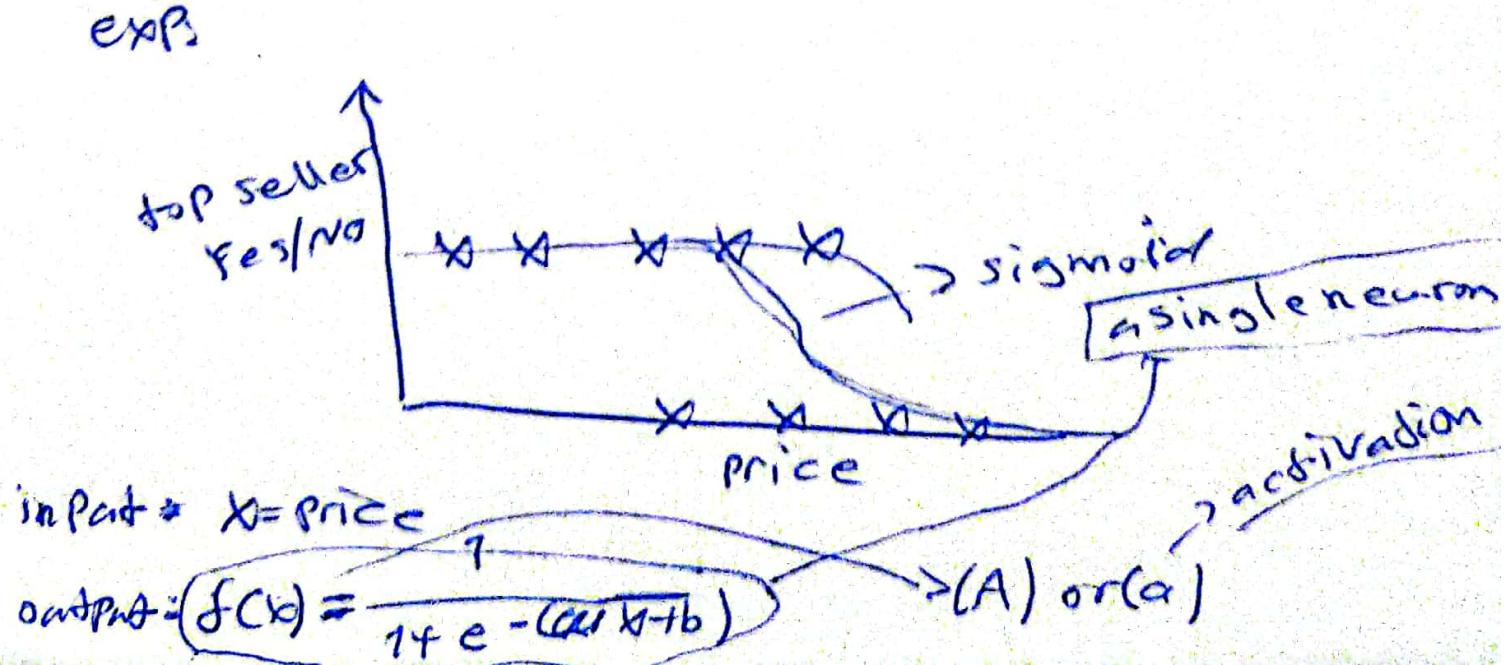
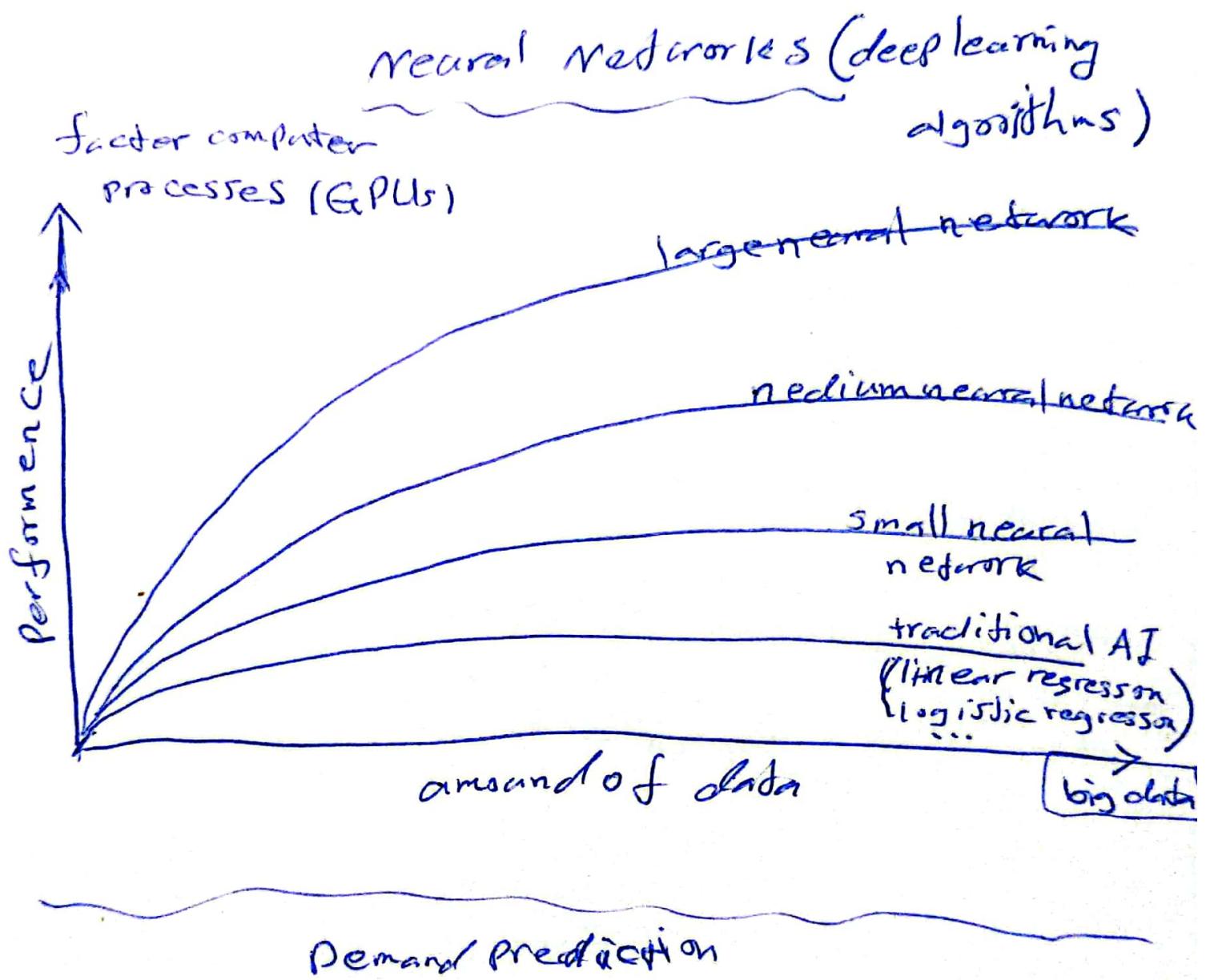
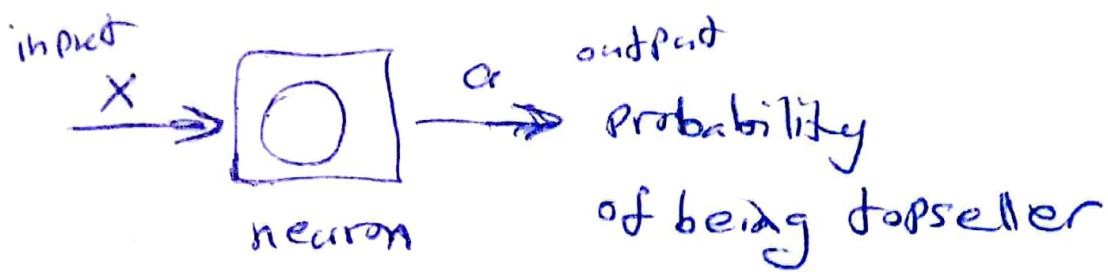


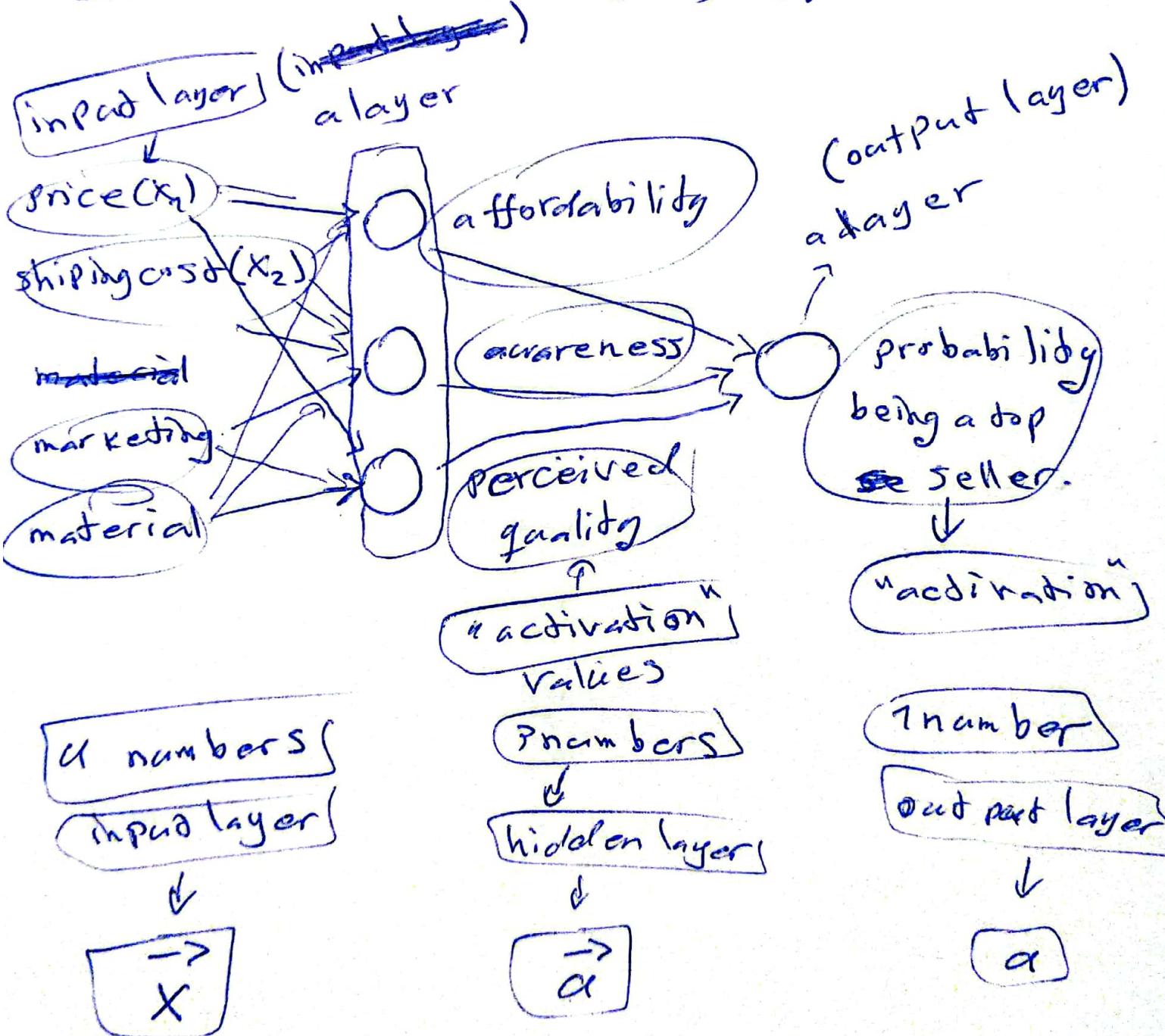
"Advanced learning Algorithms" ①



(2)

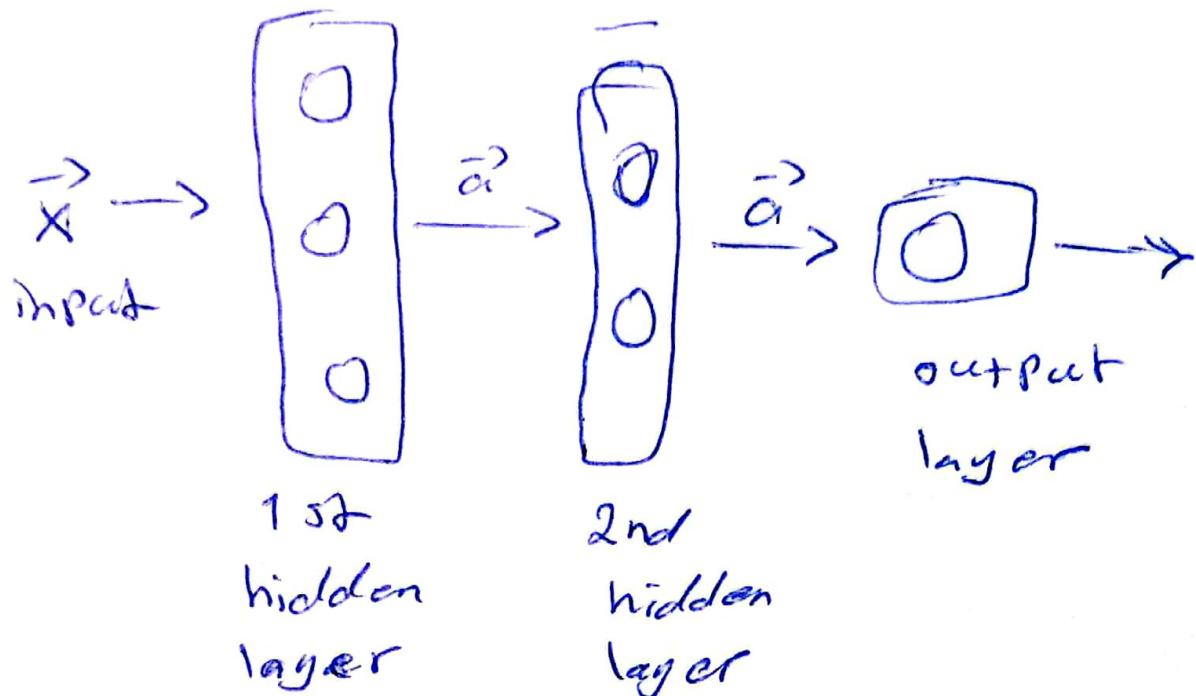


ExP: features (x_1, x_2, x_3, x_4)

