

- Supervised learning
 - regression

Machine learning (Defⁿ)

H₀: Null hypothesis (applying some tests)

H₁: Alternate hypothesis

→ LOOCV, hold-out, k-fold.

29 Jan 2024

LAB

- + uciml repos
- x wget
- x RFE
- * Random Over Sampling , SMOTE

→ normalizing and not normalizing
 → oversampling and ↓ oversampling

30 Jan 2024

feature pool → reduced pool
 (features matrix)

Classification

- pre-determined classes exist.
- we need to find a function f (hypothesis) which maps

Regression

whenever we try to fit a model and we fix some kind of condition (like using only depth = 5 for decision trees) this results in some kind of bias.

linear regression → residual error

hypothesis: given some θ

$h_{\theta}(x)$
 Hypothesis
 ↗ parameter
 ↗ label

slope
 $y = mx + c \leftarrow$ intercept
 $f(x) + b_0$

weight → $w_n + b \leftarrow$ bias

weights and bias are parameters of the model and we are interested in learning these parameters.

Inductive Bias

absolute

$$\text{error} = |h_0(\hat{x}) - y|$$

$$\Rightarrow |\hat{y} - y|$$

$$E \text{ (Lagrangian error)} = (\hat{y} - y)^2 \quad (\text{Because } |\hat{y} - y| \text{ is not differentiable at 0})$$

learning is acquiring general concepts from specific training examples.

concept learning if we have a set of examples we want to guess the label of an object using some attributes.

Task: Task is to learn to predict the value of boolean valued attribute (family car). based on the values of its other attributes.

Training set X .

Book: Tom mitchell.

$$X = \{x^t, r^t\}_{t=1}^N$$

$$r_t = \begin{cases} 1 & \text{if } x \text{ is +ve} \\ 0 & \text{if } x \text{ is -ve} \end{cases}$$

$H \rightarrow$ set of all legal hypothesis.

\uparrow
some hypothesis

$h \in H$
 \uparrow
hypothesis concept.
(approximate)

hypothesis is
consistent.

concept learning. Inferring a boolean-valued function from training examples of its input and output.

$$|H| = |X| = 2^N \quad (\text{size of instance space})$$

x_1	x_2	h_0	h_1	h_2	h_3
0	0	0	0	0	0
0	1	0	0	0	0
1	0	0	0	1	1
1	1	0	1	0	1

legal
No. of hypotheses $\rightarrow 2^{|H|}$

? → any value is acceptable. / don't care
 specify required value for the attribute.
 "∅" → no value is acceptable. [any ∅ leads to -ve example].

most general hypothesis → $\langle ?, ?, ?, ?, ? \rangle$
 most specific hypothesis → $\langle \emptyset, \emptyset, \emptyset, \emptyset \rangle$

1st feb 1994

pg - 27

as long as we assume that the hypotheses space H contains a hypothesis that describes the true target concept c and that the training

→ returns true (label) everytime
 → finds only 1 **consistent** hypothesis.
 Disadvantage of FindS → there can be some records having noise or errors.
 inconsistent sets of training examples

conjunction of monomials → {sky = sunny} \wedge {AirTemp = warm} \wedge ...

syntactically distinct hypothesis.

size of hypothesis space .

$$\begin{array}{cc} F1 & F2 \\ A & X \\ B & Y \end{array} \Rightarrow 2 \times 2$$

$$3 \cdot 2 \cdot 2 \cdot 2 \cdot 2 = 96$$

$$5 \cdot 4 \cdot 4 \cdot 4 \cdot 4 = 5120$$

distinct instances

syntactically distinct hypothesis

$$3 \rightarrow 3 + (2)$$

because of ? and ∅.

of syntactically distinct hypothesis .

$$\begin{aligned} (\text{?}, \text{X}), (\text{?}, \text{Y}), (\emptyset, \text{X}), (\emptyset, \text{Y}), (\text{A}, \text{?}), (\text{B}, \text{?}), (\text{?}, \text{?}), (\emptyset, \emptyset), \\ (\text{A}, \emptyset), (\text{B}, \emptyset), (\text{A}, \text{X}), (\text{A}, \text{Y}), (\text{B}, \text{X}), (\text{B}, \text{Y}), (\text{?}, \emptyset), (\emptyset, \emptyset) \end{aligned}$$

semantically distinct hypothesis $\Rightarrow 1 + (4 \times 3 \times 3 \times 3 \times 3 \times 3)$

Most specific remains

because this $\langle \emptyset, \emptyset \rangle$ represents all the -ve instances.

$$3 \times 3 \times 2 \times 2 \times 2 \Rightarrow$$

$$(syntactic) \quad 5 \times 5 \times 4 \times 4 \times 4 \Rightarrow$$

$$(semantic) \Rightarrow 1 + (4 \times 4 \times 3 \times 3 \times 3) \Rightarrow$$

consistent : A hypothesis h is consistent with a set of training examples D if and only if $h(x) = c(x)$ for each example $\langle x, c(x) \rangle$ in D .

consistent

BOOK example .

since ? matches with A and x matches with X , \therefore the hypothesis predicts this instance is true and the given class is also true.

since ? matches with B and x does not matches with Y , \therefore the hypothesis predicts this instance as -ve but the given class is +ve. \therefore inconsistent .

5th Feb '24

Candidate elimination algorithm

version space .

\rightarrow consistent

Ques check whether $H_1 = \langle ?, ?, \text{No}, ?, \text{Many} \rangle$ is consistent or not?

Example citations size in library price editor

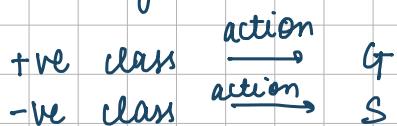
1	some	small	No	Affordable	one	No
2	many	big	No	expensive	many	Yes

$\langle ?, ?, \text{No}, ?, ? \rangle$

\hookrightarrow not consistent

S = most specific boundary / maximally specific hypothesis .

G = most general boundary / maximally general hypothesis



Candidate elimination

	sky	AirTemp	humidity	wind	water	Forecast	enjoy Sport
1	sunny	warm	Normal	strong	warm	same	yes
2	sunny	warm	high	strong	strong	same	yes
3	rainy	cool	high	strong	warm	change	no
4	sunny	warm	high	strong	cool	change	yes

$S_0 : \langle \phi, \phi, \phi, \phi, \phi, \phi, \phi \rangle$

$S_1 : \langle \text{Sunny}, \text{warm}, \text{Normal}, \text{Strong}, \text{weak}, \text{same} \rangle$

$S_2 = S_3 : \langle " ", " ", ?, " ", " ", " " \rangle$

$S_q : \langle \text{sunny}, \text{warm}, ?, \text{strong} \rightarrow ?, ?, ? \rangle$



$\langle \text{sunny}, ?, ?, \text{strong}, ?, ? \rangle$

$\langle \text{sunny}, \text{warm}, ?, ?, ?, ? \rangle$

here we pick

all of the possibilities

for five classes

we were checking here first

$G_q : \langle \text{sunny}, ?, ?, ?, ?, ? \rangle, \langle ?, \text{warm}, ?, ?, ?, ? \rangle, \langle ?, ?, ?, ?, ?, ? \rangle$

$G_3 : \langle \text{sunny}, ?, ?, ?, ?, ? \rangle, \langle ?, \text{warm}, ?, ?, ?, ? \rangle$ taking attributes which will result in a negative class.

$G_2 = G_1 : \langle ?, ?, ?, ?, ?, ? \rangle$

$G_0 : \langle ?, ?, ?, ?, ?, ? \rangle$

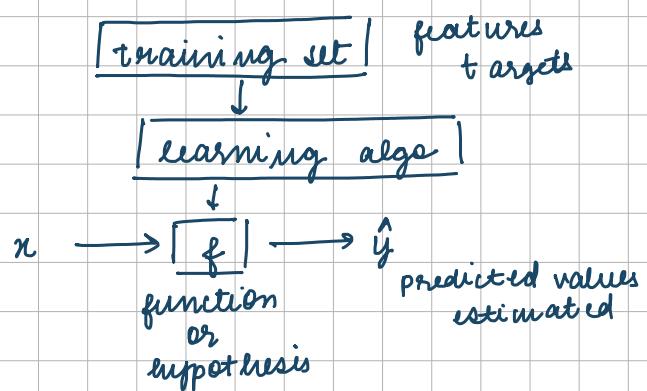
Ques

Example	Size	color	shape	Class
1	Big	Red	circle	No
2	small	Red	Triangle	No
3	small	Red	circle	Yes
4	Big	Blue	circle	No
5	Small	Blue	circle	Yes

06 Feb '24

Supervised machine learning - andrews ng

linear regression in one variable.



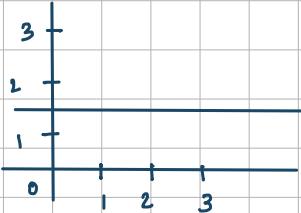
error / residual $\rightarrow \hat{y} - y \rightarrow |\hat{y} - y| \rightarrow (\hat{y} - y)^2$
can be differentiation
leading not defined.

