.array([1, 2])
w1_2 = np.array([3, 4])
w1_3 = np.array([5, 6])
w2_1 = np.array([1])
w2_2 = np.array([1])
w2_3 = np.array([1])
x1_1 = np.array([1, 2])
x1_2 = np.array([3, 4])
x1_3 = np.array([5, 6])
x2_1 = np.array([1])
x2_2 = np.array([1])
x2_3 = np.array([1])

# 2D arrays
w1 = np.array([[w1_1, w1_2, w1_3], [w2_1, w2_2, w2_3]])
w2 = np.array([[w1_1, w1_2, w1_3], [w2_1, w2_2, w2_3], [x1_1, x1_2, x1_3], [x2_1, x2_2, x2_3]])
x1 = np.array([x1_1, x1_2, x1_3])
x2 = np.array([x2_1, x2_2, x2_3])

# forward prop
z1 = np.dot(w1, x1) + b1
a1 = sigmoid(z1)
z2 = np.dot(w2, a1) + b2
a2 = sigmoid(z2)
```



Multi-class  $\Rightarrow$  more than two possible values



**Softmax**  $\Rightarrow$  generalization of logistic regression to the multiclass

example Softmax regression ( $N$  possible outputs)

$$x = \vec{w} \cdot \vec{x} + b_1$$

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

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

$$x = \vec{w} \cdot \vec{x} + b_2$$

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

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

$$x = \vec{w} \cdot \vec{x} + b_3$$

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

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

$$x = \vec{w} \cdot \vec{x} + b_4$$

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

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

(2 possible output)  
logistic regression

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

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

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

( $N$  possible output)  
Softmax regression

$$x_j = \vec{w}_j \cdot \vec{x} + b_j; j = 1, 2, \dots, N$$

$$a_j = \frac{e^{x_j}}{\sum_{k=1}^N e^{x_k}} = P(Y=j|\vec{x})$$

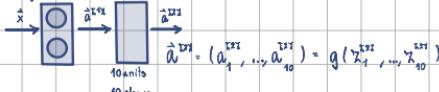
Cost Function for Softmax regression

$$a_i = \frac{e^{x_i}}{N e^{x_1} + e^{x_2} + \dots + e^{x_N}}$$

$$\text{loss}(a_1, \dots, y) = \begin{cases} -\log a_1 & \text{if } y=1 \\ -\log a_y & \text{if } y \neq 1 \end{cases}$$

Neural Network with Softmax Output

change unit of output layer equal to classes



Improved implementation of softmax

Numerical Roundoff Errors e.g.  $\frac{2}{10000} \neq \frac{1}{10000} + \frac{1}{10000} + \dots + \frac{1}{10000}$

$\hookrightarrow$  More numerically accurate implementation of logistic loss

Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

$$\text{model} = \text{Sequential}(\text{...}, \text{Dense(units=25, activation='relu')}, \text{Dense(units=15, activation='relu')}, \text{Dense(units=1, activation='sigmoid')})$$

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

Original loss

$$\text{loss} = -y \log(g(x)) - (1-y) \log(1-g(x)) \quad \text{model.compile(loss='BinaryCrossEntropy')}$$

More accurate loss (in code)

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

More numerically accurate implementation of softmax

Softmax regression

$$\text{model} = \text{Sequential}(\text{...}, \text{Dense(units=25, activation='relu')}, \text{Dense(units=15, activation='relu')}, \text{Dense(units=10, activation='softmax')})$$

Loss  $L(\vec{a}, y) = \begin{cases} -\log a_i & \text{if } y=1 \\ -\log a_y & \text{if } y \neq 1 \end{cases}$

More Accurate

$$L(\vec{a}, y) = \begin{cases} -\log\left(\frac{a_i}{\sum_{k=1}^N a_k}\right) & \text{if } y=1 \\ -\log\left(\frac{a_y}{\sum_{k=1}^N a_k}\right) & \text{if } y \neq 1 \end{cases}$$

$\text{model.compile(loss='SparseCategoricalCrossEntropy(from_logits=True)')}$

logistic regression  
(more numerically accurate)

model

$$\text{model} = \text{Sequential}(\text{...}, \text{Dense(units=25, activation='sigmoid')}, \text{Dense(units=15, activation='sigmoid')}, \text{Dense(units=1, activation='linear')})$$

from tensorflow.keras.losses import BinaryCrossentropy

loss

$\text{model.compile(..., BinaryCrossentropy(from_logits=True))}$

fit

$\text{model.fit(X, epochs=100)}$

predict

$\text{f_x = tf.nn.sigmoid(model(X)))}$

Classification with multiple outputs

$\Rightarrow$  create specific Neural Network for single output and combine them or create it in one neural network with multiple outputs

Advanced Optimization  $\Rightarrow$  another algorithm not just gradient descent

'Adam' algorithm  $\Rightarrow$  not just one  $\alpha$  (learning rate)

$\Rightarrow$  If  $w_j$  and  $b$  keeps moving in same direction  $\Rightarrow$  increase  $\alpha$

$\Rightarrow$  If  $w_j$  and  $b$  keeps oscillating  $\Rightarrow$  reduce  $\alpha$

MNIST Adam

model

$$\text{model} = \text{Sequential}(\text{...}, \text{tf.keras.layers.Dense(units=25, activation='sigmoid')}, \text{tf.keras.layers.Dense(units=15, activation='sigmoid')}, \text{tf.keras.layers.Dense(units=10, activation='linear')})$$

compile

$\text{model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=le-3)}, \text{loss=tf.keras.losses.SparseCategoricalCrossEntropy(from_logits=True)})$

fit

$\text{model.fit(X, epochs=100)}$

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

## Additional layer types

- Dense layer  $\Rightarrow$  each neuron's output is a function of all the activation outputs of the previous layer



*different part*

- Convolution layer  $\Rightarrow$  part of the previous layer's output feeds training data



## Tips

### Deciding what to try next

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features ( $x_1^2, x_1 \cdot x_2, x_1$ , etc.)