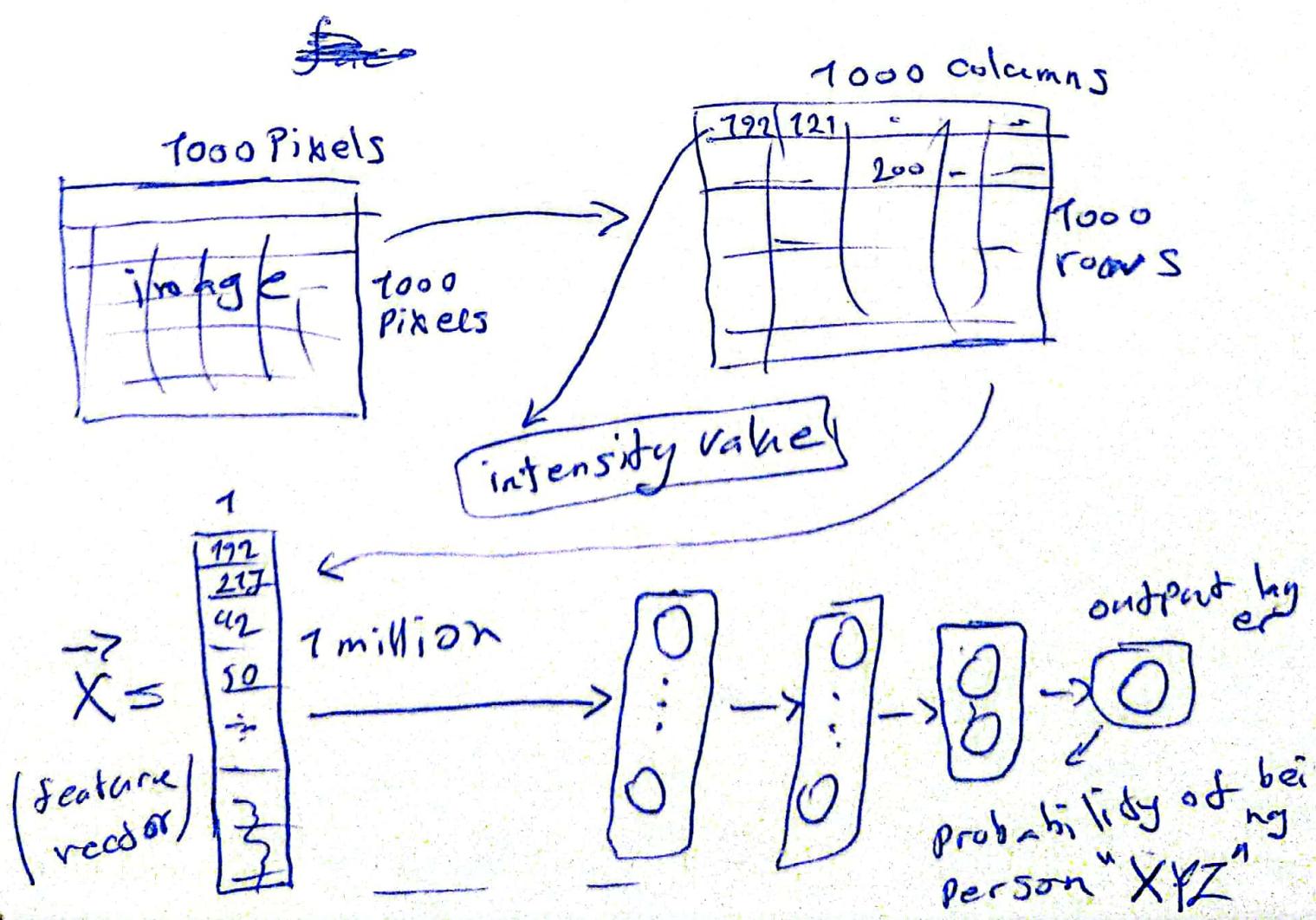


③

a multi layer perceptron (neural network)

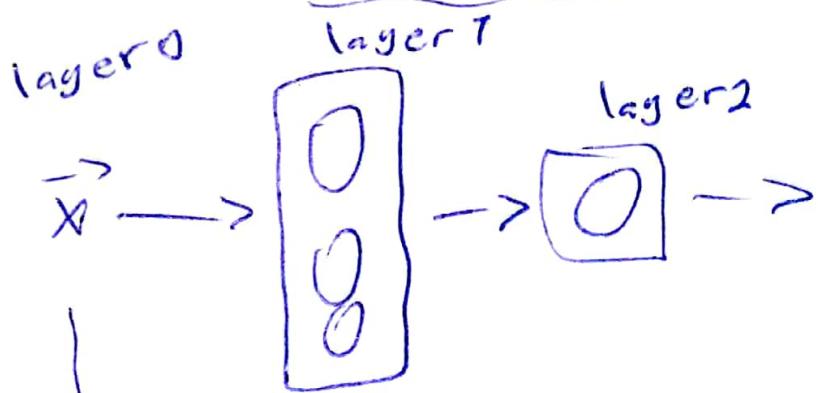


Recognizing Images

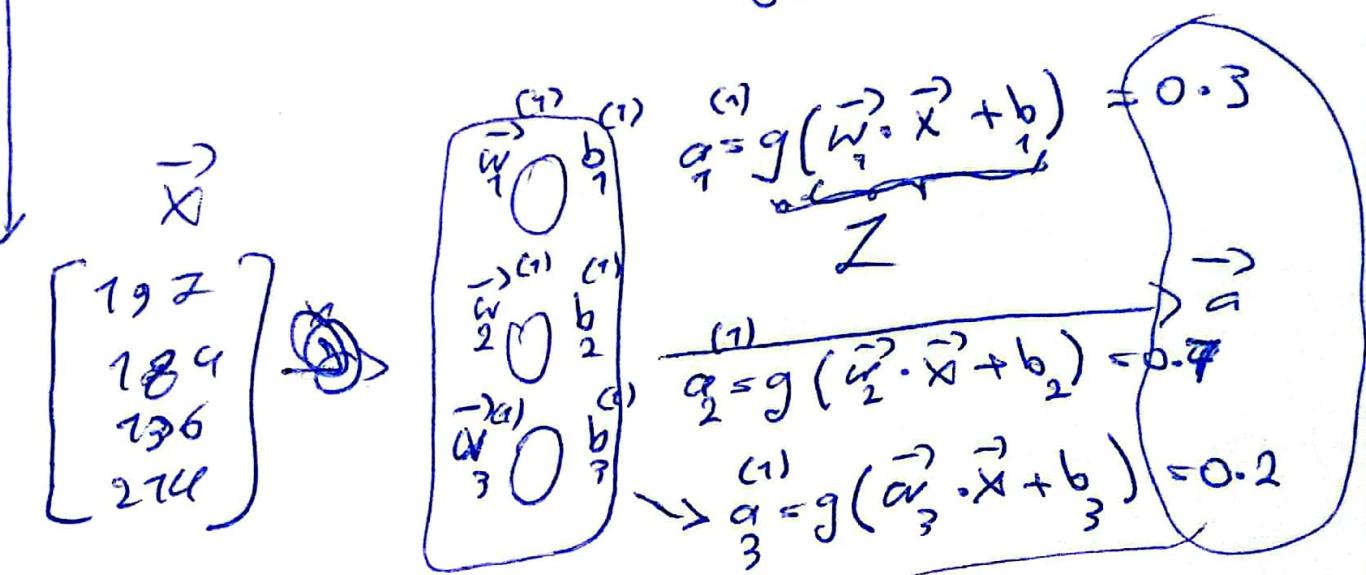


9

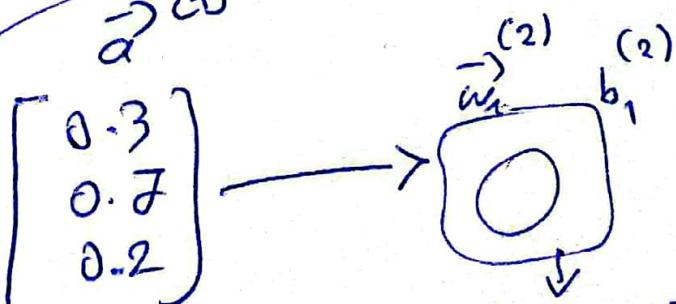
neural network layer



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



$\vec{a}^{(1)}$ → output of layer 1



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

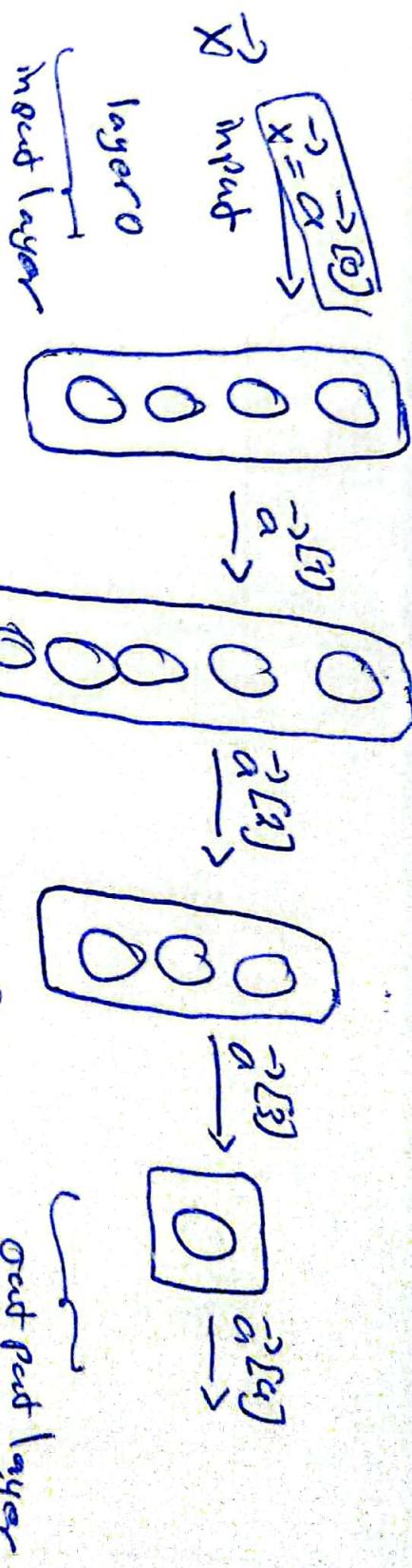
Decision logic:

Is $a^{(2)} > 0.5$?