How to represent 'f' in linear regression.

$$f_{w,b} = \underbrace{w_b + x}_{\text{or}} \rightarrow w_n + b \quad (w, b = \text{parameters})$$

$$f = \underbrace{w_b + x}_{\text{or}} \rightarrow w_n + b$$

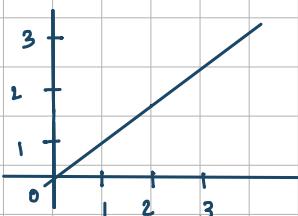
Reason to use linear regression: we want a simple function, complex functions are difficult to analyse.



$$f(x) = 0x + 1.5$$

Case I:

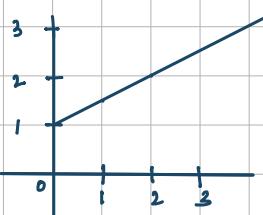
$$w = 0 \\ b = 1.5$$



$$f(x) = \frac{1}{2}x + 0$$

Case II:

$$w = 0.5 \\ b = 0$$



Case III

$$w = 0.5 \\ b = 1$$

$$f(x) = \frac{1}{2}x + 1$$

we want to find w, b so that $\hat{y}^{(i)}$ is close to $y^{(i)}$ for all $(x^{(i)}, y^{(i)})$

$$\hat{y}^{(i)} = \underbrace{f_{w,b}(x^{(i)})}_{\text{model}} = w x^{(i)} + b$$

Cost Function

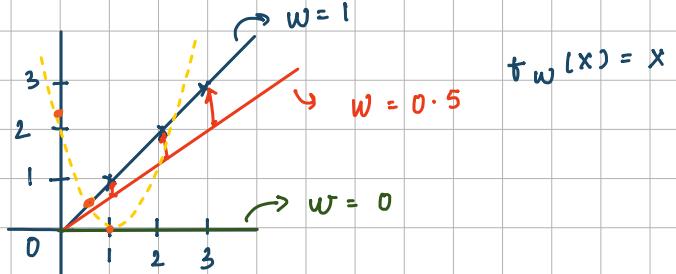
it tells us how well the model is doing so that we can adjust the parameters to improve our model (if required).

for m examples:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

our goal is to minimize this cost[↑]. w.r.t w, b .



$$\begin{aligned} w &= 1 \\ b &= 0 \end{aligned}$$

interactive plot.

$$f(x) = y \quad f(1) = 1, \quad f(2) = 2, \quad f(3) = 3$$

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

$$J(1) = \frac{1}{2 \times 3} [(-1)^2 + (2-2)^2 + (3-3)^2] = 0$$

$$J(0.5) = \frac{1}{2 \times 3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \approx 0.58$$

$$J(0) = \frac{1}{2 \times 3} [(0-1)^2 + (0-2)^2 + (0-3)^2] \Rightarrow \frac{1}{6} \times 14 \approx 2.3.$$

when slope = 0, cost function minimum.

example	size	color	shape	mass
1		Red	circle	No
2	small	Red	triangle	No
3	Small	Red	circle	Yes
4	Big	Blue	circle	No
5	small	Blue	circle	Yes

$$\langle \phi, \phi, \phi \rangle$$

$$S_0 : \langle \phi, \phi, \phi \rangle$$

$$S_1 : \langle \phi, \phi, \phi \rangle$$

$$S_2 : \langle \phi, \phi, \phi \rangle$$

$$S_3 : \langle \text{small, red, circle} \rangle$$

$s_4 : \langle \text{small, Red, Circle} \rangle$

$s_5 : \langle \text{small, ?, Circle} \rangle$

$t_3 : \langle \text{Big, ?, Triangle} \rangle, \langle ?, \text{Blue, Triangle} \rangle, \langle \text{small, ?, Circle} \rangle$
 $\langle \text{Big, ?, Triangle} \rangle, \langle ?, \text{Blue, Triangle} \rangle \times$

$t_2 : \langle \text{small, Blue, ?} \rangle \times, \langle \text{small, ?, Circle} \rangle \langle ?, \text{Blue, ?} \rangle \times$

$t_1 : \langle \text{small, ?, ?} \rangle, \langle ?, \text{Blue, ?} \rangle, \langle ?, ?, \text{Triangle} \rangle$

$t_0 : \langle ?, ?, ? \rangle$

07 Feb '24

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial J(w, b)}{\partial w} = 0 \quad \frac{\partial J(w, b)}{\partial b} = 0$$

$$\frac{\partial J}{\partial w} = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)}) = 0$$

$$= \frac{1}{m} (w x^{(i)} + b - y^{(i)}) (x^{(i)}) = 0 \quad \text{--- (1)}$$

$$\frac{\partial J}{\partial b} \frac{1}{2m} \sum_{i=1}^m (w x^{(i)} + b - y^{(i)})^2 = 0$$

$$\frac{1}{m} (w x^{(i)} + b - y^{(i)}) = 0 \quad \text{--- (2)}$$

By (2)

$$w(x^{(1)} + x^{(2)} + \dots + x^{(m)}) + b - (y^{(1)} + y^{(2)} + \dots + y^{(m)}) = 0$$

$$m w \bar{x} + b \sum_{i=1}^m \frac{m \bar{y}}{mb} = 0$$

$$b = \bar{y} - w \bar{x} \quad \text{--- (3)}$$

By (1)

$$\sum_{i=1}^m w(x^{(i)})^2 + \sum_{i=1}^m b x^{(i)} - \sum_{i=1}^m y^{(i)} x^{(i)} = 0$$

$$\sum_{i=1}^m w(x^{(i)})^2 + b m \bar{x} - \sum_{i=1}^m y^{(i)} x^{(i)} = 0$$

Substitute (3).

$$\sum_{i=1}^m w(x^{(i)})^2 + (\bar{y} - w \bar{x}) m \bar{x} - \sum_{i=1}^m y^{(i)} x^{(i)} = 0$$

$$w \left(\sum_{i=1}^m (x^{(i)})^2 - m \bar{x}^2 \right) - \sum_{i=1}^m y^{(i)} x^{(i)} + m \bar{y} \bar{x} = 0.$$

$$w = \frac{\sum_{i=1}^m y^{(i)} x^{(i)} - m\bar{y}\bar{x}}{\sum_{i=1}^m (x^{(i)})^2 - m\bar{x}^2}$$

$$= \frac{m\bar{y} - m\bar{x}\bar{y}}{m\bar{x}^2 - m(\bar{x})^2}$$

$$= \frac{\bar{y}\bar{x} - \bar{x}\bar{y}}{\bar{x}^2 - (\bar{x})^2} \quad \text{--- (5)}$$

$$= \frac{\text{cov}(x, y)}{\text{var}(x)}$$

Gradient Descent

12 Feb '24

Multi-variable Regression

$$y = f(x_1, x_2, x_3, \dots, x_n)$$

x_1 = size in feet²

x_2 = no. of bedrooms

x_3 = age of the house

x_4 = no. of floors

x_5 = price (\$)

x_1	x_2	x_3	x_4	x_5
2104	5	45	1	460
1916	3	40	2	232

$$\begin{aligned} j &= 1, 2, 3, 4 \\ &= 1 \dots 4 \end{aligned}$$

m = no. of training samples.

n = # of features

$$x_2^{(2)} = 3$$

initially: $f_{w,b}(x^{(i)}) = w x^{(i)} + b$

for multiple variable

$$\vec{y} = f_{\vec{w}, b}(\vec{x}^{(i)}) = \vec{w} \cdot \vec{x}^{(i)} + b$$

$$= w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + b$$

$$\vec{x}^{(1)} = [2104 \ 5 \ 45 \ 1].$$

$$f_{(\vec{w}, b)}(\vec{x}) = 0.1 x_1 + 4 x_2 + (-2) x_3 + 10 x_4 + \boxed{80}.$$

for any house which doesn't even consist of any feature this value represents base value of the house (in this case land).

$0.1 x_1$ = for every additional feet² the price \uparrow by 0.1×1000 . i.e. £100.

Cost function: $J(\vec{w}, b)$ or $J(w_1, w_2, \dots, w_n, b)$.

$$= \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$$

for features:

$$n \geq 2$$

repeat until convergence {

$$j=1, w_1 = w_1 - \alpha \frac{\partial}{\partial w_1} \left(\frac{1}{2m} \sum_{i=1}^m (\vec{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 \right)$$

$$= w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\vec{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$j=2, w_2 = w_2 - \alpha \frac{1}{m} \sum_{i=1}^m (\vec{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_2^{(i)}$$

⋮

$$j=n, w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (\vec{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)}$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (\vec{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

why feature scaling in multiple-variable regression is important?

size in feet²

$$x_1 [200-2000]$$

no. of bedrooms

$$x_2 [0-5]$$

Case I.

$$w_1 = 50$$

$$w_2 = 0.1$$

$$b = 50$$

$$x_1 = 2000$$

$$x_2 = 5$$

$$w_1 x_1 + w_2 x_2 + b = \hat{y}$$

$$\hat{y} = \$500,000$$

↖ (given).

$$\Rightarrow 50 \times 2000 + (0.1) \times 5 + 50$$

$$\Rightarrow 100,000 + 0.5 + 50$$

$$\Rightarrow 100,050.5 \times 1000$$



Case II

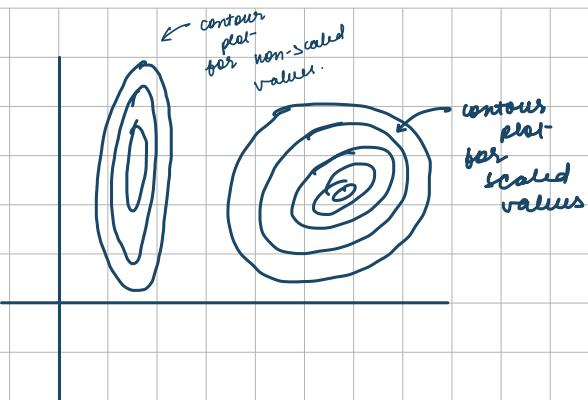
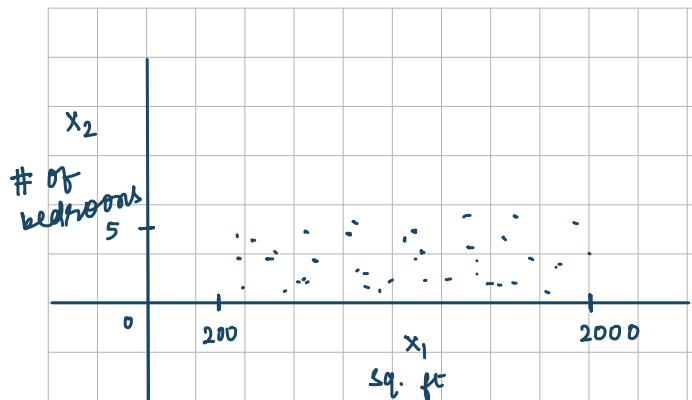
$$w_1 = 0.1 \quad w_2 = 50 \quad b = 50$$

$$\Rightarrow 0.1 \times 2000 + 50 \times 5 + 50$$

$$\Rightarrow 500 \times 1000$$

more precise.

w_1 was having \uparrow values and assigning $w_1 = 50$ for it resulted in even higher values and the contribution of x_2 was negligible in 1st case.



