26/05/23

## DAY 5:

I. Machine Learning

Applications

1 sound / (speech recognition)

$\rightarrow$ +ve voltage wave.
(emotions, accent, speed detected)

(FEATURE EXTRACTION)

of an AI system:
written problem v/s expected outcome

↓

this accuracy is very high.

naturally → no idea how its done.

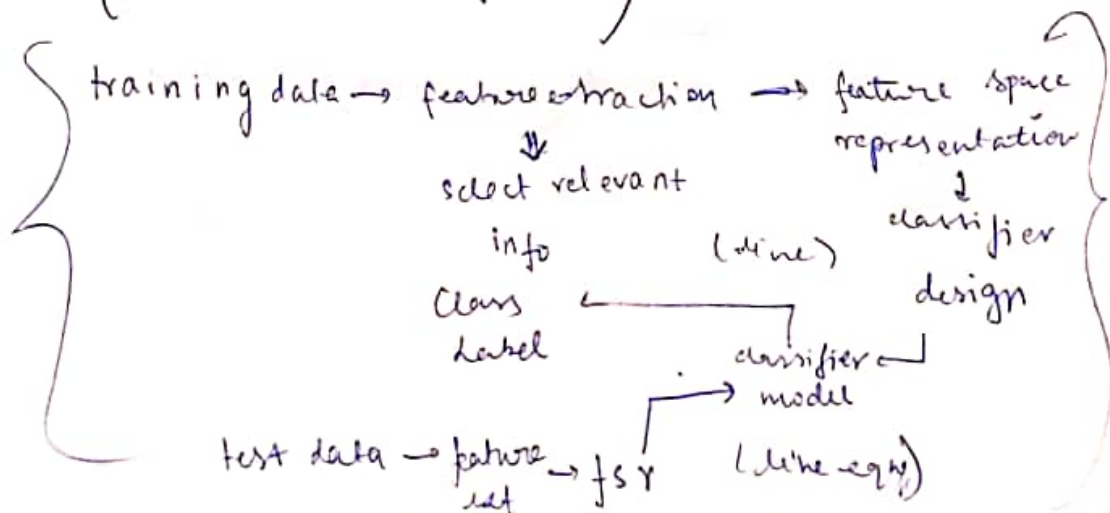$$y = f(x)$$

learn $f$
code $f$

Classification → int $y$
clustering → unsupervised.
Regression → real no. $y$

Time series prediction → e.g. predict future weather based on past 5 days.

* (formulate real world problem) *

training data → feature extraction → feature space representation

⇊

select relevant info        (line)        classifier design

Class label ← classifier

→ model

test data → feature set → $f's$ $y$    (line eqn)

not corrections, its _accuracy_.

<div align="center">IMP.</div>

either decide an eqn that divides | fit a model for different

generative | discriminative.                    classes.

① Classifier— k-NN ( find k nearest neighbours

image → feature vector

↳ most common class is
the answer)
(majority)
most frequent label.

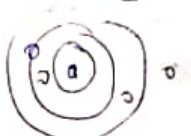for multiple classes → voronoi tessellation. ( between

piecewise linear boundary → final product.
(course ⊥ bisector)

nearest
neighbours)

sklearn neighbours.

Training is trivial . ( no learning effectively).
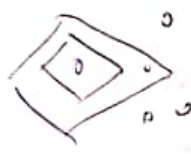
⚹ Bayesian classifier known to be most ideal

what is near→                    (theoritical proof)

L2 ¹) Euclidean Distance → $\sqrt{\sum(x_1-x_2)^2}$ $\sum(x_1-x_2)^2$

Iso surfaces → ○ ○ draw circles.

L1 ²) Manhattan distance (no sq.).

$\sum |x_1-x_2|$

iso surface

draw rhombus.

L0   distance?

3) Minkowski Dist

$\sum(x_1-x_2)^r$

r = 2  = ○     r = ∞ = ☐

r = 1 = ◇

r = 1.25 = ◇

If features have different variances
⤷ some feature will dominate
⤷ so req. to dominate.



normalization.