- If Yes: $\hat{Y} = 1$
- If No: $\hat{Y} = 0$

(optional) Output = 0.84

(5) Exp



* in counting the layer

the input layer (layer 0) is not counted So this model has

4 layers

(sigmoid or activation function)

$$\left[\begin{array}{l} a_1^{(3)} = g(\vec{w}_1^{(3)} \cdot \vec{a}^{(2)} + b_1^{(3)}) \\ a_2^{(3)} = g(\vec{w}_2^{(3)} \cdot \vec{a}^{(2)} + b_2^{(3)}) \\ g^{(3)} = g(\vec{w}_3^{(3)} \cdot \vec{a}^{(2)} + b_3^{(3)}) \end{array} \right]$$

$$\vec{a}^{(3)} = \begin{bmatrix} a_1^{(3)} \\ a_2^{(3)} \\ a_3^{(3)} \end{bmatrix}$$

$$g(z) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

The generalized form

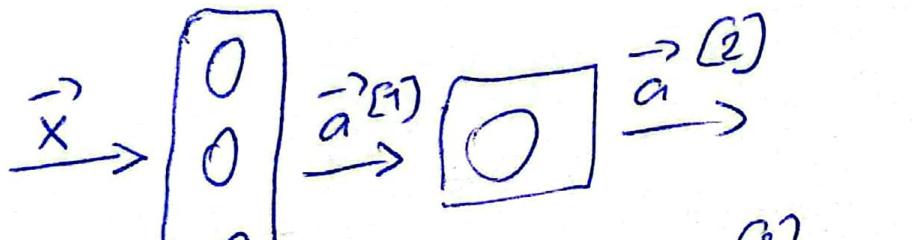
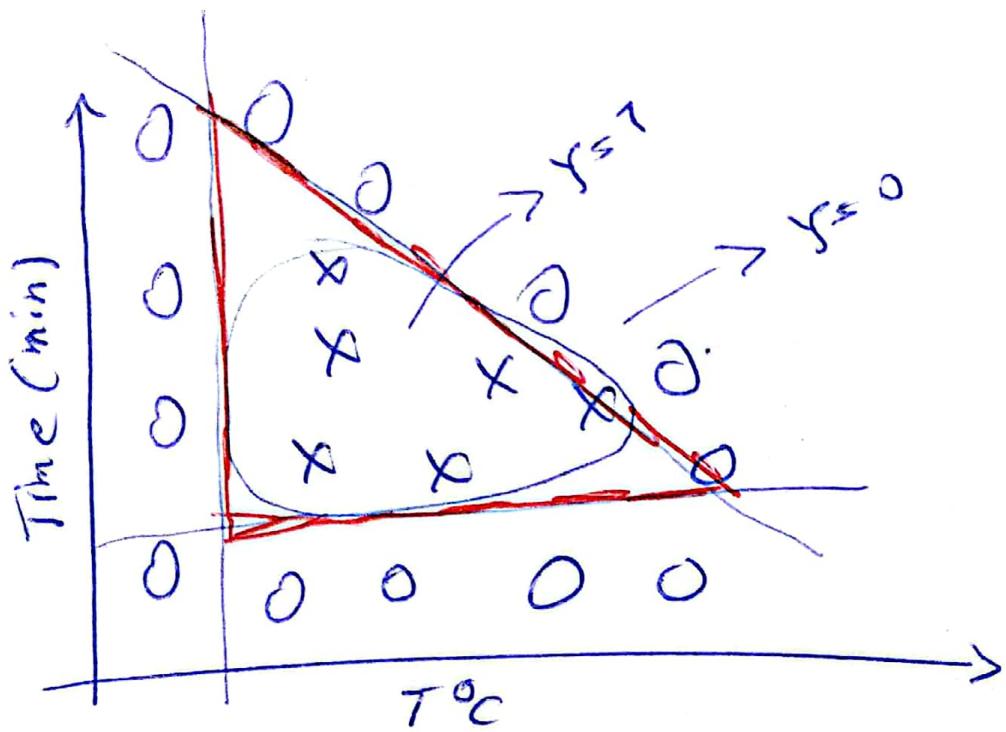
$$a_j^{(n)} = g(\vec{w}_j \cdot \vec{a}^{(n-1)} + b_j^{(n)})$$

6

"TensorFlow"

(Exp)

coffee roasting



is $a_1^{(2)} > 0.5$?

yes

No
 $y = 0$

Tensor code

$\vec{a}^{(1)} = \tanh$

$\vec{x} = \text{np.array}([200, 17])$

layer-1 = Dense(units=3, activation="sigmoid")

$a_1 = \text{layer-1}(\vec{x}) \xrightarrow{\text{sigmoid}} \begin{bmatrix} 0.2 \\ 0.3 \\ 0.5 \end{bmatrix}$

↓ next page

layer-2 = Dense(units=7, activation="sigmoid") 7

$a_2 = \text{layer-2}(a_1)$

if $a_2 \geq 0.5$:

$y_{\text{hat}} = 1$

else:

$y_{\text{hat}} = 0$

exp2

$x = \text{narray}([0.0, \dots, 2^{25}, \dots, 2^{a_0}, \dots, 0])$

layer-1 = Dense(units=25, activation="sigmoid")

$a_1 = \text{layer-1}(x)$

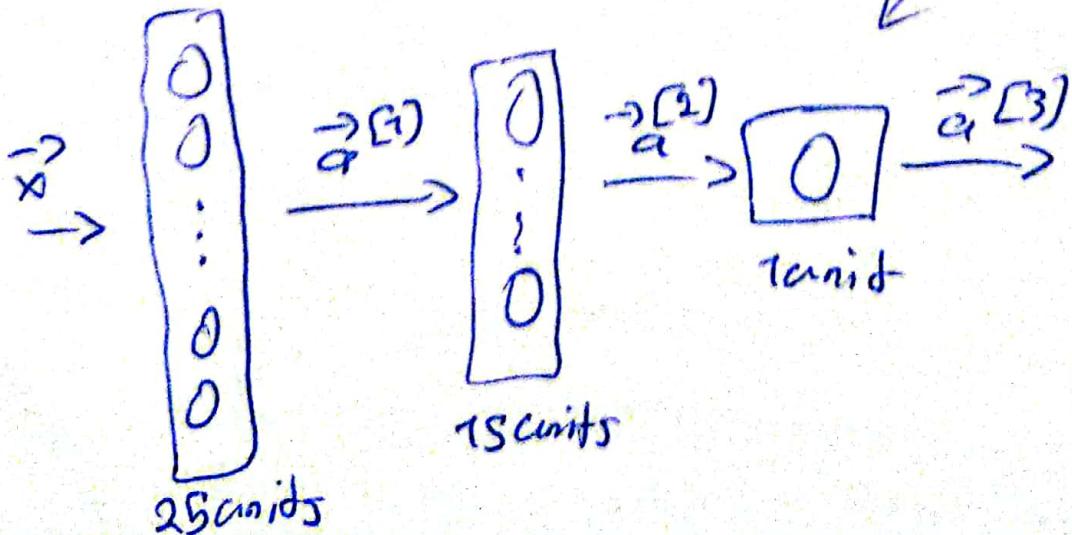
layer-2 = Dense(units=75, activation="sigmoid")

$a_2 = \text{layer-2}(a_1)$

layer-3 = Dense(units=1, activation="sigmoid")

$a_3 = \text{layer-3}(a_2)$

is $\cancel{a_3} \geq 0.5 \}$



"data in Tensorflow"

- Tensorflow is created by Google brain Team.

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

$$x.\text{shape} = (\cancel{1 \times 2})$$

$$\boxed{x = \text{np.array}([[], [1, 2]]) (2 \times 1)}$$

$$\rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(2 \times 3)} \Rightarrow \text{np.array}([[], [1, 2, 3], [4, 5, 6]])$$

↓
↳

$$\rightarrow \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{bmatrix}_{(4 \times 2)} \Rightarrow \text{np.array}([[], [[1, 2], [3, 4], [5, 6], [7, 8]]])$$

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

$$\boxed{\begin{array}{l} \downarrow \\ \text{1D array} \end{array}} \rightarrow \text{a linear array}$$

no rows or columns

accepted

Tensorflow does not recognize 1D array
the input must be in 2D.

exp

9

$x = np.array([[-1, 1], [2, 1]])$

```
layer1 = Dense(3, activation = "sigmoid")
```

$$a_1 = \text{layer-1}(x)$$

$$\rightarrow [[0.2, 0.7, 0.3]]_{(7 \times 3)} \text{ matrix}$$

if we print at:

outputs: `tf.Tensor([[[0.2, 0.7, 0.3]]], shape=(1,3), dtype=fl
oating-point
order)`

`a1.numpy()`

→ converts tensor to numpy ~~matrix~~
array

```
array([6.2, 0.7, 0.3]), dtype= float32)
```

"Building a neural network" in TensorFlow

$\rightarrow \text{layer-1} = \text{Dense}(\text{units=3, activation = "sigmoid"})$

→ layer-2 = **Dense**(units=1, activation = "sigmoid")

```
→ model = Sequential([layer1, layer2])
```

```
x = np.array([[[1, 2], [2, 4], [3, 6], [7, 8]]]) (ax2)
```

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

1D array of length 4.

model.compile(...)

Inference

model.fit(X, y) $\xrightarrow{\text{Inference}}$ model.predict(X-new)

simplify the creation of a model in TensorFlow (10)

model = Sequential([

Dense(anids=3, activation="sigmoid"),

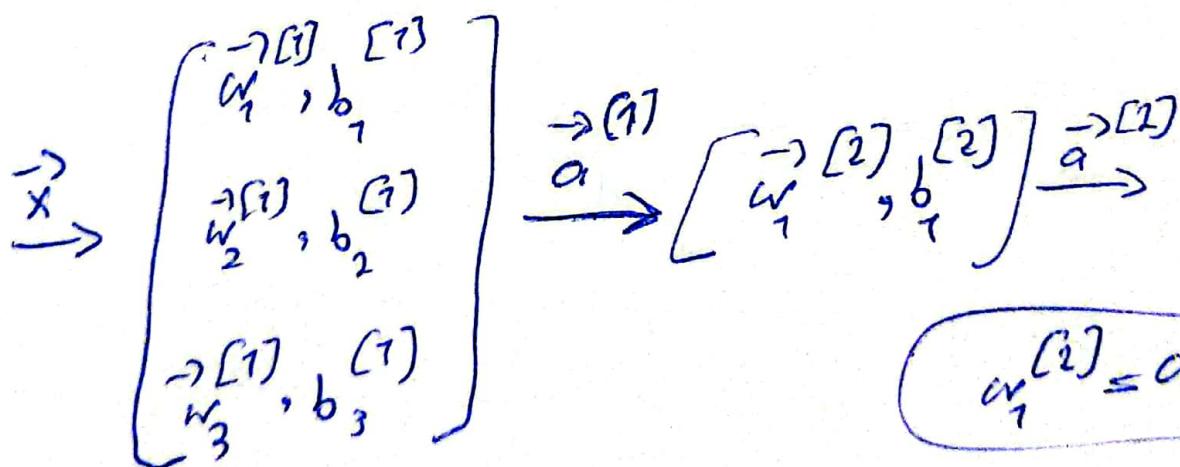
Dense(anids=1, activation="sigmoid")])

model.compile(...)

model.fit(X, Y)

model.predict(X-new)

"forward propagation in a single layer"



$$a_1^{(2)} = a_2 - 7$$

$\vec{x} = np.array([200, 77])$ 1D array

$$a_1^{(1)} = g(\vec{w}_1^{(1)} \cdot \vec{x} + b_1^{(1)})$$

\downarrow
or $\vec{a}^{(0)}$

$w1-1 = np.array([1, 2])$

$b1-1 = np.array([-7])$

$$ZT-1 = \text{np.dot}(wT-1, X) + bT-1$$

$$aT-1 = \text{sigmoid}(ZT-1)$$

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

$$wT-2 = \text{np.array}([-3, 4])$$

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

$$ZT-2 = \text{np.dot}(wT-2, X) + bT-2$$

$$aT-2 = \text{sigmoid}(ZT-2)$$

$$a_3^{[T]} = g(\vec{w}_3^{[T]} \cdot \vec{x} + b_3^{[T]})$$

$$wT-3 = \text{np.array}([5, -6])$$

$$bT-3 = \text{np.array}([2])$$

$$ZT-3 = \text{np.dot}(wT-3, X) + bT-3$$

$$aT-3 = \text{sigmoid}(ZT-3)$$

$$aT = \text{np.array}([aT-1, aT-2, aT-3])$$

$$a_2^{[2]} = g(\vec{w}_1^{[2]} \cdot \vec{aT} + b_1^{[2]})$$

$w2-1 = np.array([7, 8, 9])$

$b2-1 = np.array([3])$

$I2-1 = np.dot(w2-1, a1) + b2-1$

$a2-1 = \text{Sigmoid}(I2-1)$

"General implementation of forward propagation"

activation from the previous layer

def dense(a-in, w, b):

units = w.shape[1]

a-out = np.zeros(units)

for j in range(units):

w = w[:, j]

I = np.dot(w, a-in) + b[j]

a-out[j] = g(I)

return a-out

w = np.array([

[1, -3, 5],

[2, 4, -6]]

(2x3)

b = np.array(

[-1, 7, 2])

1D array

a-in = np.array([-2, 4])

$\frac{\partial f}{\partial x} = X$

def sequential(X):

a1 = dense(X, w1, b1)

a2 = dense(a1, w2, b2)

a3 = dense(a2, w3, b3)

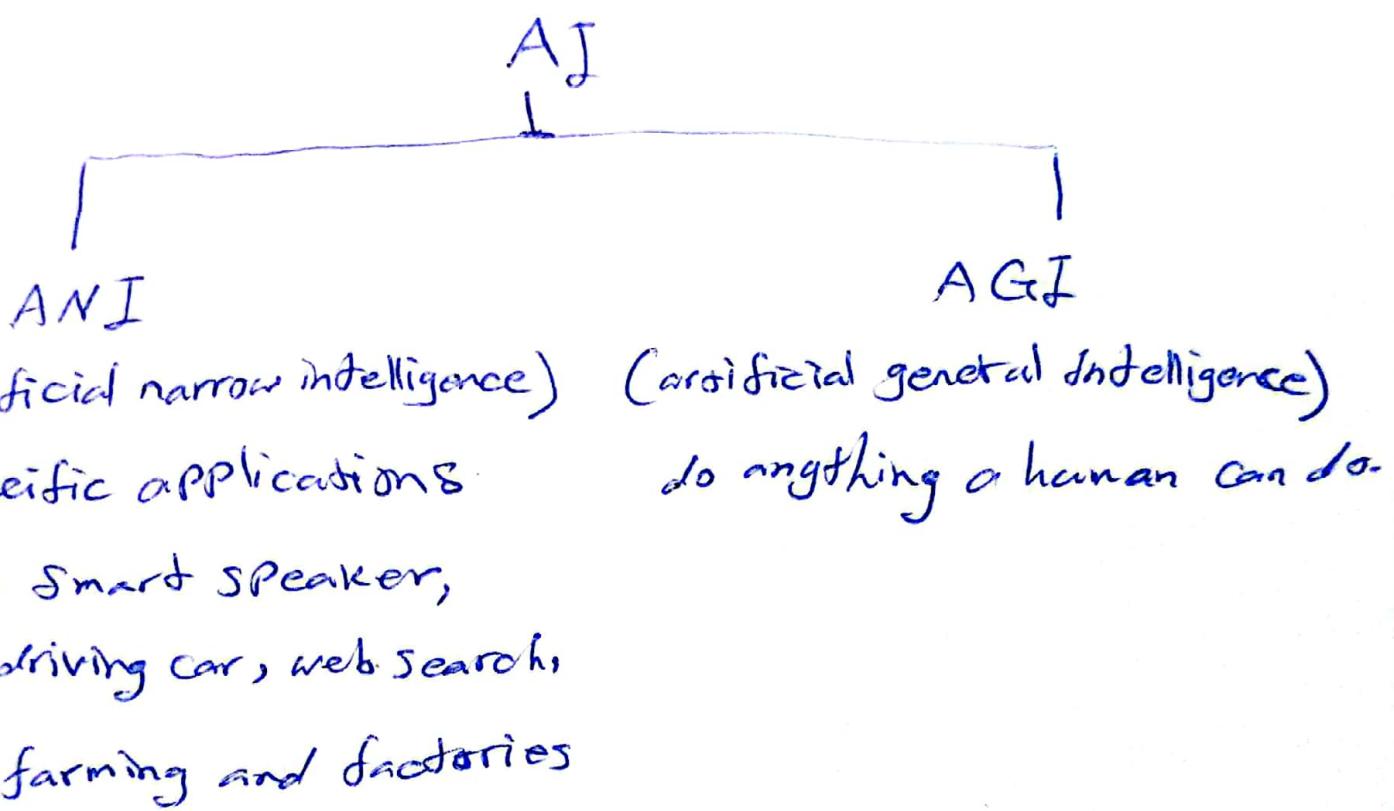
a4 = dense(a3, w4, b4)

f-X = a4

> return f-X

speculations on Artificial General

Intelligence (AGI)



The "one learning algorithm" hypothesis.