13 Feb '24

A very small change in w_i can have a very large impact on the estimated which creates a large impact on the cost.

Because the contours are too tall and somewhat thin, the gradient descent may end up bouncing back and forth for a long time before it reaches the global minimum.

min - max :

normalization

$$\frac{x_i - \text{min}}{\text{max} - \text{min}}$$

$$300 \leq x_1 \leq 2000$$

$$0.15 \leq x_1 \leq 1$$

scaled

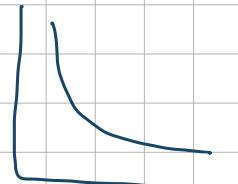
mean normalization :

$$\frac{x_i - \mu}{\text{max} - \text{min}}$$

z-score :

$$\frac{x_i - \mu}{\sigma}$$

if the learning curve shows any ↑ in if the learning rate is chosen or high



for early stopping use epsilon $\epsilon = 10^{-3}$

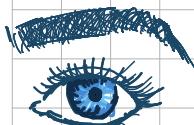
$J(w, b) \downarrow \leq \epsilon$ in one iteration ≤ 0.82

* Feature Engineering

height and width \Rightarrow area
 2 attributes 1 attribute.

created a new feature using existing features.

R² metric: (coefficient of determination).



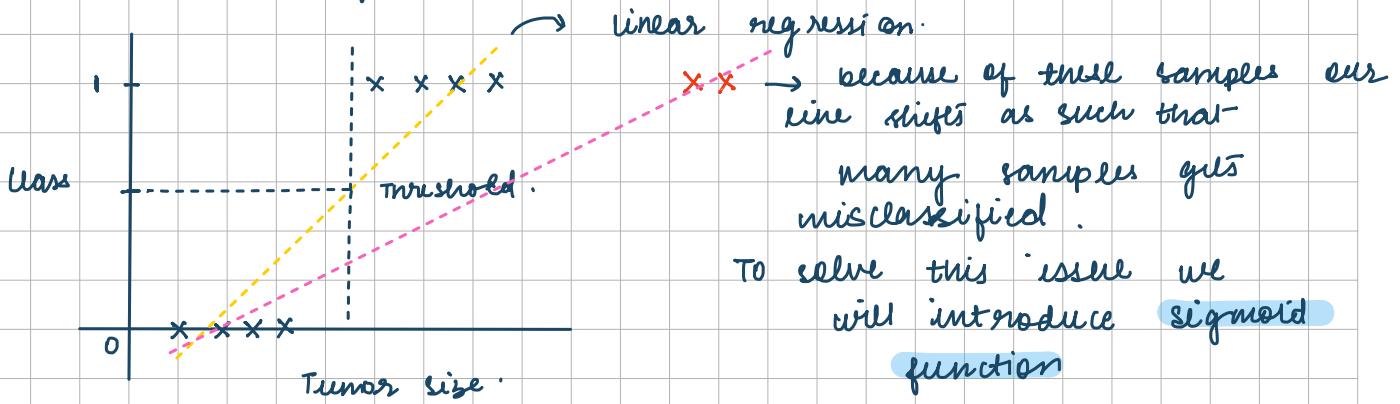
$$R^2 = 1 - \frac{\text{sum squared regression (SSR)}}{\text{total sum of squares (SST)}}$$

$$= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y}_i)^2}$$

14 Feb '24

* logistic regression (classification)

- Binary classification

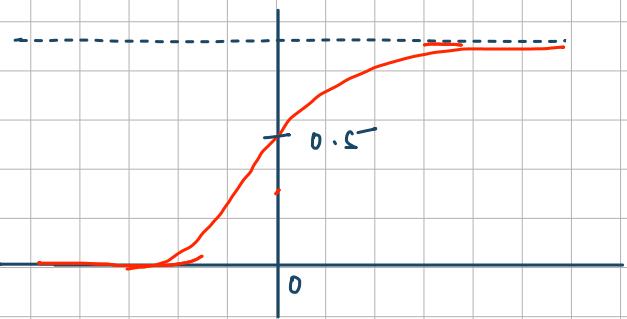


$$\vec{x} \rightarrow \boxed{\vec{w} \cdot \vec{x} + b = z} \rightarrow \boxed{g(z)} \rightarrow \hat{y}$$

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

$$\hat{y} = t_{\vec{w}, \vec{b}}$$

logistic/sigmoid function : $g(z) = \frac{1}{1+e^{-z}}$



$$\hat{y} \rightarrow p(y=1 | X=x)$$

$$p(Y=1 | X=x) + p(Y=0 | X=x) = 1$$

Case I : $z = 10000$ (high value).

$$g(z) = \frac{1}{1 + e^{-10000}} \rightarrow \text{this will be a } \downarrow \text{ value } \approx 0.$$

for higher values of z , $g(z)$ will approach 1.

Case II : $z = 0$

$$g(z) = \frac{1}{1 + e^0} = \frac{1}{1+1} = \frac{1}{2} = 0.5$$

Case III : $z = -10000$ (large -ve no.)

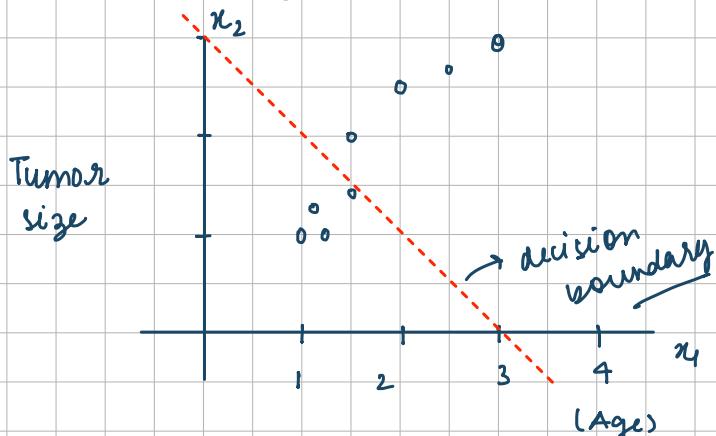
$$g(z) = \frac{1}{1 + e^{(-10000)}} \approx \frac{1}{\infty} = 0$$

$$\begin{aligned} f_{w,b} &\geq 0.5 & \hat{y} &\rightarrow 1 \\ &< 0.5 & \hat{y} &\rightarrow 0 \end{aligned}$$

problem : multiple local minima in SSE.

solution :

Decision Boundary



$$\begin{aligned} w_1x_1 + w_2x_2 + b \\ w_1=1, w_2=1, b=-3 \end{aligned}$$

when $w_1x_1 + w_2x_2 + b = 0$
↓
decision boundary.

When $z = -ve$
class $\Rightarrow -ve$
When $z = +ve$
class $\Rightarrow +ve$.

patient age (x_1)	tumor size (x_2) (in cm)	y
1 0	1	0
1	1.2	0
1.2	1.1	0
1.3	1.5	0
2	1.5	1
2.5	2	1
2.6	2.5	1
3.0	3	1

19 Feb '24

$$y^{(i)} = 1$$

for one example

$$l(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

$$= -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

cost function -

lost

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

convex cost f(n).

gradient descent

repeat until convergence.

{
j = 1, 2, ... n → no. of features.

$$w_j = w_j - \alpha \left(\frac{\partial J(\vec{w}, b)}{\partial w_j} \right) \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$b = b - \alpha \left(\frac{\partial J(\vec{w}, b)}{\partial b} \right) \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

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

$$J(\vec{w}, b) = \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log \left(\frac{1}{1 + e^{-z}} \right) + (1 - y^{(i)}) \log \left(\frac{e^{-z}}{1 + e^{-z}} \right) \right]$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m \log(1 + e^{-z}) + (1 - y^{(i)})[-z \log e - \log(1 + e^{-z})] \right]$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m \log(1 + e^{-z}) - z \log e - \log(1 + e^{-z}) + y^{(i)} z \log e + y^{(i)} \log(1 + e^{-z}) \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[-\log(1 + e^{-z}) - z + y^{(i)} z \right]$$

$$\begin{aligned}
 &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} z - [\log(1 + \frac{1}{e^z}) + \log(e^z)]] \\
 &= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} z - \log(e^z (1 + \frac{1}{e^z}))) \\
 &= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} z - \log(1 + e^z))
 \end{aligned}$$

$$\frac{\partial J}{\partial w_j} = \frac{\partial}{\partial w} \left(-\frac{1}{m} \sum_{i=1}^m (y^{(i)} z - \log(e^z (1 + \frac{1}{e^z}))) \right)$$

=

Muditā

21 Feb '24

- * ↑ data (not always possible).
- * hold-out method (not effective with less data).
- * K-fold cross validation
- * leave one out (less data).