- Try decrease / increase learning rates

## Evaluating a model

+ separate dataset into training set, test set, validation set

### compute test error

$$J_{\text{test}}(\hat{w}, b) = \frac{1}{M_{\text{test}}} \sum_{i=1}^{M_{\text{test}}} (\hat{w}^T x_{\text{test}}^{(i)} + b - \hat{y}_{\text{test}}^{(i)})^2$$

### Model selection and training / cross-validation / test sets

#### Model selection

$$\hat{f}_{\hat{w}, b}(x) = W_1 x + b$$

$$\hat{f}_{\hat{w}, b}(x) = W_1 x + W_2 x^2 + b$$

⋮

*collection of test set / evaluation set*

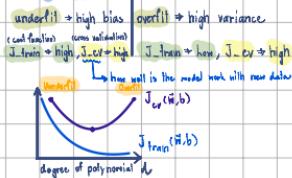
*training / cross-validation / test sets*

training set  $\approx 60\%$

cross-validation  $\approx 20\%$

test set  $\approx 20\%$

## Diagnosing bias and variance



### Regularization and bias/variance

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

$$f_{\hat{w}, b} = \hat{w}^T x + b$$

$$x^A \quad x^B$$

High bias

Large λ

High variance

λ?

### How to choose the regularization parameter?

Try  $\lambda$  from low to high

### Establishing a baseline level of performance

$J_{\text{train}}$   $J_{\text{cv}}$   $J_{\text{test}}$   
→ compare training error, cross-validation to Human level performance or competing algorithms performance

example	Baseline performance	$\rightarrow 10.6\%$	$10.6\%$
	Training error ( $J_{\text{train}}$ )	$\rightarrow 10.6\%$	$10.6\%$
	Cross Validation error ( $J_{\text{cv}}$ )	$\rightarrow 14.5\%$	$15.0\%$

*High variance*      *High bias*

### Learning curves

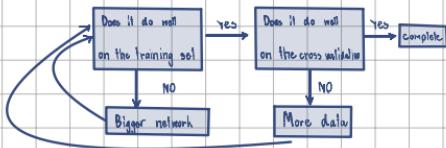


### Deciding what to try next revisited

- Get more training examples
- Try smaller sets of features
- Try Increase  $\lambda$
- Try getting additional features
- Try adding polynomial features
- Try decrease  $\lambda$

## Bias / Variance and neural networks

large neural networks are low bias machines



## Iterative loop for ML development

(model, data, compute, etc.)

Choose architecture

Diagnostics  
(bias, variance, error analysis)

Train model

**Adding data** → add more data of the types where error analysis has indicated

**Data augmentation** → modifying an existing example to create new training example

e.g. original sound + noise

**Data synthesis** → using artificial data inputs to create a new training example

e.g. change type of font

## Tweaking the data used by your system

1) Conventional model-centric approach (Improve code)

2) Data-centric approach (Improve data)

$$AI = \text{Code} + \text{Data}$$

## Transfer learning: using pretrained model

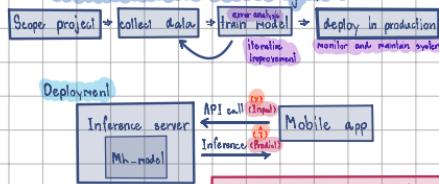
supervised pretraining + train network with large data

fine tuning → train the network on your own data

2 options  
→ 1) only train output layers parameters (small training set)

2) train all parameters (large training set)

## Full cycle of a machine learning project



93.4% or 93.2% or 93.1% (real-life test?)

## Error metrics for skewed datasets

→ how to find best algorithm (good prediction or useful prediction)

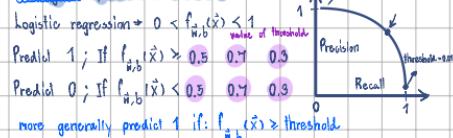
### Precision / recall

$y = 1$  in presence of rare class

		actual class			
		1	0		
predicted class	1	True positive False negative	False positive True negatives	Precision → check the ratio that model predicted	
	0	0	10	True positive = $\frac{15}{20} = 0.75$	True positive + False neg = 20
		25	75	Recall → check how much model predicting correctly	
				True positive = $\frac{15}{25} = 0.6$	True pos + False neg = 25

### Selecting threshold

#### Trading off precision and recall



### F1 score → compare precision/recall

$$\text{F1 score} = \frac{1}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} = \frac{1}{\frac{\text{precision} + \text{recall}}{\text{precision} \cdot \text{recall}}}$$

	Precision	Recall	Average (all intervals)	F1 score
Algorithm 1	0.5	0.4	0.45	0.444 (bad)
Algorithm 2	0.7	0.1	0.6	0.175 (lower value → then the higher value)
Algorithm 3	0.02	1.0	0.501	0.0982