"vectorization"

for loops vs. vectorization

$x = np.array([[200, 17]])$ 2D array

$w = np.array([[-1, -3, 5], [-2, 4, -5]])$ (2×3) 2D array

$b = np.array([-1, 1, 2])$ (1×3) 2D array

def dense(A_in, W, B):

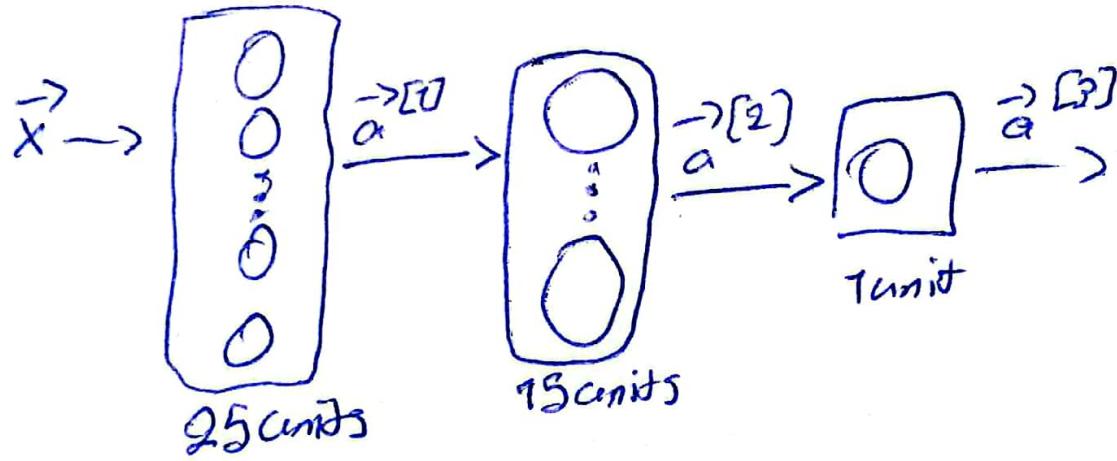
$$Z = np.matmul(A_{\text{in}}, W) + B$$

$$A_{\text{out}} = g(Z)$$

return A_out $\rightarrow [[1, 0, 1]]$

"Tensorflow implementation"

Image recognition examples



given set of (x, y) examples

How to build and train this?

[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=3, activation="sigmoid")])

① builds the layers

(15)

```

from tensorflow.keras.losses import
BinaryCrossentropy
model.compile(loss=BinaryCrossentropy())
model.fit(X, Y, epochs=100)

```

?] ② chooses the loss function
train the model

Training details

logistic Regression

$$\begin{aligned}
 I &= np.dot(w, x) + b && \text{define model} \quad ① \\
 f-x &= 1/(1+np.exp(-I)) && f_{\vec{w}, b}(\vec{x}) = ? \\
 \text{loss} &= -g \times np.log(f-x) - (1-y) \times np.log(1-f-x) && L(f_{\vec{w}, b}(\vec{x}), y) = ? \quad \begin{array}{l} \text{(logistic loss)} \\ \text{loss function} \end{array}
 \end{aligned}$$

$$\begin{aligned}
 w &= w - \alpha \times \frac{\partial J}{\partial w} && \text{minimize } J(\vec{w}, b) \quad ③ \\
 b &= b - \alpha \times \frac{\partial J}{\partial b} && (\text{GD})
 \end{aligned}$$

train a neural network in Tensorflow

model = Sequential([Dense(...), Dense(...), ..., Dense(...)])

model.compile(loss=BinaryCrossentropy()) ②

model.fit(X, y, epochs=100) ③

76

logistic loss called as BinaryCrossentropy in
Tensorflow.

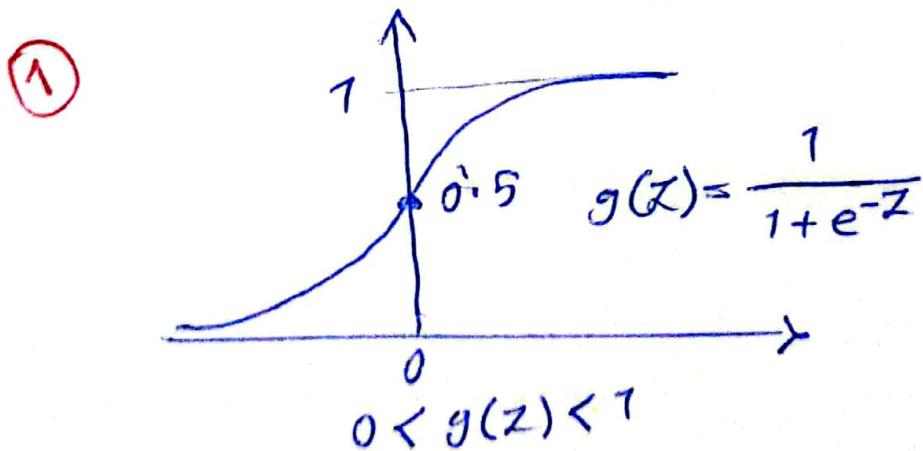
$$L(f(\bar{x}), y) = -y \log(f(\bar{x})) - (1-y) \log(1-f(\bar{x}))$$

mean squared error called as MeanSquaredError()
in Tensorflow (regression models)

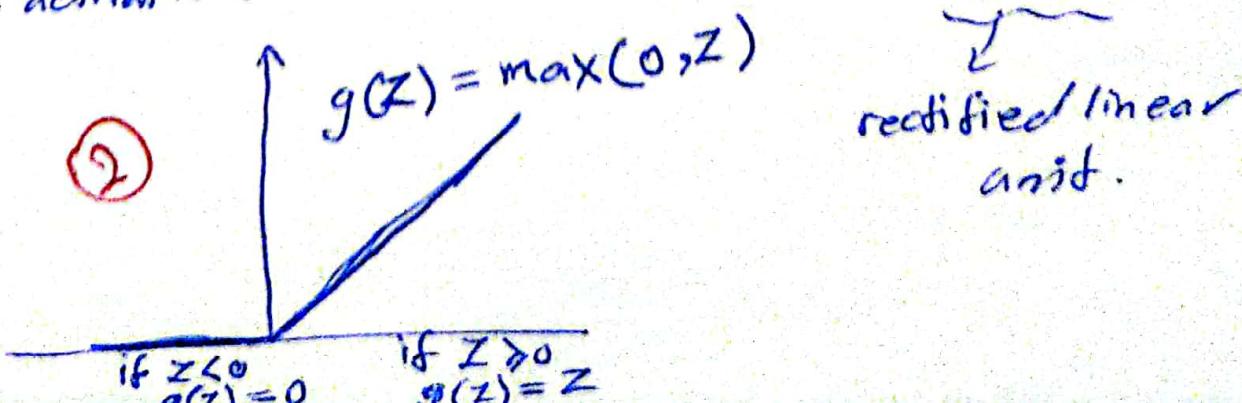
model.compile(loss=MeanSquaredError())

~~gradient descent called as~~
"Activation functions"

Sigmoid activation: $q_2^{[1]} = g(\vec{w}_2^{[1]} \cdot \vec{x} + b_2^{[1]})$



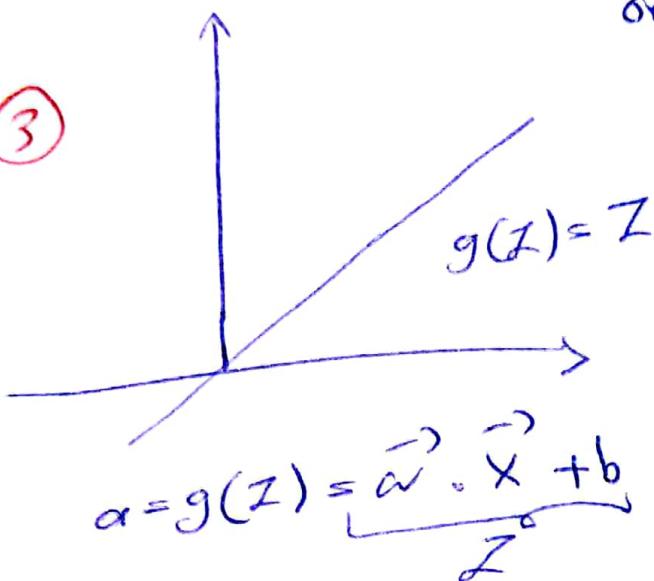
a common activation function in neural networks (ReLU)



linear activation function

or no activation function

③



softmax activation function ④

choosing activation function

- one can choose different activation functions in different layers in a network.

$$y = 0/1$$

outPut layer

classification
(Binary)

regression

Sigmoid

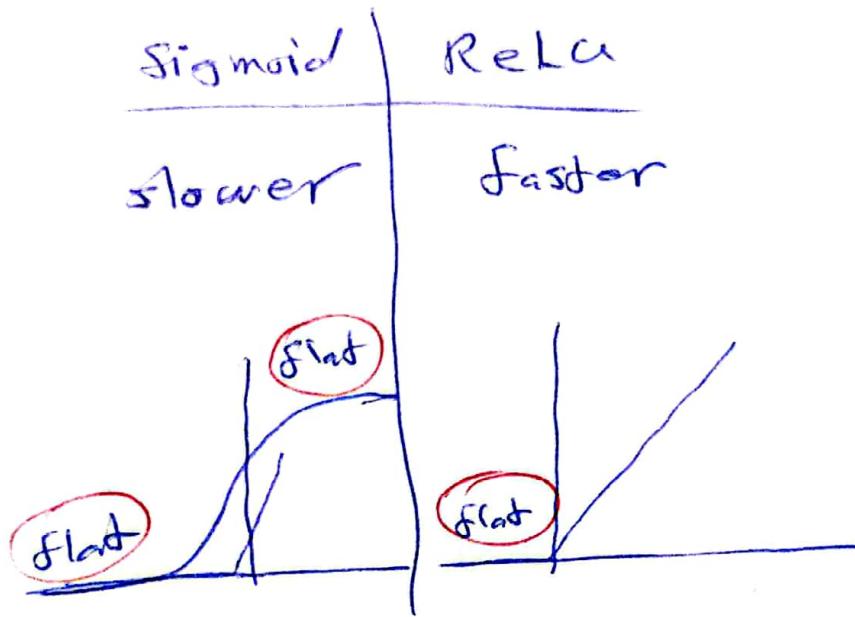
$$(Y = +/-)$$

χ -linear activation function

ReLU $Y = 0 \text{ or } +$

for example the price of a product can never be negative, so ReLu is the best choice here.