} validation method
mean accuracy $\pm \frac{1}{\sqrt{n}} \cdot \text{variance}$,
confidence interval

SMOTE (over-sampling) (used for imbalanced data).

↳ **Regularization** (way to overcome overfitting)

$m = \text{no. of attributes}$

Regularization \rightarrow hyper parameter

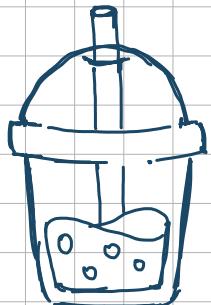
- Lasso regression (L1)
- Ridge " (L2)

$$\lambda \sum_{i=1}^n w_i^2$$

$$\lambda \sum_{i=1}^n |w_i|$$

(penalty term)

\downarrow
→ penalty to w corresponding to higher values



we will do L2 regularization

$$J(\vec{w}, b) = \underbrace{\frac{1}{2m} \sum_{i=1}^m f_{\vec{w}, b}(\vec{x}^{(i)}) + y^{(i)}_i}_\text{↑ degree of fit / error of data.}^2 + \underbrace{\frac{1}{2m} \sum_{i=1}^m w_i^2}_\text{penalty term. tries to keep wj small.}$$

how to choose λ ?

$$-\lambda > 0$$

$$\lambda = 10^{-8}$$
$$w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$
$$\approx 0 \quad \approx 0 \quad \approx 0 \quad \approx 0 \quad \approx b$$

(simpler model \Rightarrow underfitting)

$$\lambda = 0.0001$$

gradient descent

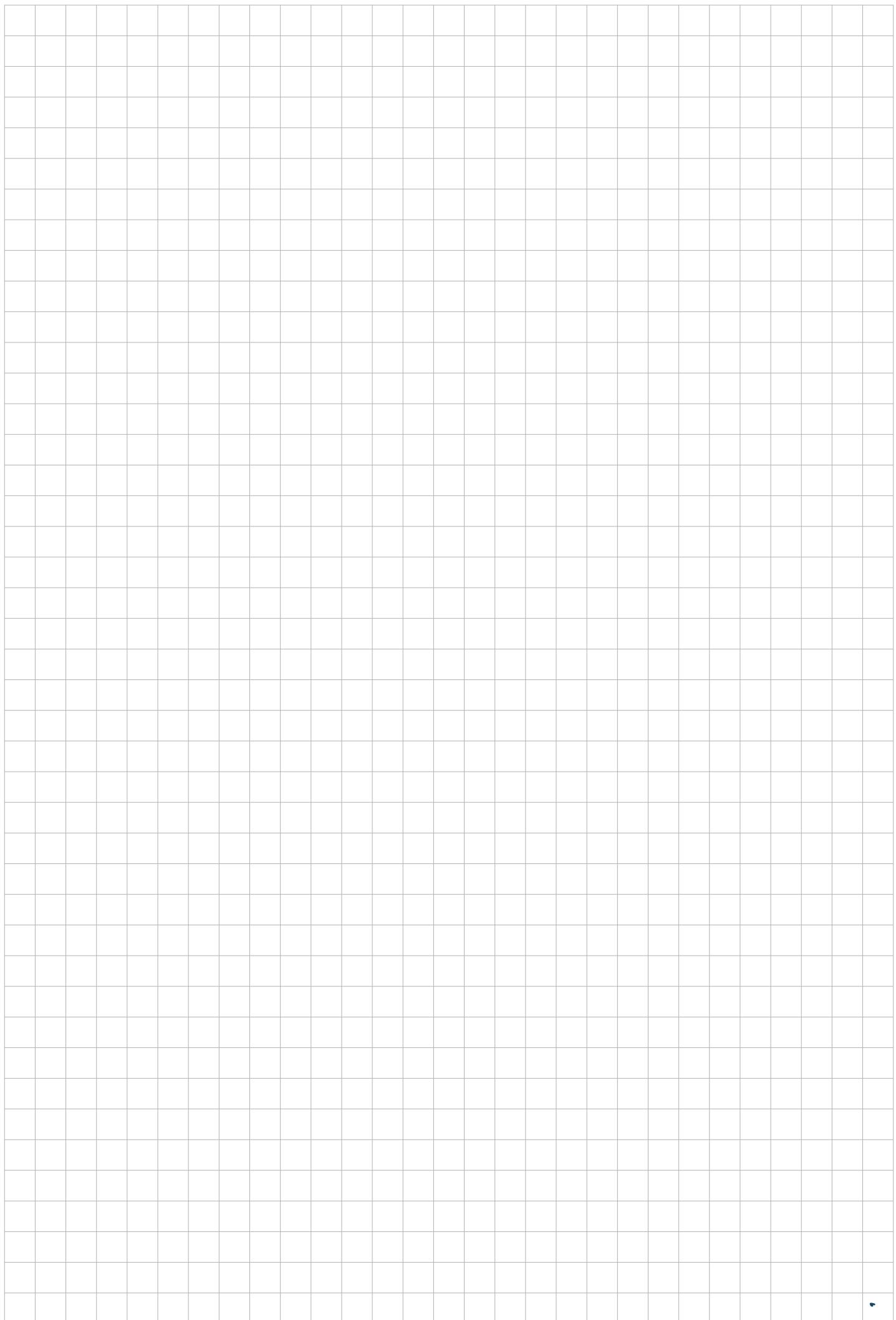
$$\min_{\vec{w}, b} J(\vec{w}, b) + \text{regularization term}$$

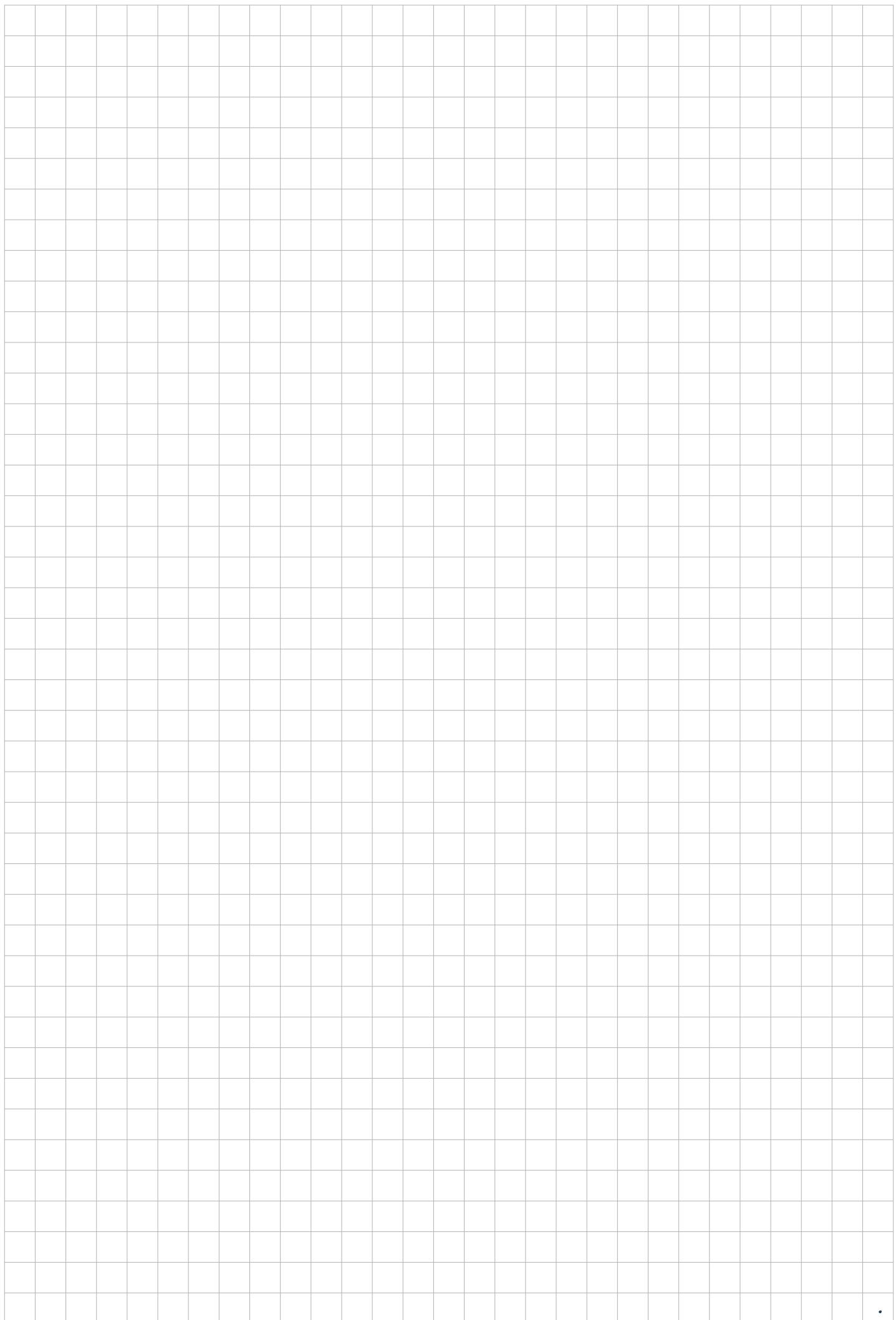
shrink / shrink \Rightarrow higher value rapidly decreases.

