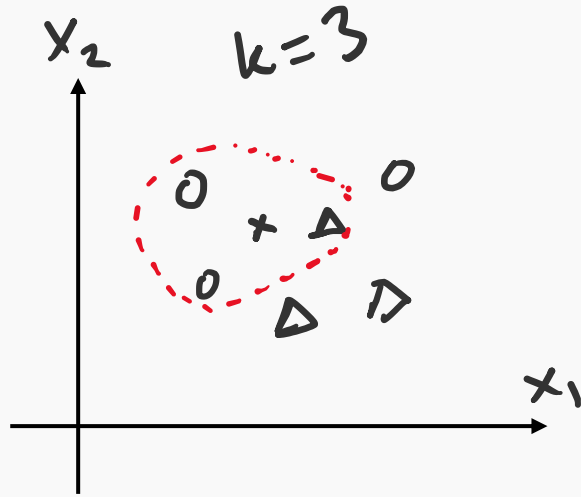


kokchun giang

leverage the distance to your
neighbours for classification or
regression using **KNN**



find closest distance to your **k neighbours**



○ label 0
Δ label 1
x test sample

⇒ x classified
as 0 here

KNN

k nearest neighbours

1. compute distances d_i between test sample & all training samples
2. sort d_i descending order
3. choose k nearest points

4. majority voting
in regression case
⇒ mean of
k closest
neighbours values

k hyperparameter
→ elbow plot
→ k-fold cross validation

there are several ways to measure **distances**

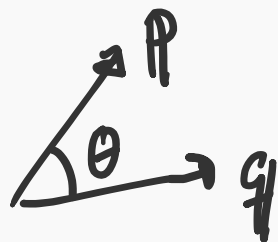
distance depends
on similarity
measure

→ Euclidean distance

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + \dots + (p_n - q_n)^2}$$

→ cosine similarity

$$S_c(p, q) = \cos \theta$$



measures angle θ
between 2 vectors

$$-1 \leq S_c \leq 1 \quad \leftarrow \text{closest}$$

↑
opposite

$$\theta = 0^\circ \Rightarrow \cos \theta = 1$$

$$\theta = 180^\circ \Rightarrow \cos \theta = -1$$

freq not important
here

cosine similarity
used extensively in
NLP as vectors
are very sparse
2) computationally
efficient

Ex document similarity