Hidden Layer: ReLU is by far the most common choice.



exp

model = Sequential([

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

Dense (units=75, activation="relu"),

Dense (units=1, activation="sigmoid")

])

"relu"
or
"linear"

if classification
else: relu or
linear

"why activation function"

$$\vec{a}^{[1]} = \vec{w}_1 \cdot \vec{x} + b_1^{[1]}$$

$$\vec{a}^{[2]} = \vec{w}_2 \cdot \vec{a}^{[1]} + b_2^{[2]}$$

$$= \vec{w}_2 \cdot (\vec{w}_1 \cdot \vec{x} + b_1^{[1]}) + b_2^{[2]}$$

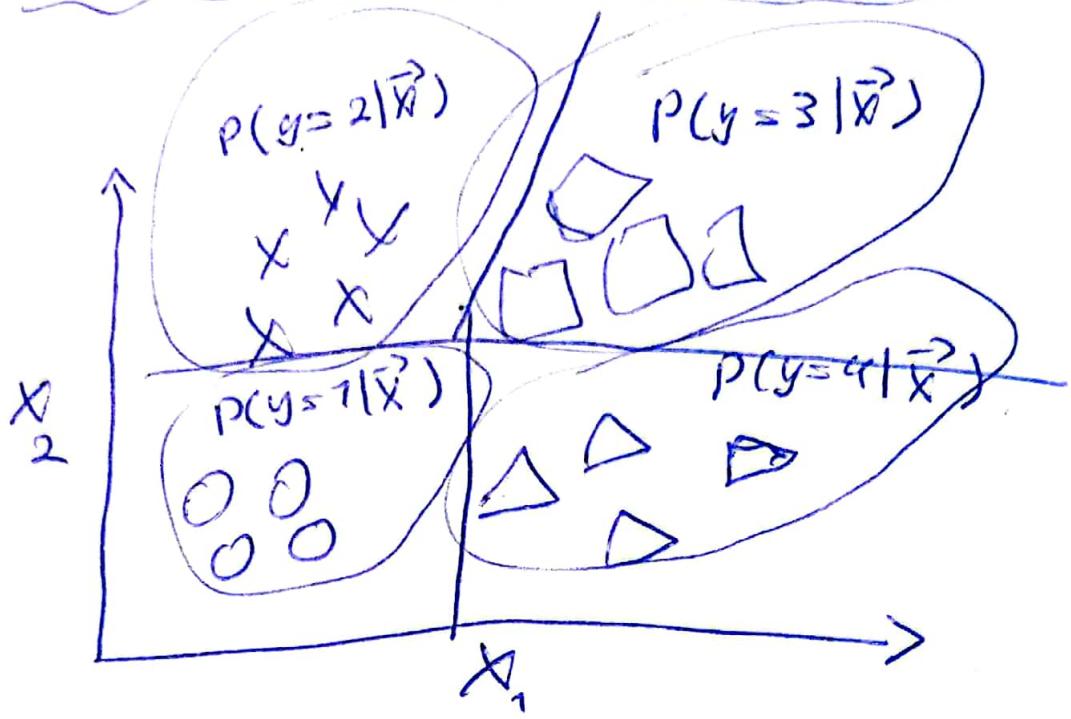
$$\vec{a}^{[2]} = \underbrace{(\vec{w}_2 \cdot \vec{w}_1)}_{\vec{w}} \cdot \vec{x} + \underbrace{\vec{w}_2 \cdot b_1^{[1]} + b_2^{[2]}}_b$$

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

$f(x) = \vec{w} \vec{x} + b$ linear regression

- if all layers of a neural network use the linear activation function the output is the same to the linear regression.
- if all hidden layers of a NN use the linear activation function and the output uses a sigmoid activation function, the output is equivalent to the logistic regression.
- Don't use linear activation in hidden layers
(use ReLu)

multiclass classification



"softmax regression algorithm"

- a generalization of logistic regression to the
multiclass classification.

logistic regression (2 possible output values)

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

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

$$\begin{cases} p(y=1|\vec{x}) \\ p(y=0|\vec{x}) \end{cases}$$

softmax regression (4 possible outputs)

$$X Z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

$$\square Z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$O Z_2 = \vec{w}_2 \cdot \vec{x} + b_2$$

$$\Delta Z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$\alpha_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=1|\vec{x})$$

$$\alpha_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=2|\vec{x})$$

$$\alpha_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=3|\vec{x})$$

$$\alpha_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=4|\vec{x})$$

softmax regression (N possible outputs)

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

$$\vec{w} = w_1, w_2, \dots, w_N$$

$$\vec{b} = b_1, b_2, \dots, b_N$$

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

$$\text{note: } \alpha_1 + \alpha_2 + \dots + \alpha_N = 1$$

cost for softmax function

logistic regression

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

$$\alpha_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x})$$

$$= P(y=0|\vec{x})$$

$$\alpha_2 = 1 - \alpha_1$$

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

if $y=1$ $\underbrace{a_1}_{a_2}$
 if $y=0$

softmax regression

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

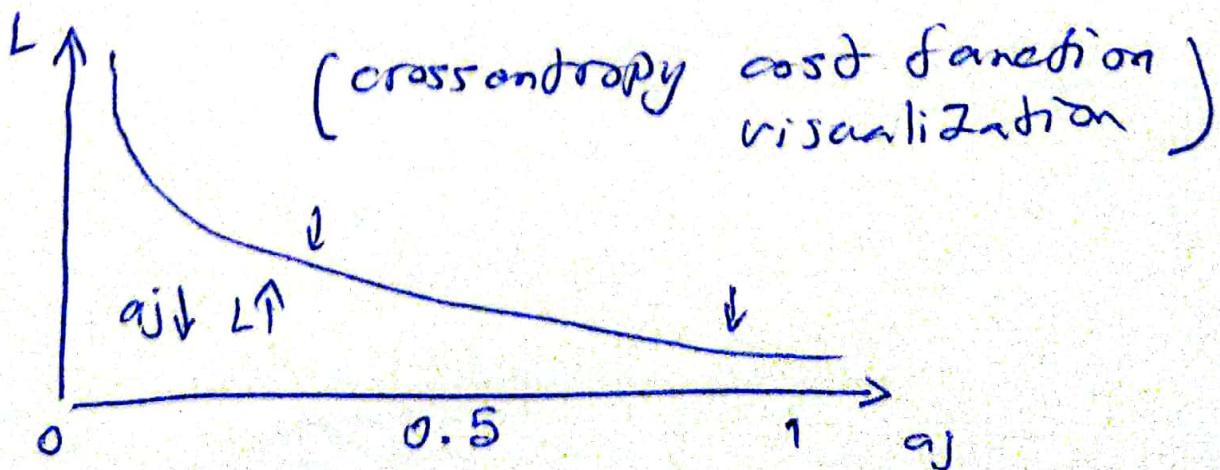
⋮

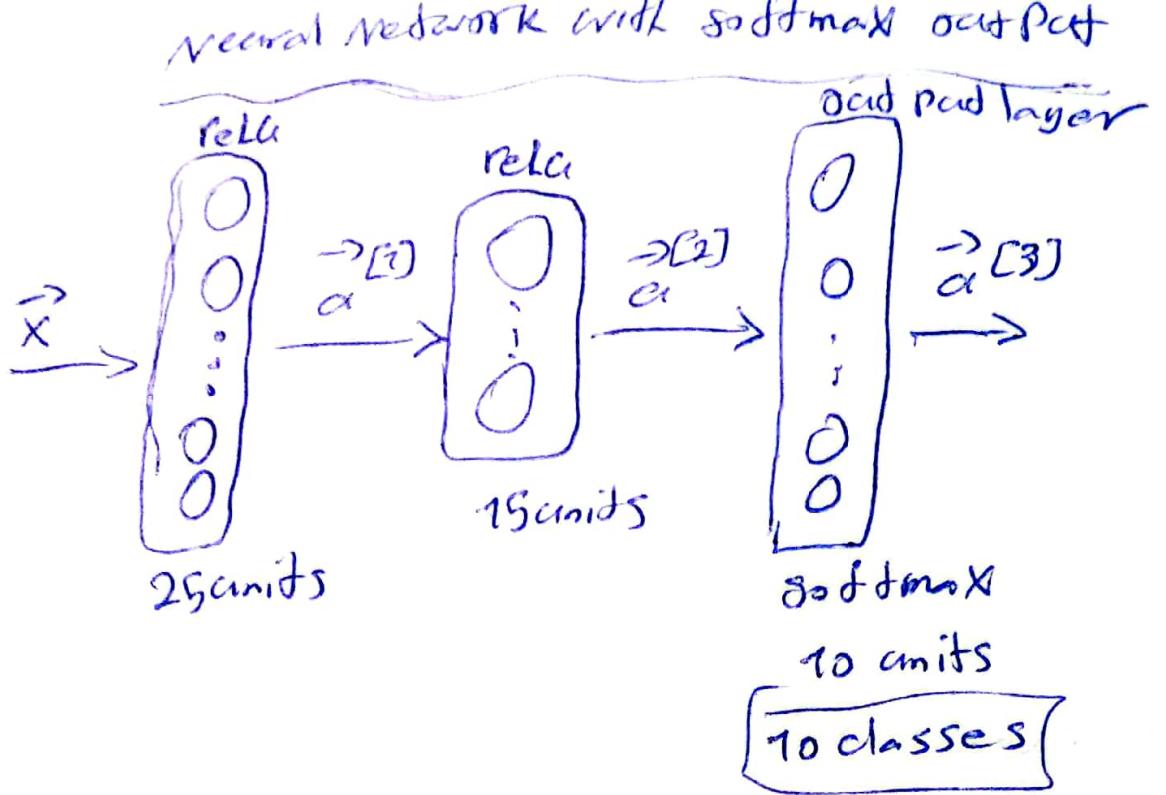
$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y=N | \vec{x})$$

crossentropy

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

$$\text{loss} = -\log a_j \text{ if } y=j$$





$$\stackrel{(2)}{z_1} = \stackrel{(1)}{\alpha_1} \cdot \alpha + b_1 \quad \stackrel{(3)}{z_1} = \stackrel{(2)}{\alpha} \cdot \alpha + b_1$$

$$\stackrel{(3)}{\alpha_1} = \frac{\alpha^{z_1}}{e^{z_1} + \dots + e^{z_{10}}} = P(y=1|\vec{x})$$

$$\stackrel{(2)}{z_{10}} = \stackrel{(1)}{\alpha_{10}} \cdot \alpha + b_{10} \quad \stackrel{(3)}{z_{10}} = \stackrel{(2)}{\alpha} \cdot \alpha + b_{10}$$

$$\stackrel{(3)}{\alpha_{10}} = \frac{\alpha^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}} = P(y=10|\vec{x})$$

$$\stackrel{(3)}{\alpha} = (\alpha_1^{(3)}, \dots, \alpha_{10}^{(3)}) = g(z_1^{(3)}, \dots, z_{10}^{(3)})$$

softmax in tensor flow

import tensorflow as tf

from tensorflow.keras import Sequential

from tensorflow.keras.layers import Dense

model = Sequential([

Dense(25, activation="relu"),

Dense(15, activation="relu"),

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

J)

from tensorflow.keras.losses import

SparseCategoricalCrossentropy

model.compile(loss=SparseCategoricalCrossentropy())

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

improved implementation of softmax

more numerically accurate implementation of logistic
losses

model = Sequential([

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

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

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

linear

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

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

train logits model(X)

softmax $\delta - x = \text{tf.nn.sigmoid}(\text{logit})$

Prediction

model = Sequential([

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

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

Dense(units=10, activation="softmax") J)

model.compile(loss='SparseCategoricalCrossentropy', from_logits=True) 25

model.fit(X, Y, epochs=100) → train

logits = model(X) → predict
↑
not a_1, \dots, a_{10}
is z_1, \dots, z_{10}

$f(x) = \text{tf.nn.softmax}(\text{logits})$

multilabel classification

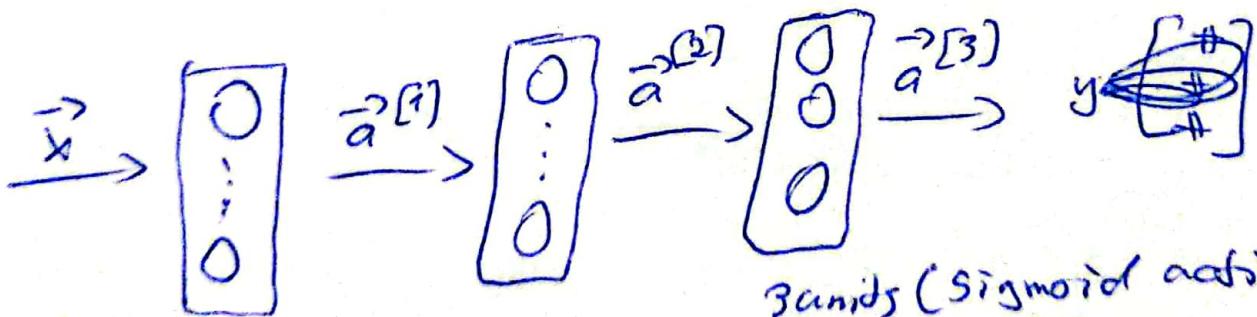
(exp)

in a picture

	<u>Image 1</u>	<u>Image 2</u>	<u>Image 3</u>
is there a car?	Yes	No	Yes
is there a bus?	No $y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$	No Yes $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$	Yes Yes $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$
is there a pedestrian?	Yes	Yes	No

In this case train a model with three outputs:

outPut



3 units (sigmoid activation)

outPut → $a^3 = \begin{bmatrix} a_1^{(3)} \\ a_2^{(3)} \\ a_3^{(3)} \end{bmatrix}$

car
bus
pedestrian

Advanced Optimization

Gradient descent

$$w_j = \arg - \alpha \frac{d}{dw_j} j(\vec{w}, b)$$



"Adam" algorithm: adjust the learning rate automatically.

Adam: Adaptive moment estimation

not just one alpha,

(exp)

different alpha

$$w_1 = w_1 - \alpha_1 \frac{d}{dw_1} j(\vec{w}, b)$$

:

$$w_{10} = w_{10} - \alpha_{10} \frac{d}{dw_{10}} j(\vec{w}, b)$$

$$b = b - \alpha_{11} \frac{d}{db} j(\vec{w}, b)$$

if w_j (or b) keeps moving in same direction,

increase α_j

if w_j (or b) keeps oscillating, reduce α_j

model = Sequential [

Dense (units=25, activation="sigmoid"),

Dense (units=15, activation="sigmoid"),

Dense (units=10, activation="linear").