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

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

note: no summation





26 Feb '24

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

activation value of layer l , unit neuron j

parameters of layer l , neuron j

activation function (here sigmoid)

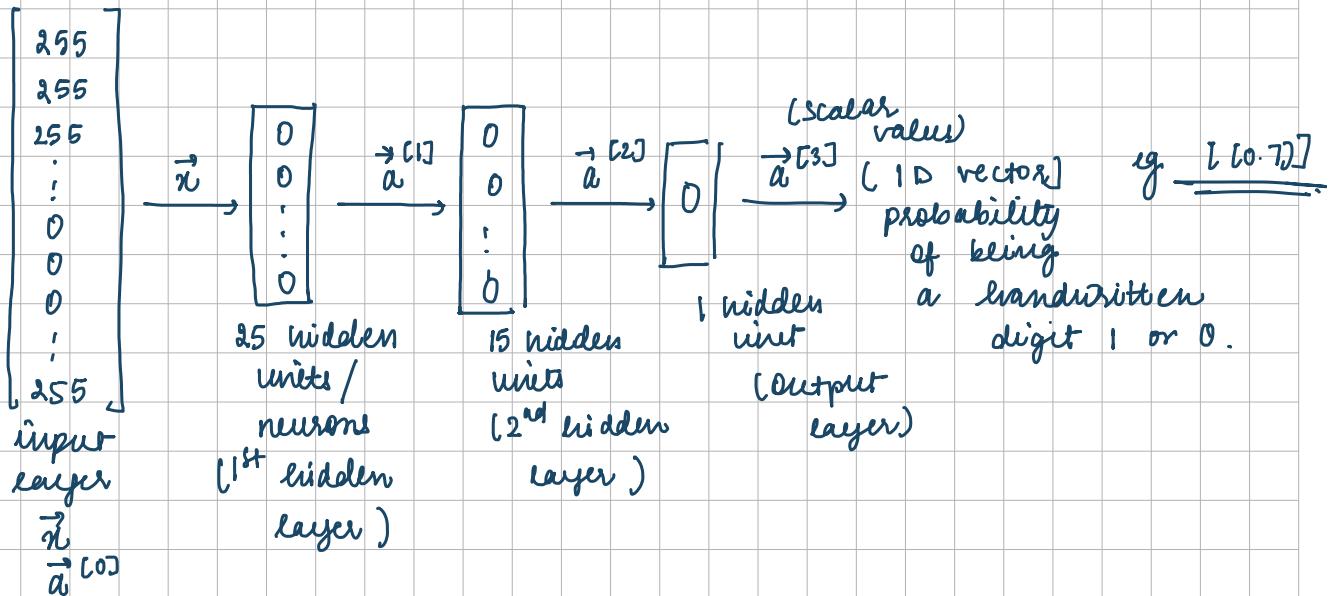
output of layer $l-1$ (i.e. previous layer)

FORWARD PROPAGATION

propagate the information to the next layer.

MNIST

Handwritten digit recognition \rightarrow 0 or 1
(8x8 binary images)



In order to get probability, we use sigmoid function.
(ranges from 0 to 1)

(ranges from 0 to 1)

$$\vec{a}^{[l]} = \left[\begin{array}{c} g(w_1^{[l]} \cdot \vec{a}^{[0]} + b_1^{[l]}) \\ \vdots \\ g(w_{25}^{[l]} \cdot \vec{a}^{[0]} + b_{25}^{[l]}) \end{array} \right] \rightarrow \vec{a}_1^{[l]}$$

$$\vec{a}^{[2]} = \begin{bmatrix} g(w_1 \cdot \vec{a}^{[1]} + b_1) \\ \vdots \\ g(w_{15} \cdot \vec{a}^{[1]} + b_{15}) \end{bmatrix}$$

$$\vec{a}^{[3]} = [g(w_1 \cdot \vec{a}^{[2]} + b_1), g(w_2 \cdot \vec{a}^{[2]} + b_2)]$$

$\vec{a}^{[3]} \geq 0.5$
yes / \ no
 $\hat{y} = 1$ $\hat{y} = 0$

Model Training Steps are:

- Create a model.

$$z = \text{np.dot}(w, x) + b$$

$$f(x) = \frac{1}{1 + \text{np.exp}(-z)}$$

```
import tensorflow as tf
from tensorflow.keras import
sequential
from tensorflow.keras.layers import
dense.
```

$\vec{a}^{[1]} \text{ (one after one)}$

 $\text{model} = \text{sequential}([\text{dense}(\text{units}=25, \text{activation}=\text{'sigmoid'})]$
 $\quad \quad \quad \text{dense}(\text{units}=15, \text{activation}=\text{'sigmoid'})]$
 $\quad \quad \quad \text{dense}(\text{units}=1, \text{activation}=\text{"sigmoid"})])$

↑ training example ↓ avg. loss for all examples.

- Specify loss and cost function

$$\text{cost} = \frac{1}{m} \sum_{i=1}^m (-y_i \log(f(x_i)) - (1-y_i) \log(1-f(x_i)))$$

[cross entropy]

for binary classification: binary cross entropy.

Code:

```
from tensorflow.keras.losses import binary_crossentropy.
model.compile(loss=binary_crossentropy(),)
    ↳ loss function
    for all the layers.
```

- Gradient descent / Back propagation
(minimizing the lost)

repeat {

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

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

}

code:

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

→ 100 forward 100 backward.

linear activation function → also known as no activation function.

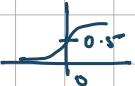
27 feb '24

Tuesday

choice of activation function

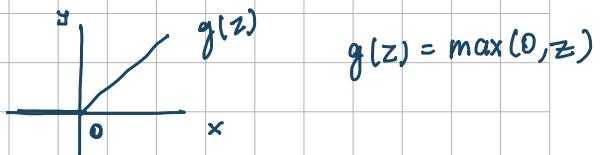
Case I: **Binary answer**

→ sigmoid



$$\frac{1}{1+e^{-z}}$$

Case II: **Non-negative number** indicating the degree of activation
→ Rectified Linear Unit (ReLU)

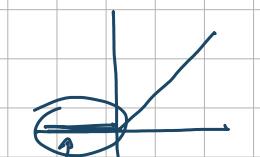
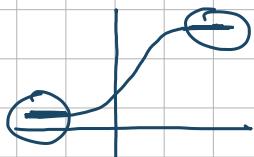


Regression.
Linear Activation
function)

$$y = +\text{ve} / -\text{ve} / ^0 .$$

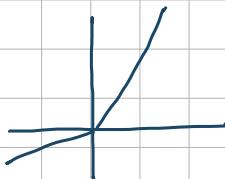
Case III: **Regression** → $y = 0 \text{ to } +\infty$.
(ReLU)

- most common choice of the activation function at hidden layer is ReLU.



(ReLU converges faster).

problematic
for gradient-
descent.



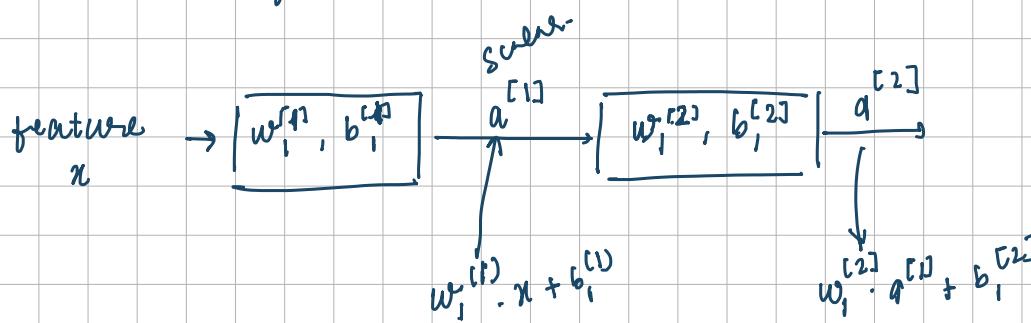
leaky
ReLU.

sequential([Dense(units=25, activation='relu')
dense(units=15, activation='relu')
dense(units=1, activation='sigmoid')])

* Why do we need activation functions?

- non-linearity.

-



$$a_1^{[1]} = w_i^{[1]} \cdot x + b_i^{[1]} \quad \text{--- } ①$$

$$a_1^{[2]} = w_i^{[2]} \cdot a_1^{[1]} + b_i^{[2]} \quad \text{--- } ②$$

substitute ① in ②

$$a_1^{[2]} = w_i^{[2]} (w_i^{[1]} \cdot x + b_i^{[1]}) + b_i^{[2]}$$

$$a_1^{[2]} = w_i^{[2]} \cdot w_i^{[1]} \cdot x + w_i^{[2]} \cdot b_i^{[1]} + b_i^{[2]}$$



28 Feb '24

Multiclass Classification

softmax:
Regression

$$P(y=1 | \vec{x}) \quad z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

The estimated chance of $y=1$ given x . softmax activation function transforms the raw output of a neural network into a vector of probabilities.

restricts the value of P from 0 to 1 .

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

↑ feature set.
Total probability / relative probability

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^n e^{z_k}} = P(y=j | \vec{x}) \quad [\text{Generalized form}]$$

Note: softmax regression with two classes is same as logistic regression.

softmax Regression is a generalization of logistic regression.

Cost Function:

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

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

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

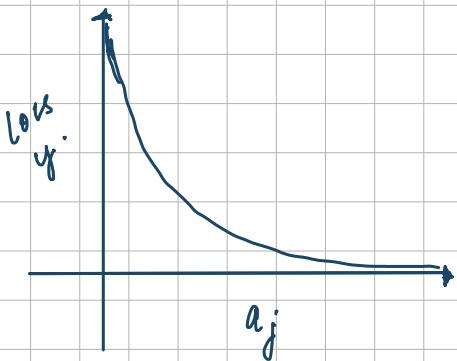
softmax:

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}}$$

:

$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}}$$

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



sparse categorical cross entropy.