4) Mahalnobis Dist.

$$d(P, Q) = (P-Q)^T S^{-1} (P-Q)$$

weighted distance ⟶ for accounting for features that change less/more.
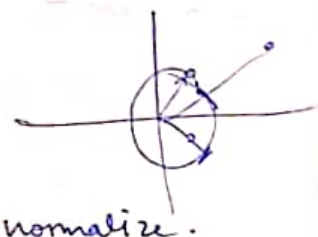
5) Hamming distance.

(binary vectors)

$$d(P, Q) = \sum I(p_i \neq q_i)$$

⤷ Indicator funct. (T=1, F=0)

6) cosine distance

$$d(P, Q) = \cos \theta = \frac{\vec{P} \cdot \vec{Q}}{|\vec{P}||\vec{Q}|} = \frac{P \cdot Q}{\|P\| \cdot \|Q\|}$$

(hypersphere - 4D+)



normalize.

feature vectors need not be of same length.

non metric, edit, jaccard dist.

② Linear Classifier

linear decision boundaries.

pt, line, plane, hyperplane.

$w_0 + w_1 x_1 + w_2 x_2 \quad \cdots \quad w_d x_d = 0$
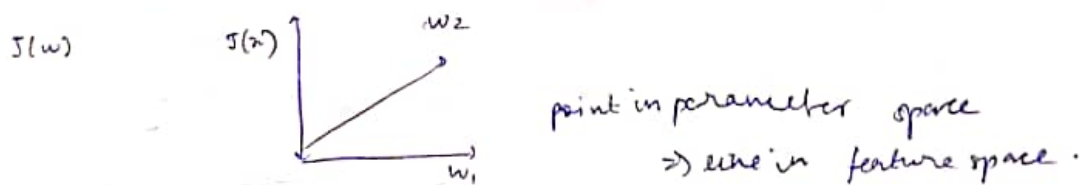
(d dimension space line).

$f(x) = 0$.

goal → learn parameters $w_i$ so training samples
are separated.

linear boundary → better to compute.
understanding is simplifying.

Minimizing loss func.

$J(w)$         $J(x)$



point in parameter space
⇒ line in feature space.

iteratively modify $w_i$ to reduce $J(w)$.
what modifications reduce $J(w)$?

$\dfrac{\partial J}{\partial w_i} \nearrow$    } gradient descent
       slope          slope

$\dfrac{\partial J}{\partial w}$   } gradient. (ascent)
       vector          (vector)

$\dfrac{-\partial J}{\partial w}$   } gradient descent.

$w = \dfrac{\partial J}{\partial w}$.

of algorithm.     $w^{t+1} = w^t - \eta \dfrac{\partial J}{\partial w}$

Loss function of GD)
↳ no. of misclassified samples.
    (not differentiable).

↳ $, y_i(w^T x_i) > 0$ for all samples)
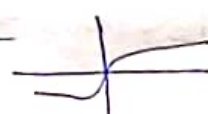
$J(w) = -y_i(w^T x_i)$ (for misclassified samples)

$\nabla J = -y_i x_i$

called perceptron update rule.

so

$w^{t+1} = w^t + \eta \sum_{i \in X} y_i x_i$ → already negative.

$\begin{array}{c|c} \cdot & 0 \\ \hline 0 & 1 \end{array}$ } linear classifier doesn't work, no way to separate.

$x_1$   $w_1$
$x_2$   $w_2$          , bias.
 :       :                              , output.
 :       :      $\sum_{i=1}^{n}(w_i x_i) + b$  if $> T$ } $\phi()$  $= y$
                                  ↓              ↳ activation
                              threshold.              function.

hyperbolic tangent
    tanh ⎯⎯⎯⎯ } most popular

ReLU ⎯⎯⎯

determine $\underline{w, b}$

i/o                    o/p
layer      (neuron) layer
           hidden
           layer.

neural networks help in creating cloud
boundaries (highly complex).

determining w/b is very very difficult.

Only 1 hidden layer suffices .

>2 hidden layers => vanishing gradient descent.
(back propagation).

hyperparameter = no. neurons
                    ↓
            user pics,

MLP

weights in     weights all



I/O        H          matrix    o/p   non
                      mult            linearity.

matrix                  apply non linearity.
multiplication.