])

model.compile(optimizer=tf.keras.optimizers.

Adam(learning_rate=1e⁻³), loss=

SparseCategoricalCrossentropy(from_logits=True))

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.

convolutional layers { faster computation
need less training data
(less prone to overfitting)

derivative

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

Say $w = 3 \quad J(w) = 3^2 = 9$

If we increase w by a tiny amount $\epsilon = 0.001$

How does $J(w)$ change?

$w = 3 + 0.001 \quad J(w) = w^2 = 9.006001$

if $w \uparrow \epsilon$

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

$$\frac{d}{dw} j(w) = \delta$$

if formal definition of derivative

if $w \uparrow \epsilon$ causes $J(w) \uparrow K \times \epsilon$ then

$$\frac{d}{dw} j(w) = K$$

Gradient Descent

repeat {

$$w_j = w_j - \alpha \frac{d}{dw_j} j(\hat{w}, b)$$

}

(29)

if derivative is small, then this update step will make a small update to w_j and ~~even~~ if it is large,

then the update step makes a ~~big~~ large update to w_j .

$$w=2 \rightarrow J(w) = w^2 = 4 \rightarrow w \uparrow 0.001$$

$$J(w) = J(2.001) = 2.000001$$

$$J(a) \uparrow \underbrace{6 \times 0.001}_{\downarrow}$$

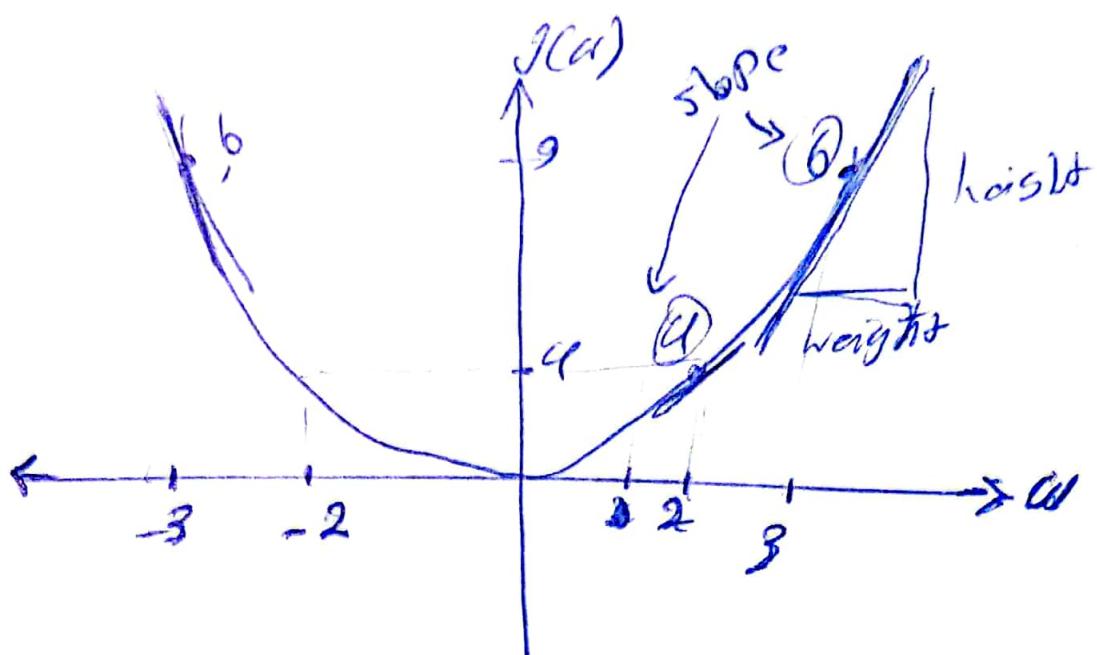
$$\frac{d}{dw} J(w) = ④$$

$$a = -3 \quad J(w) = w^2 = 9 \quad w \uparrow 0.001$$

$$J(w) = J(-2.999) = 8.999001$$

$$J(w) \uparrow \underbrace{6 \times 0.001}_{-0.006}$$

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



$$\frac{d}{dw} J(w) = 2w$$

import SymPy

$J, w = \text{SymPy.symbols('}j\}, w\})$

$$J = w^2$$

calculate derivative
with

$$\frac{dJ}{dw} = \text{SymPy.diff}(J, w) = 2w \quad \text{SymPy.}$$

calculate the derivative of J

with respect to w .

$$\text{dJ-dar.subs}(L(w, 2)) = 4$$

plug a number in

the derivative

if $j(w)$ is a function of one variable (w),

$$\frac{d}{dw} j(w) \rightarrow d \rightarrow \text{derivative notation}$$

if $j(w_1, w_2, \dots, w_n)$ is a function of more than one variables,

$$\frac{\partial}{\partial w_i} j(w_1, w_2, \dots, w_n) \rightarrow \partial \rightarrow \text{partial derivative notation}$$

"computation graph"

$$\vec{x} \rightarrow \boxed{0} \xrightarrow{a}$$

$$\begin{aligned} &\text{linear activation} \\ &\alpha = g(z) = z \\ &\alpha = \alpha x + b \end{aligned}$$

$$J(a, b) = \frac{1}{2}(a - y)^2$$

$$\text{if } w=2, b=8, x=-2, y=2$$

$$\begin{aligned} \alpha &\rightarrow C = \alpha x \xrightarrow{-4} [a = c + b] \xrightarrow{-4} d = a - y \\ &\quad \begin{matrix} -4 \\ 8 \\ b \end{matrix} \end{aligned}$$

$$2 \rightarrow J = \frac{1}{2} d^2 \xrightarrow{2} J$$

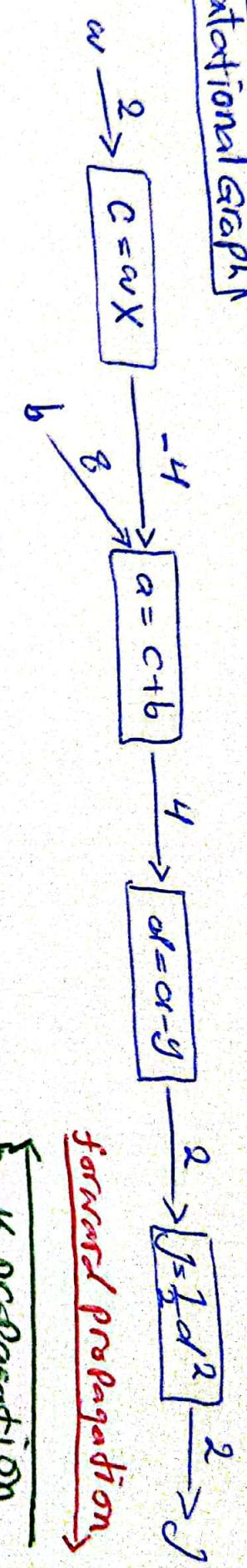
computation graph

(32)

$$w=2 \quad b=8 \quad x=-2 \quad y=2 \quad J = \frac{1}{2} (a - y)^2$$

~~forwards~~

Computational Graph



forward propagation
back propagation

Debugging a learning algorithm

(33)

if you implemented regularized linear regression.

on a dataset

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

but it makes unacceptably large errors in predictions. what do you try next?

- Get more training example
- Try smaller set of features
- Try getting additional features
- Try adding Polynomial features ($x_1^2, x_2^2, x_1 x_2$, etc)
- Try decreasing λ
- Try increasing λ

Machine learning diagnostic

- A test which you can run to gain insight into what is / is not working with a learning algorithm, to gain guidance into improving its performance.

Evaluating a model

(34)

1) Train and test set: split data into train and test sets.

compute the accuracy of the ^{Prediction coming} test set.

2) model selection

$$1. \hat{f}_{\bar{w}, b}(\bar{x}) = w_1 x + b$$

$$2. \hat{f}_{\bar{w}, b}(\bar{x}) = w_1 x + w_2 x^2 + b$$

$$3. \hat{f}_{\bar{w}, b}(\bar{x}) = w_1 x + w_2 x^2 + w_3 x^3 + b$$

$$\dots \\ 10. \hat{f}_{\bar{w}, b}(\bar{x}) = w_1 x + w_2 x^2 + \dots + w_{10} x^{10} + b$$

try different models and evaluate them and choose the one ~~which~~ which gives you the best performance.

(35)

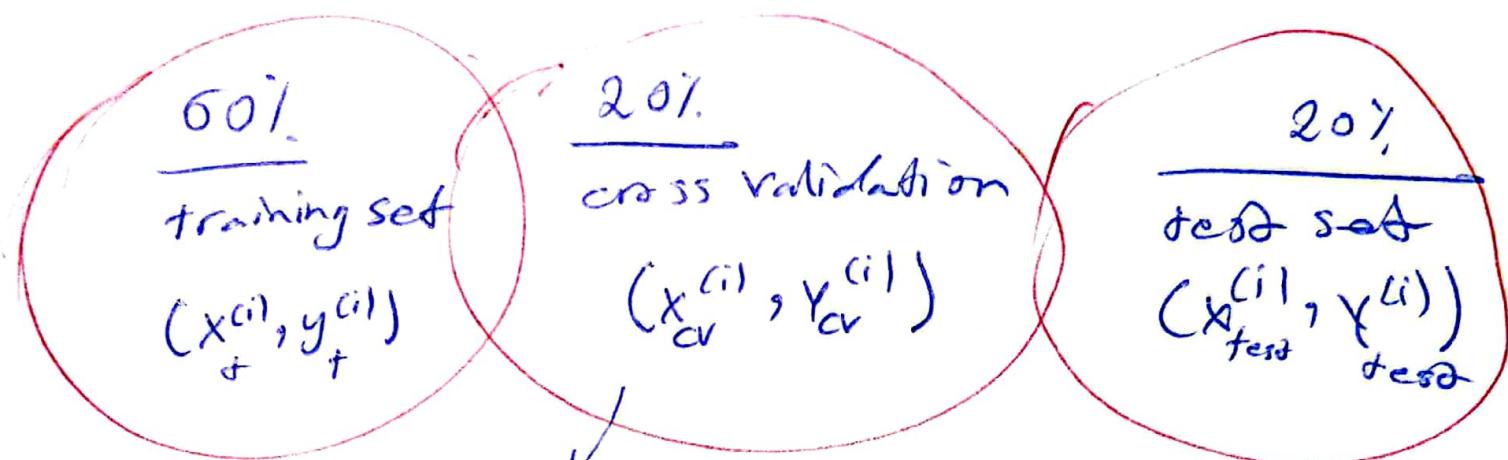
(Validation set, development set)



development set

③ Training / cross validation / test set.

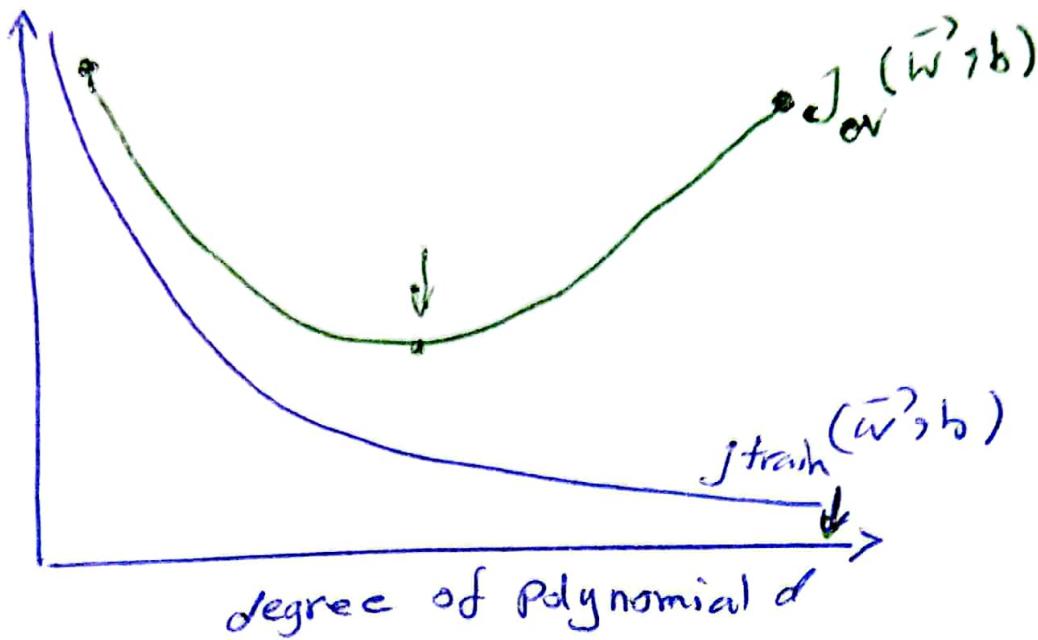
(exp)



use this set to ~~evaluate~~
select the best model
(model selection) and then
use the test set to evaluate
the performance of the
model.

Bias/Variance

- High bias (underfit) $\rightarrow J_{train}$ is high \wedge J_{cv} is high \wedge underfitted
- High variance (overfit) $\rightarrow J_{train}$ is low \wedge J_{cv} is high \wedge overfitted
- right model: J_{train} is low \wedge J_{cv} is low



- How to detect High bias or High Variance of model?