29 Feb '24

multilabel classification

↳ medical data ?

Specify the model

```
import tensorflow as tf  
from tensorflow 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())
```

- Train the model to minimize $J(\vec{w}, b)$

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

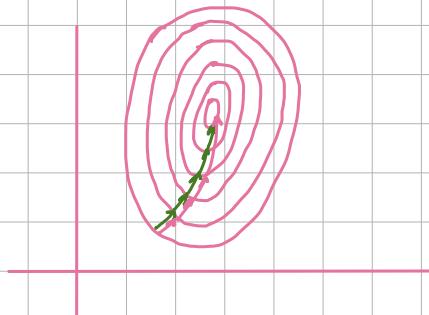
10

ADAM OPTIMIZER (Adaptive Movement Estimation)

some optimization to 'x'.



if the gradient descent is going in the right direction we ↑ alpha.



if the learning rate is relatively high than

* if the parameter value keeps on moving in roughly the same direction, let us ↑ the learning rate for that parameter.

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

$$w_2 = w_2 - \alpha_2 \frac{\partial J(\vec{w}, b)}{\partial w_2}$$

:

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

4%

conversely if the parameter value keeps oscillating back and forth, let us reduce α_j for that parameter j .

Recommender Systems



content based filtering
collaborative filtering
(mai or mere
like minded people)

↳ hybrid

2006
netflix challenge

4 80 000 users
18 000 movies
Rati ?
1984 Animal Farm

1 April '24

why vertical partitioning?
" horizontal " ?

train test? ↗

unbiased estimate of how
well the model is doing?

hyper-parameters.

development set / validation set →

Traditionally the dataset is divided into training and test-set and generally when the data is scarce, we often use 70/30 or 80/20 split.

In the era of big data we generally have 98/2 or 90/10 split. The work flow in this case would be train your algo on training set, → hyperparameter tuning 2nd use your dev set and end out cross validation set. to see which of many diff. models performs best on your dev set.



grid search cv

↳ cross validation.

pycaret

↑ project

3rd → take the best model you have found and evaluate it on your test set in order to get an unbiased estimate of how well your algo is doing.

JPEG
JPEG 2000

mismatched train - test distribution

Explainable AI
Post hoc

training set → high resolution, professional images.
test set → user uploaded blurred images.

HSDLSS

09 April '24

eigenvalues and eigenvectors of a 3×3 matrix

$$\det(A - \lambda I) = 0$$

$$Av_1 = \lambda_1 v_1$$

$$Av_2 = \lambda_2 v_2$$

$$A \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

basis \rightarrow

These vectors are linearly independent and other vectors depend on these.

feature extraction

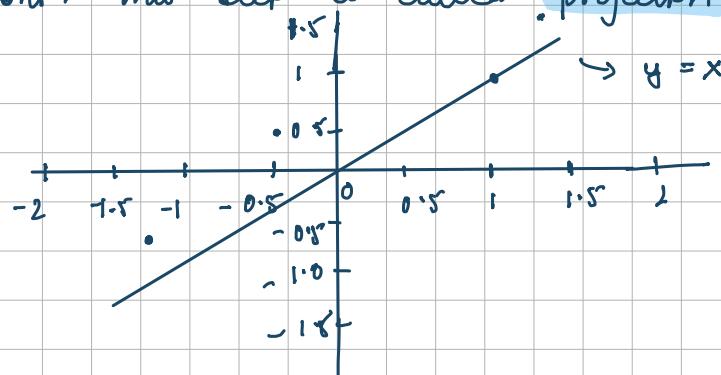
use of dimensionality Reduction!

- it leads to smaller datasets.
- it would be easier to visualize.

PCA allows us to reduce the dimension of our data. but it also preserves much of the info which we might lose if we simply started deleting columns.

Our aim is to move your data points into a vector space with fewer dimensions. This step is called projection.

	X	Y
A	1.0	1.0
B	1.2	1.6
C	-0.5	0.2
D	-1.3	-0.6



Projecting a data to a line.

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \frac{1}{\sqrt{[1][1]}} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \frac{1}{\sqrt{2}}$$

$$\sqrt{1^2 + 1^2} = \sqrt{2}$$

(norm).

$$\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \frac{1}{\sqrt{2}} = \frac{1+1}{\sqrt{2}} = \frac{2}{\sqrt{2}} = \sqrt{2} = 1.4142$$

$$\begin{bmatrix} 1.2 & 1.6 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \frac{1}{\sqrt{2}} = \frac{(1.2 + 1.6)}{\sqrt{2}} = \frac{2.8}{\sqrt{2}} = 1.9799$$

$$[-0.5 \ 0.2] \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \frac{1}{\sqrt{2}} = -\frac{0.3}{\sqrt{2}} = -0.2121$$

$$[-1.3 \ -0.6] \begin{bmatrix} 1 \\ 1 \end{bmatrix} \times \frac{1}{\sqrt{2}} = -\frac{1.9}{\sqrt{2}} = -1.344.$$

In general any matrix A can be projected onto the direction given by a vector V .

$$\downarrow A_{px_2} = A \cdot \frac{V}{\|V\|_2} \begin{bmatrix} \frac{V_1}{\|V\|_2} & \frac{V_2}{\|V\|_2} \end{bmatrix}_{c \times 2}$$

Projected space.

$$r =$$

c = no. of dimensions

Mean, variance and co-variance.

mean (x_i, y_i)

$$\mu_n = \bar{x} = \text{mean}(x) = \frac{1}{n} \sum_{i=1}^n x_i$$

$\forall i = 1 \dots n$
(the avg. of the data).

$$\mu_y = \bar{y} = \text{mean}(y) = \frac{1}{n} \sum_{i=1}^n y_i$$

(\bar{x}, \bar{y}) define the mean

Variance

how much spread your data has?

$$\text{var}(x) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_x)^2$$

$$\text{var}(y) = \frac{1}{n-1} \sum_{i=1}^n (y_i - \mu_y)^2$$

$x \uparrow y \downarrow \quad x \uparrow y \uparrow$

This relation is

Captured by

co-variance. $\text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$

$$A = \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_n & y_n \end{bmatrix}$$

$$C = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

$$\mu = \begin{bmatrix} \mu_x & \mu_y \\ \vdots & \vdots \\ \mu_x & \mu_y \end{bmatrix}$$

$$C = \frac{1}{n-1} (A - \mu)^T (A - \mu)$$

↓
matrix to
be projected.

$\rightarrow A$ ke dimension
jitna he μ ka
dimension bnaa
raha hai

A	$A - \mu$	$(A - \mu)^T$
x	x	
10	2	
12	4	
6	-2	
6	-2	
5	-3	
14	6	
8	0	
3	-5	
<hr/>	<hr/>	
$\mu_x = 8$	$\mu_y = 6$	

cov?

10 April '24

$$C = \begin{bmatrix} 9 & 4 \\ 4 & 3 \end{bmatrix}$$

$$\begin{aligned}\det(A - \lambda I) &= 0 \\ (9-\lambda)(3-\lambda) - 16 &= 0 \\ \lambda^2 - 12\lambda + 11 &= 0\end{aligned}$$

eigen value = 11, 1
eigen vector = $\begin{bmatrix} 2 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

$$\lambda = 11, 1$$

① create a matrix

$$X = \begin{bmatrix} x_{11} & \dots & x_{15} \\ x_{n1} & \dots & x_{n5} \end{bmatrix} n \times 5$$

② centre the data

"find: $L(X-\mu)$ "

$$(X - \mu) = \begin{bmatrix} x_{11} - \mu_1 & \dots & x_{15} - \mu_5 \\ \vdots & & \vdots \\ x_{n1} - \mu_1 & & x_{n5} - \mu_5 \end{bmatrix}$$

③ calculate the covariance matrix

$$C = \frac{1}{n-1} (X - \mu)^T (X - \mu)$$

$$= \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) & \text{cov}(x_1, x_3) & \text{cov}(x_1, x_4) & \text{cov}(x_1, x_5) \\ & \text{var}(x_2) & & & \\ & & \text{var}(x_3) & & \\ & & & \text{var}(x_4) & \\ & & & & \text{var}(x_5) \end{bmatrix}_{5 \times 5}$$

④ find the eigenvalues and eigen vectors
sort eigenvalues in decreasing order.

⑤ create a projection matrix.

$$V = \begin{bmatrix} \frac{v_1}{\|v_1\|_2} & \frac{v_2}{\|v_2\|_2} \end{bmatrix}$$

⑥ Project the centered data.

$$X_{PCA} = (X - \mu) \cdot V$$

(corresponding 2 principal components)

Ques

(2, 1) (3, 5) (4, 3) (5, 6) (6, 7) (7, 8)

10

$$x = \begin{bmatrix} x & y \\ 2 & 1 \\ 3 & 5 \\ 4 & 3 \\ 5 & 6 \\ 6 & 7 \\ 7 & 8 \end{bmatrix}$$

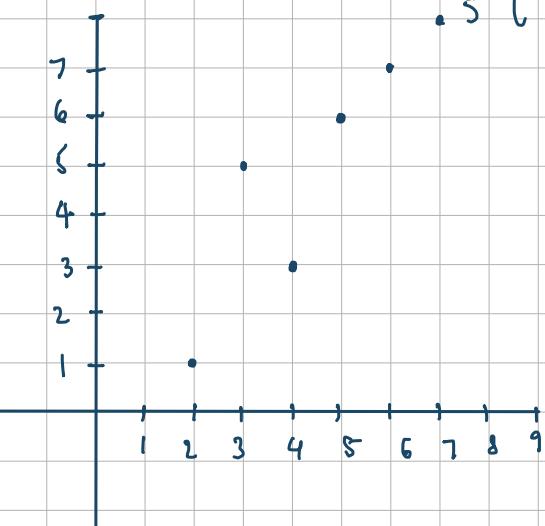
$$\bar{x} = \frac{27}{6} = 4.5$$

$$\bar{y} = \frac{30}{6} = 5$$

$$x - \mu = \begin{bmatrix} -2.5 & -4 \\ -1.5 & 0 \\ -0.5 & -2 \\ 0.5 & 1 \\ 1.5 & 2 \\ 2.5 & 3 \end{bmatrix}$$

$$C_2 = \frac{1}{5} \begin{bmatrix} -2.5 & -1.5 & -0.5 & 0.5 & 1.5 & 2.5 \\ -4 & 0 & -2 & 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} -2.5 & -4 \\ -1.5 & 0 \\ -0.5 & -2 \\ 0.5 & 1 \\ 1.5 & 2 \\ 2.5 & 3 \end{bmatrix}$$

$$\cdot \frac{1}{5} \begin{bmatrix} 17.5 & 22 \\ 22 & 34 \end{bmatrix} \Rightarrow \begin{bmatrix} 3.5 & 4.4 \\ 4.4 & 6.8 \end{bmatrix}$$



$$|C - \lambda I| = 0$$

$$\begin{vmatrix} 3.5 - \lambda & 4.4 \\ 4.4 & 6.8 - \lambda \end{vmatrix} = 0$$

$$(3.5 - \lambda)(6.8 - \lambda) - (4.4)^2 = 0$$

$$23.8 - \lambda^2 + 10.3\lambda - 19.36 = 0$$

$$\lambda_1 = 9.849 \Rightarrow \begin{bmatrix} 0.693 \\ 1 \end{bmatrix} \Rightarrow$$

$$\lambda_2 = 0.45 \Rightarrow$$

Variance captured by λ_1 ,

$$\frac{9.84}{9.84 + 0.45} \Rightarrow \frac{9.84}{10.29} = 0.9562$$

~~9.84
0.45
10.29~~

$$\begin{bmatrix} -2.5 & -4 \end{bmatrix} \begin{bmatrix} 0.57 \\ 0.82 \end{bmatrix} \Rightarrow -4.73$$

$$\begin{bmatrix} -1.5 & 0.7 \end{bmatrix} \begin{bmatrix} 0.57 \\ 0.82 \end{bmatrix}$$

15 April 2024

Bayesian classifiers

Class conditional
independence.

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)}$$

X be a data tuple.

$$\frac{|C_i, D|}{|D|}$$

Class conditional independence.

$$P(X|C_i) = P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_d|C_i).$$



— why naive bayes works even with such high assumption?
if there exist many attributes which are correlated
then we also have some attributes $\perp \!\!\! \perp$ net $\perp \!\!\! \perp$.



should follow the normal distribution



should class

a b c d e f g h i j k

l m n o p q r s t

u v w x y z

18 April 2024

End-to-End CNN Transfer Learning.

23 April 2024

Padding

intuition 1: if you don't wish to shrink the image
 intuition 2: if the corner pixel consists of more
 padded convolution information, we need padding.

Output size without padding

$$(n - f + 1) \times (n - f + 1)$$

$$(n + 2p - f + 1) \times (n + 2p - f + 1)$$

↑ ↑
image padding filter

what should be the filter size?



$$n + 2p - f + 1 = n$$

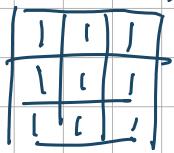
can find padding/filter size using this if we wish to have the same size of output and input.

same padding → output has size same as [padding = on]
 valid padding → no padding ($p = 0$),

Strided convolution

Stride: By default stride is 1 but we can have (step) 2, 3 or any number.

But why?



→ if we have the same information then by ↑ the size of stride we can capture the information and also compress.

$$\Rightarrow \left[\frac{n + 2p - f + 1}{s} \right] \times \left[\frac{n + 2p - t + 1}{s} \right]$$

$s = \text{stride} \cdot [\text{steps for columns as well as rows}]$

→ to use corner pixels we need this.

\times	1	2	3	4	5	6	7	8	9
=									
	2	3	7	4	6	2	9		
	6	6	9	8	7	4	3		
	3	4	8	3	8	9	7		
	7	8	3	6	6	3	4		
	4	2	1	8	3	4	6		
	3	2	4	1	9	8	3		
	0	1	3	9	2	1	4		

*	3	4	4
	1	0	2
	-1	0	3

$$g = 2$$

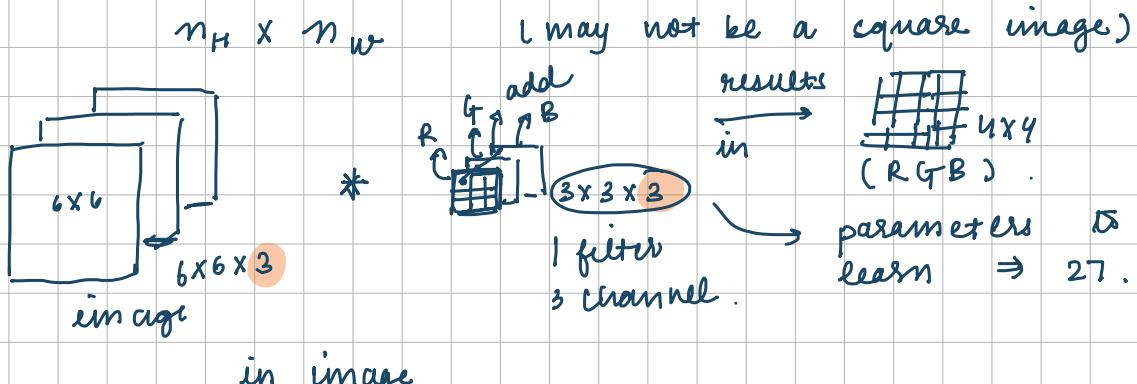
$$\begin{array}{r}
 9 + 16 + 32 + \quad \Rightarrow \quad 91 & 100 & 88 \\
 7 + 6 - 4 + \quad \quad \quad 69 & 91 & 127 \\
 3 \quad \quad \quad \quad \quad 44 & 72 & 74
 \end{array}$$

$$\begin{array}{r}
 & & 1 & 12 \\
 & & 46 & + 6 \\
 & 70 & 16 & \\
 \underline{-}3 & & 52 & \\
 & + 18 & 28 & \\
 & + 24 & 46 & \\
 \hline
 & 70 & 46 & \\
 \end{array}$$

$b + 12 + 28 +$
 $b + 18 + (-8)$
 $+ 8 \times 3'$
 \Rightarrow
 6×3
 $18 + 8 + 86 +$
 $7 + 6 - (8) +$
 21

24 April 2024

Convolution over Volumes



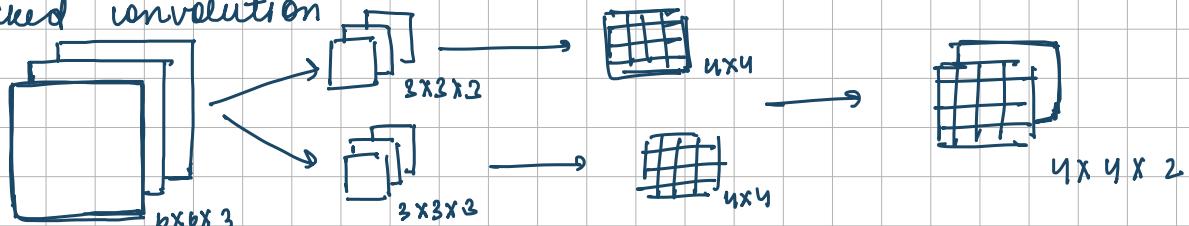
no. of channels should match the no. of channels in filter
with even filters \rightarrow aliasing. \leftarrow results in
images can't be padded in case of even filter
output image!

$$n_H \times n_w \times n_c \quad (n_H - f + 1) \times (n_w - f + 1).$$

\downarrow channels

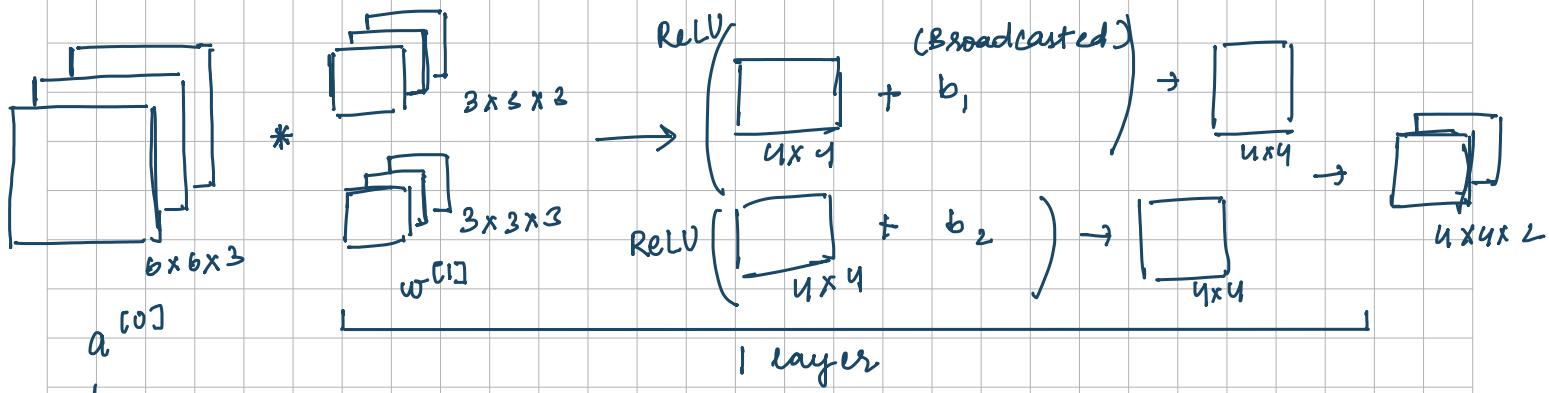
Take each of the 27 no. of the filter and multiply them with corresponding no. from Red, Green and Blue channel.