High bias (underfit)

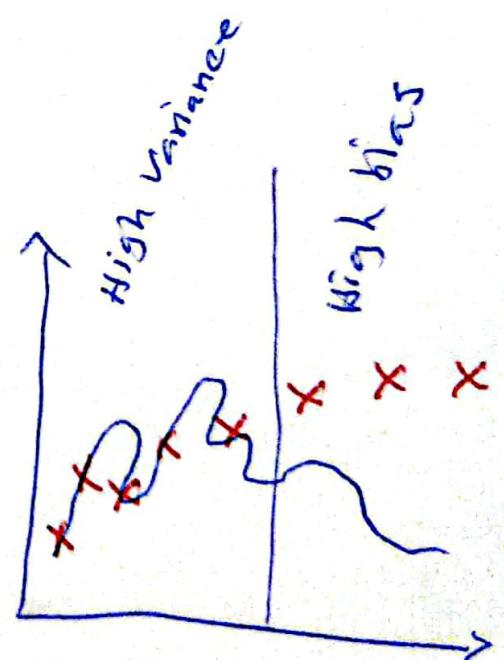
J_{train} will be high

$$(J_{\text{train}} \approx J_{\text{cv}})$$

High variance(overfit)

$$J_{\text{cv}} \gg J_{\text{train}}$$

$$(J_{\text{train}} \text{ may be low})$$



High bias and High variance

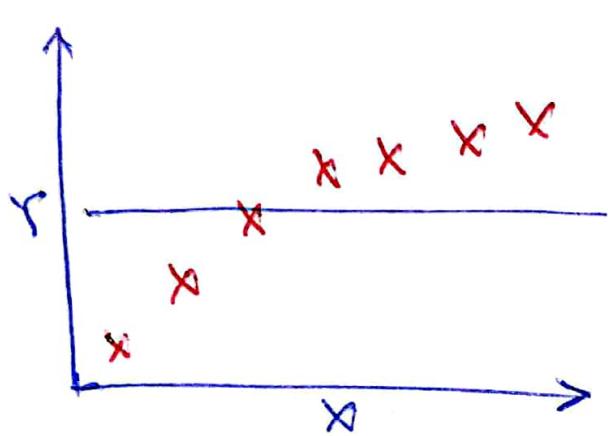
J_{train} will be high. and $J_{\text{cv}} \gg J_{\text{train}}$

regularization and bias/Variance

(exp)

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

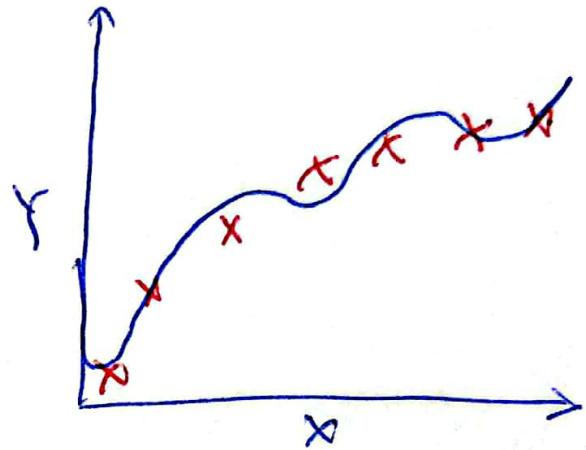
$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$



$\lambda = \text{very high (10000)}$

(High bias)

$J_{\text{train}}(\vec{w}, b)$ is large

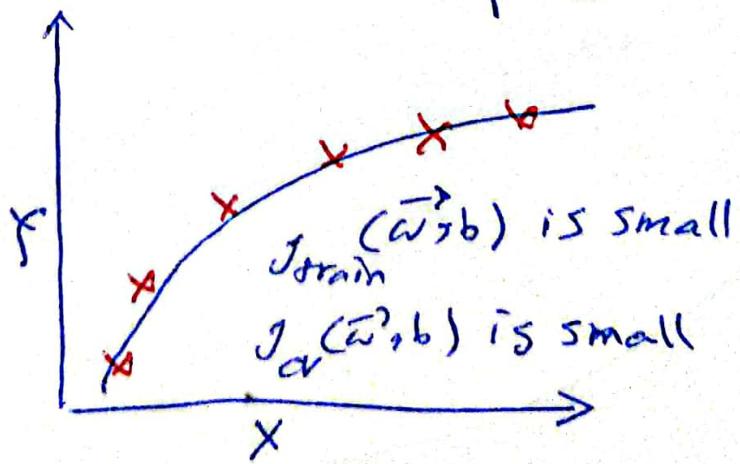


small $\lambda (\lambda = 0)$

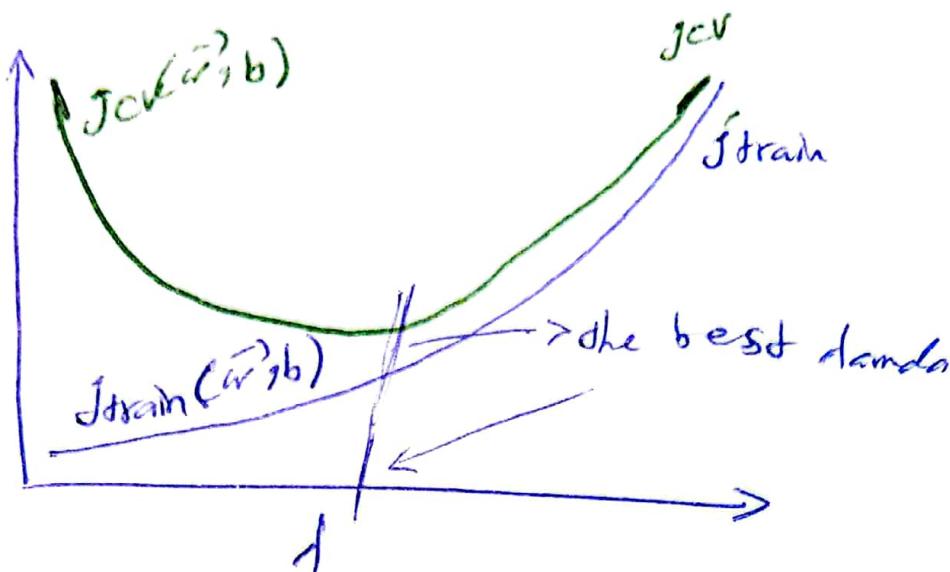
$J_{\text{train}}(\vec{w}, b)$ is small.

$J_{\text{cv}}(\vec{w}, b)$ is large.

(High variance)



Intermediate λ



Baseline performance

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

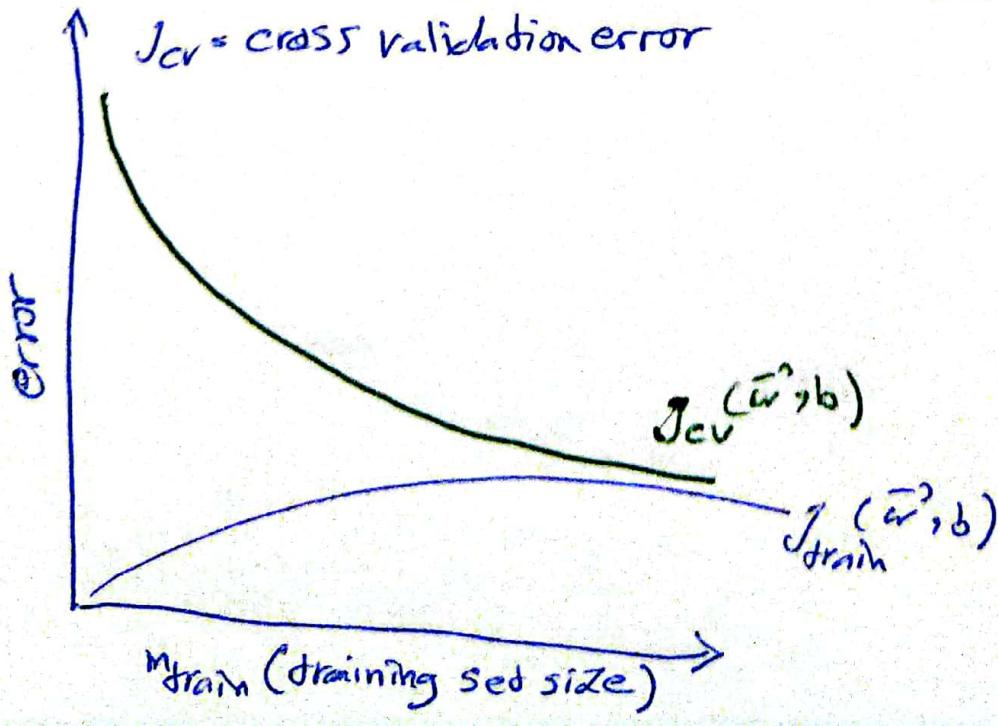
- Human level Performance.
- Competing algorithms Performance.
- Guess based on experience.

learning curves

J_{train} = training error

J_{cv} = cross validation error

$$\text{model: } f_{\bar{w}, b}(x) = c_1 x + c_2 \bar{x} + b$$



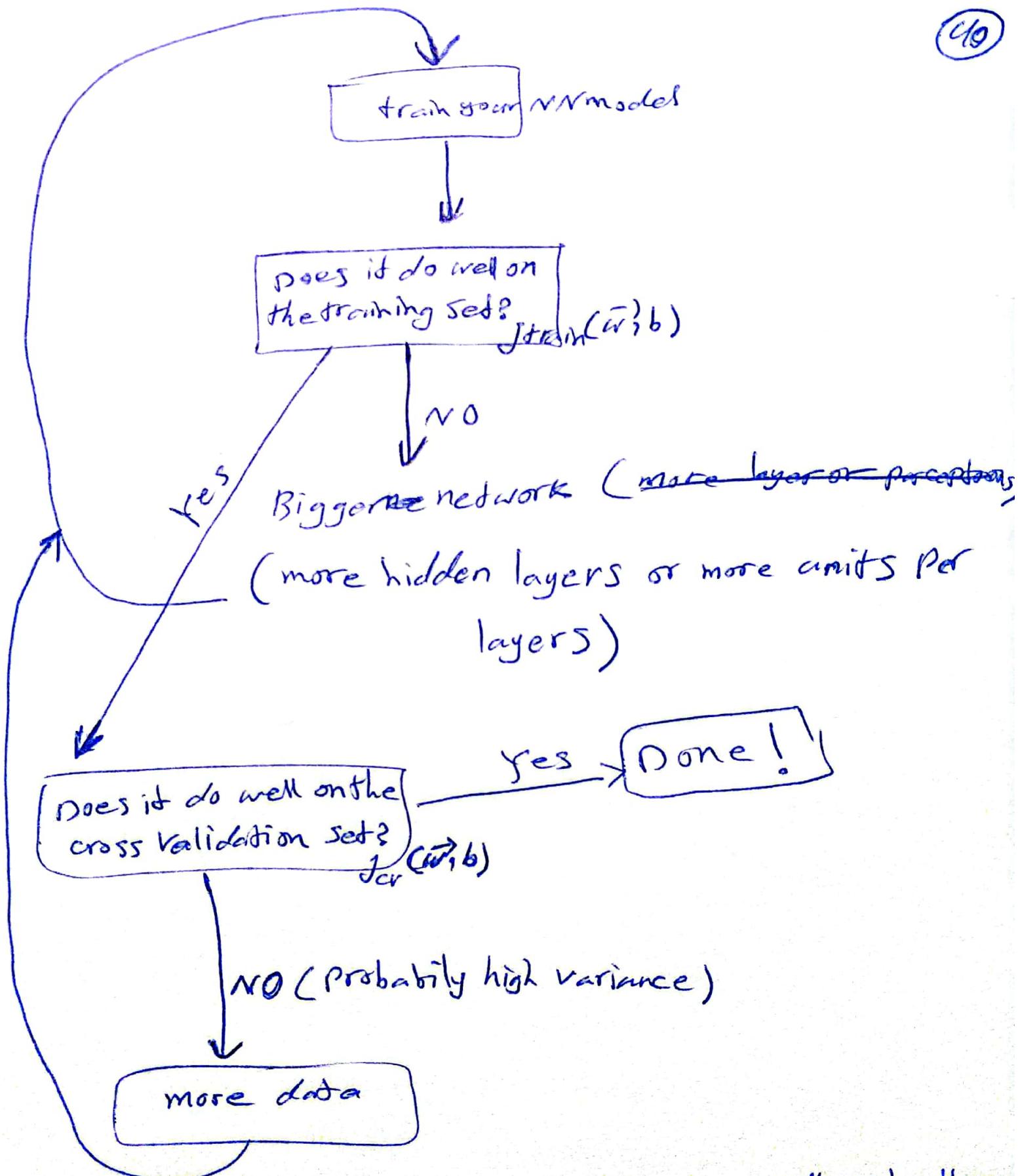
- if a learning algorithm suffers from high variance (overfit), getting more training data is likely to help.
 - if a learning algorithm suffers from high bias (under fit) getting more training data does not help.
-

if a model get a very high error in predictions. what do you try next?

- Get more training example (~~high bias~~)
 - Try smaller sets of features (~~High variance~~)
 - Try getting additional features (fixes High bias)
 - Try adding Polynomial features (x_1^2, x_2^2, \dots) (fixes High bias)
 - Try decreasing λ (fixes high bias)
 - Try increasing λ (fixes high variance)
-

bias/variance and NNs.

- large neural networks are low bias machines, means if you make your model large enough, at the most cases you get a well fitted model.

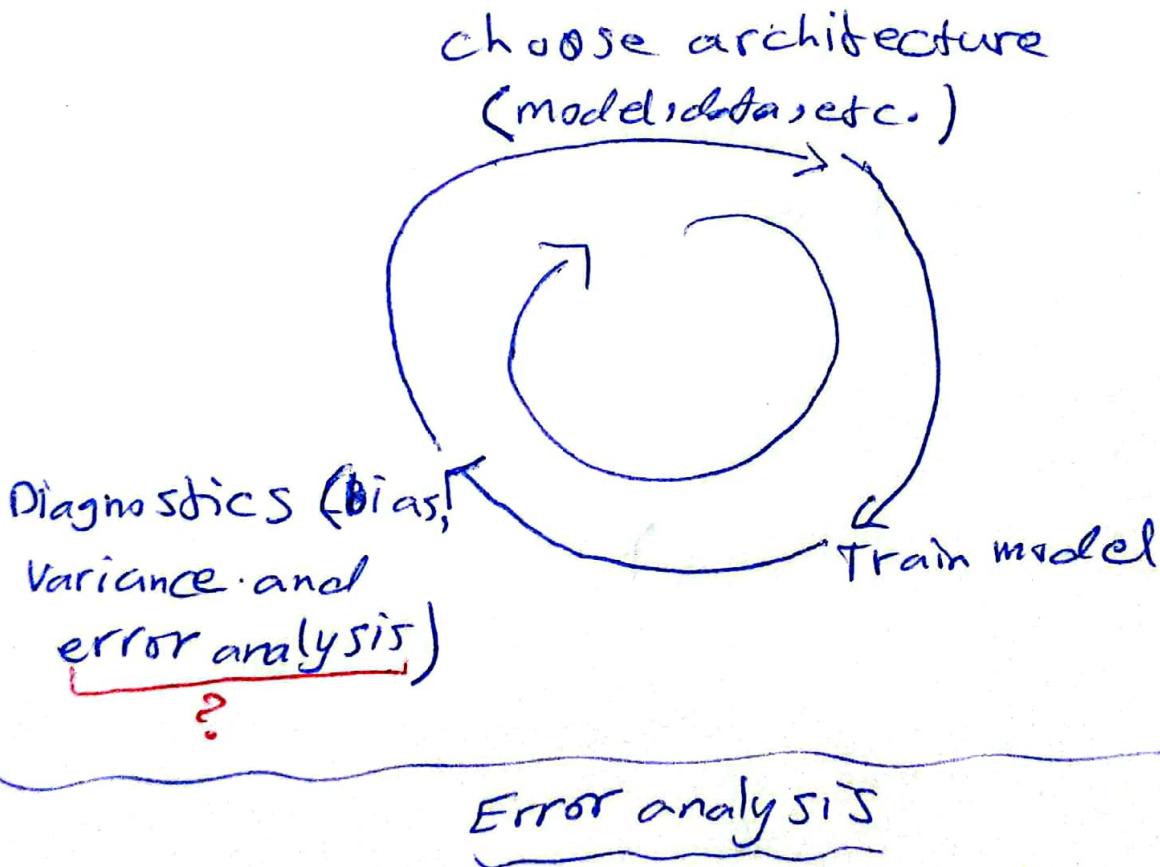


- A large neural network will usually do as well or better than a smaller one so long as regularization is chosen appropriately.

regularized model in tensorflow

layer = Dense(units=10, activation="relu",
kernel_regularizer=L2(0.01))

I Iterative loop of ML development



Error analysis

- Try to analyse misclassified examples (manually collect them and analyse them and try to make decisions based on them to improve your the performance of your model.

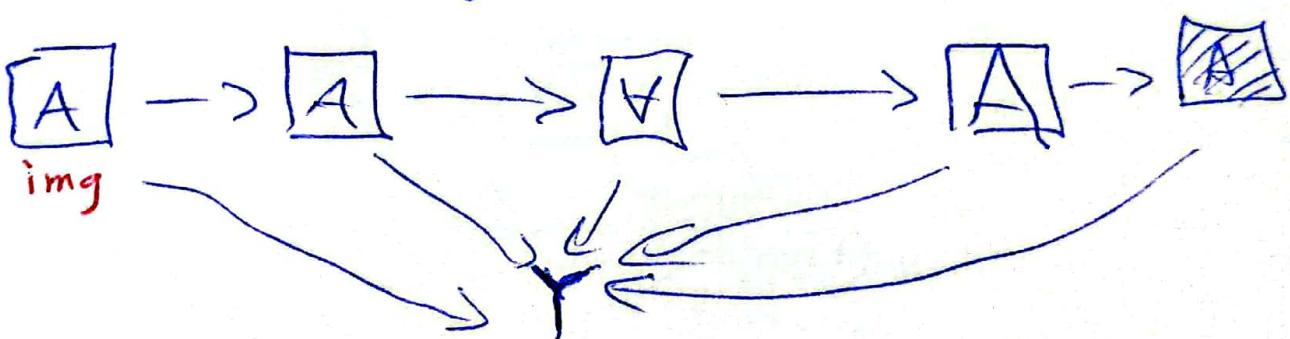
Adding data

- Add more data of everything. E.g., "Honeypot" project.
- Add more data of the types where error analysis has indicated it might help.
- Beyond getting brand new training examples (X, y) , another technique: Data augmentation

?

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

exp.



- Data synthesis: using artificial data inputs to create a new training example.

Engineering the data used by your system

conventional model-centric
approaches

AI = Code + Data
model
work on this

Data-centric approaches

AI = code + Data
data
work on this

Transfer-learning

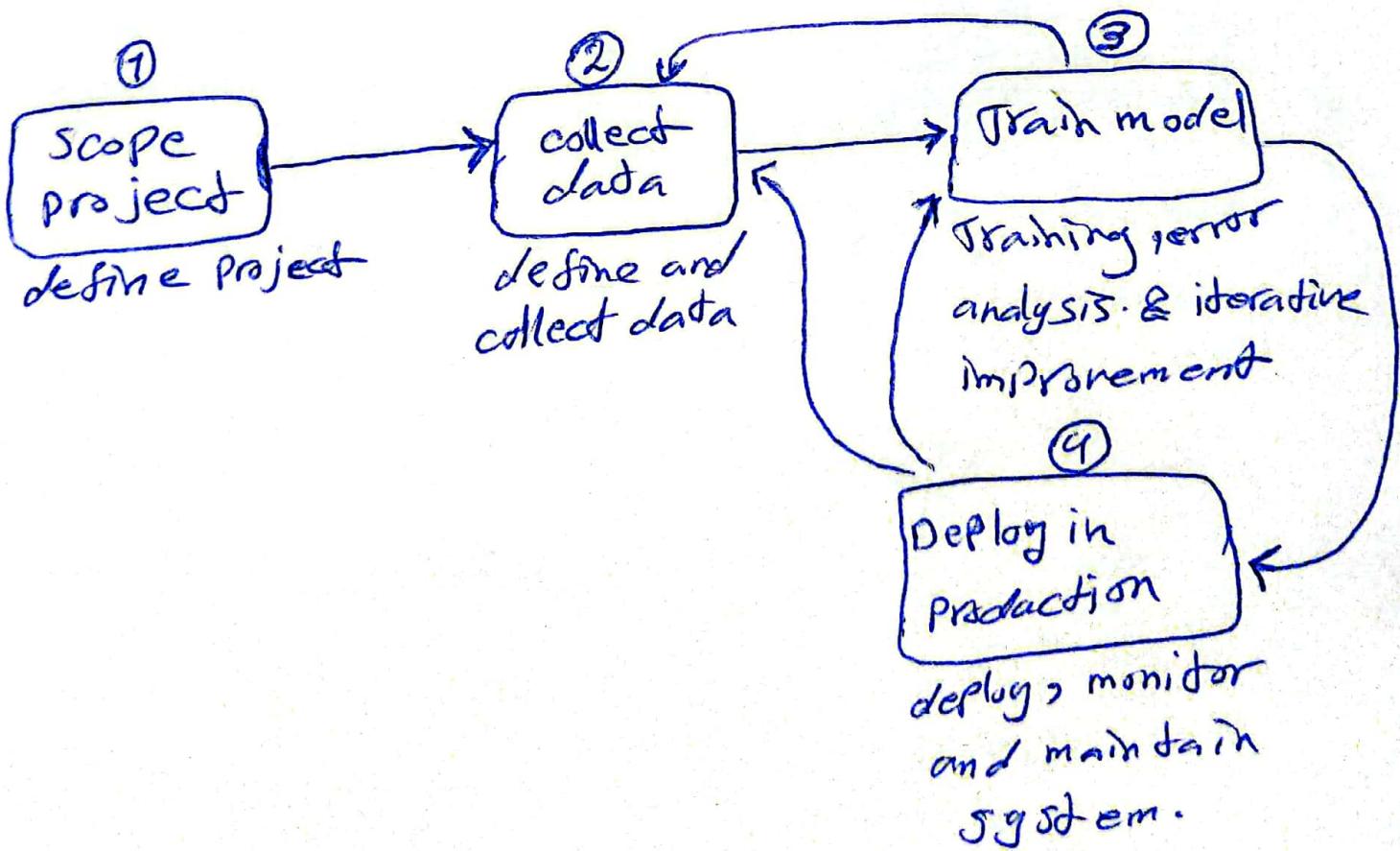
43

-using data from a different task.

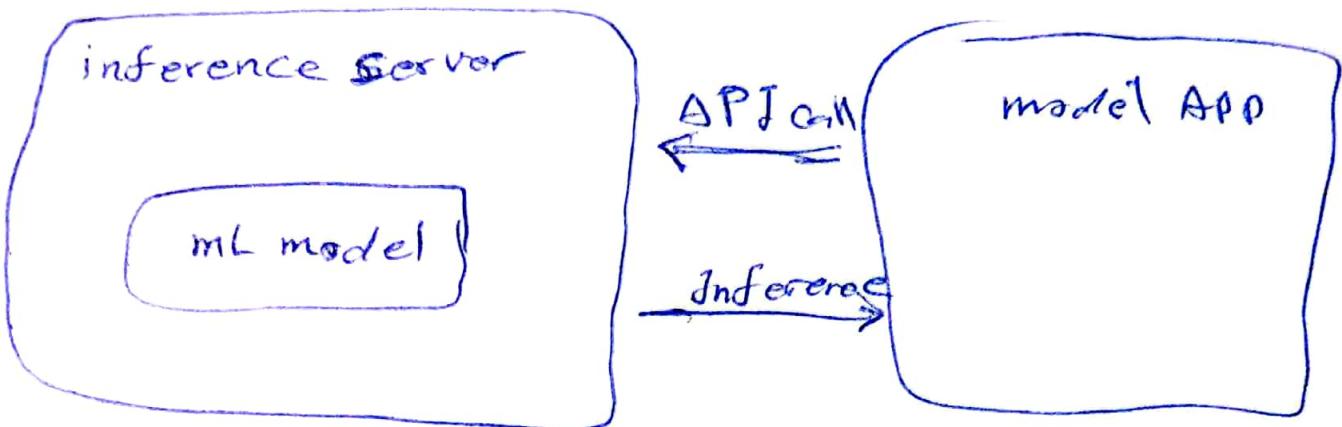
steps

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

full cycle of a ML project



Deployment



- Software engineering may be needed for

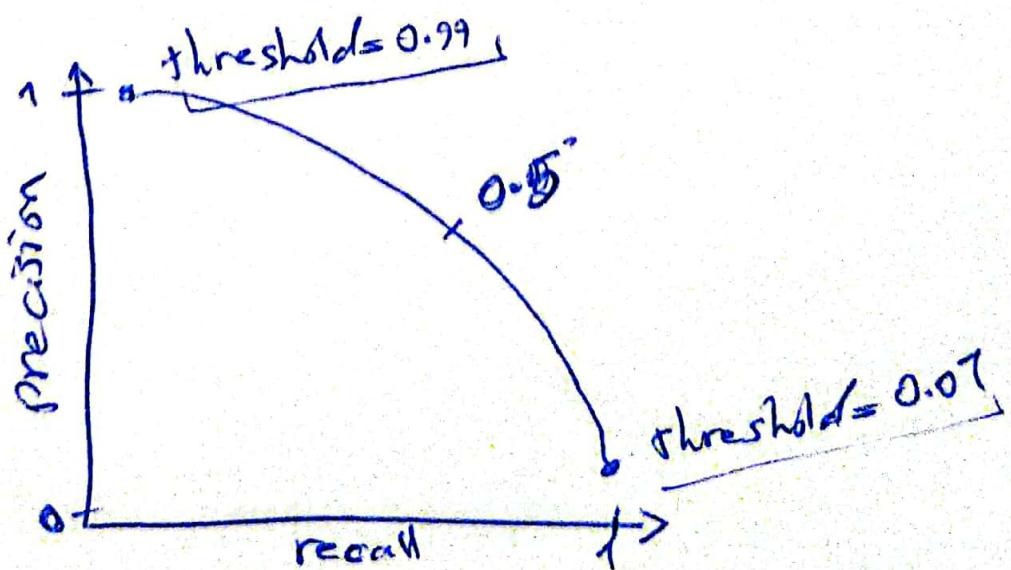
Precision/Recall

		1	0
1	True positive	False positive	
	False negative	True negative	

$$\text{Precision} = \frac{\text{True positive}}{\text{Predicted positive}}$$

$$= \frac{\text{True positive}}{\text{True-pos} + \text{false-pos}}$$

$$\text{recall} = \frac{\text{True_pos}}{\text{actual_pos}} = \frac{\text{True_pos}}{\text{True_pos} + \text{false_neg}}$$



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

P = precision

R = recall

Harmonic mean

(95)