Stacked convolution



$$(n_H - f + 1) \times (n_W - f + 1) \times n'_C \hookrightarrow \# \text{ of filters.}$$

How to represent a layer of convolution neural network?



pooling layer (no parameter)

- if you have 10 filters that are $3 \times 3 \times 3$ in size of 1 layer of neural network?
 $\Rightarrow 3 \times 3 \times 3 \times 10 + 10 \text{ bias}$

If we are at layer l of CNN, then:

- $f^{[l]}$ = filter-size at layer l
- $p^{[l]}$ = padding at layer l
- $s^{[l]}$ = stride at layer l
- $n_C^{[l]}$ = no. of filters at layer l .

Input size: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$$n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]} + 1}{s^{[l]}} \right\rfloor$$

$$n_W^{[l]} = \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]} + 1}{s^{[l]}} \right\rfloor$$

25 April 2024

$$\text{input} : n_H^{[L-1]} \times n_w^{[L-1]} \times n_c^{[L-1]}$$

$$\text{output} : n_H^{[L]} \times n_w^{[L]} \times n_c^{[L]}$$

$$\text{filter size} : f^{[L]} \times f^{[L]} \times n_c^{[L]}$$

Activations :

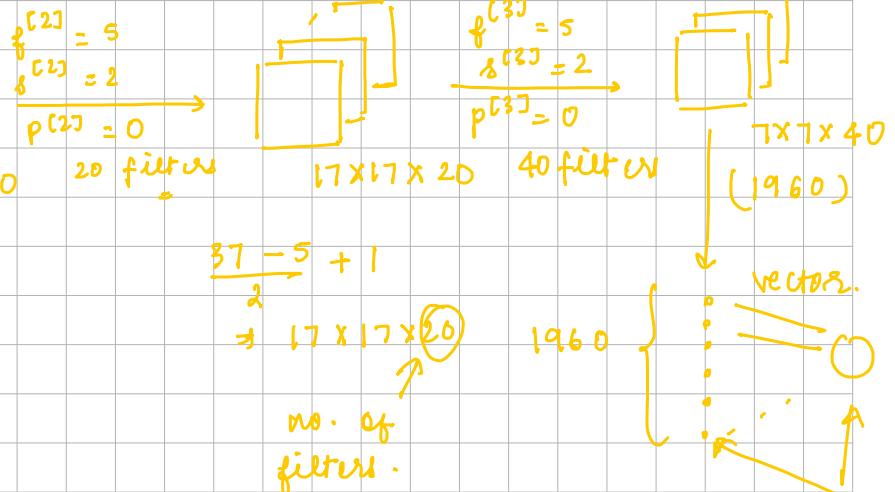
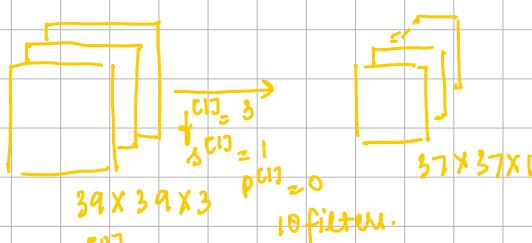
$$a^{[L]} : n_H^{[L]} \times n_w^{[L]} \times n_c^{[L]}$$

$$\text{weights} : f^{[L]} \times f^{[L]} \times n_c^{[L-1]} \times n_c^{[L]} \quad \hookrightarrow \# \text{ of filters}$$

→ detecting multiple edge detection. [multiple filters].

$$\text{bias} : n_b^{[L]} \quad (1, 1, 1, n_c^{[L]}) \quad \hookrightarrow \text{for broadcast.}$$

A simple convolution neural network.



$$39 + (-3) + 1$$

$$36 + 1 = 37$$

CONV

↳ hyperparameters

[



pooling layer
fully connected layer (FC)
max pooling
average pooling

$$\begin{matrix} 1 & 3 & 2 & 1 \\ 2 & 9 & 1 & 1 \\ 1 & 4 & 2 & 3 \\ 5 & 6 & 1 & 2 \end{matrix}$$

pooling layer

$$g = 2$$

$$f = 2$$



output of
avg. pooling
layer

output
of pooling
layer

$$\begin{matrix} 3.75 & 1.25 \\ 4 & 2 \end{matrix}$$

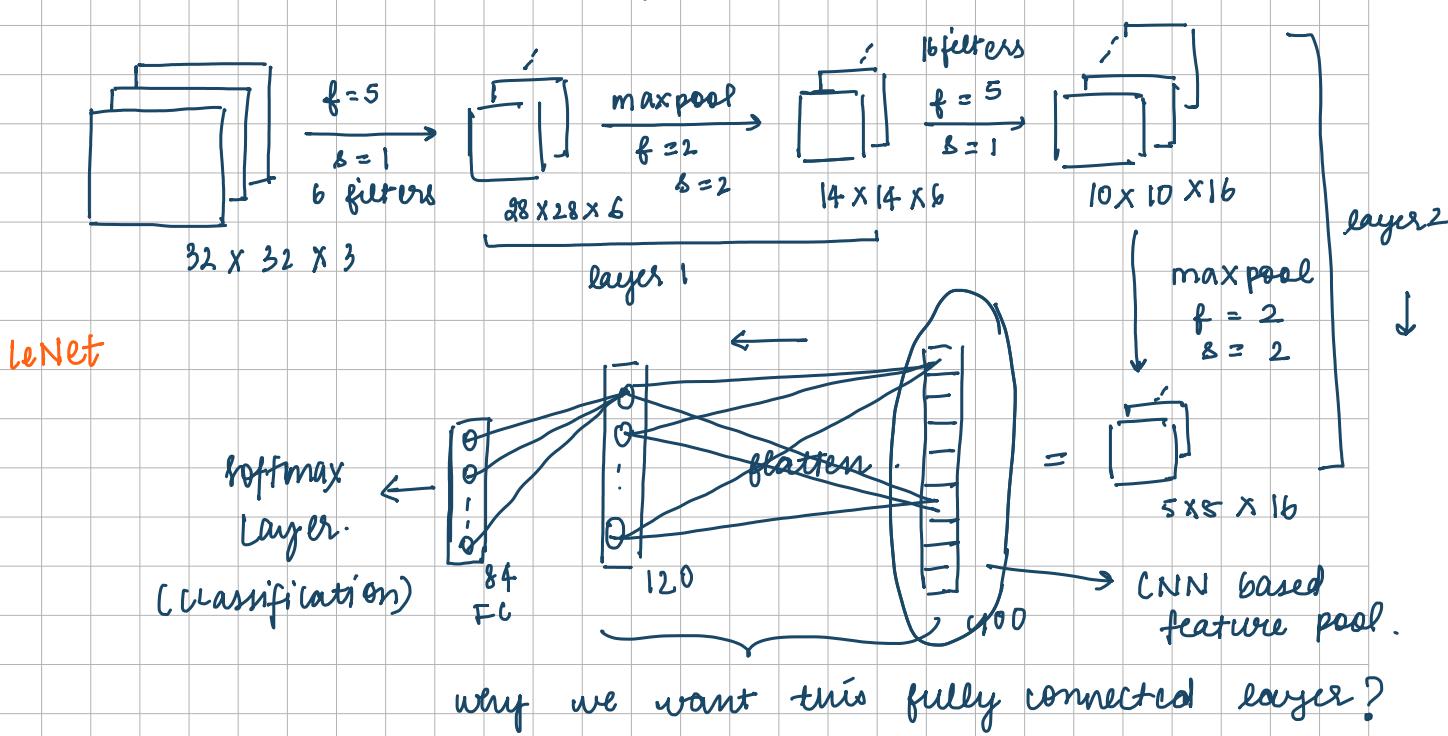
$$\begin{matrix} 9 & 2 \\ 6 & 3 \end{matrix}$$

disadvantage: loss of information

disadvantage:

→ in pooling layer there is no parameter to learn.

29 April '24



No. of trainable parameters?

Input	Activation shape -	Activation size	# of parameters
CONV1 [$f=5, s=1$]	$(32, 32, 3)$	3072	0
POOL1 [$f=2, s=2$]	$(28, 28, 8)$	6272	608
CONV2 [$f=5, s=1$]	$(14, 14, 8)$	1568	0
POOL2 [$f=2, s=2$]	$(10, 10, 16)$	1600	3216
FC1	$(120, 1)$	120	0
FC2	$(84, 1)$	84	$120 \times 84 + 84$
Softmax	$(10, 1)$	10	$84 \times 10 + 10$

no. of filters $\left(\frac{\text{size of filter}}{\text{size of filter}} \times \text{no. of channel} + 1 \right)$ (prev layer).

$$8 \times (5 \times 5 \times 3 + 1)$$

$$16 \times (5 \times 5 \times 8 + 1)$$

- why convolution?

- reduced no. of parameters (parameter sharing)
- sparsity of connection

$$1000 \times 1000 \times 3$$

feature detector (say a vertical edge detector) that is useful in 1 part of the image is probably useful for the other part of the image.

in each layer, each output value depends only on a smaller no. of inputs (prevents overfitting.)

* CNN are translation invariance.