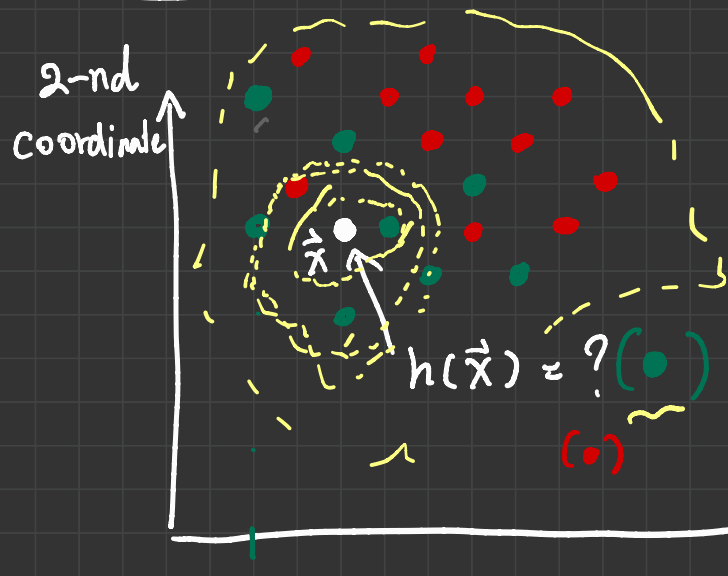


k-Nearest Neighbors (kNN)

เพื่อนบ้านที่ใกล้ที่สุด



$k = 2$

training: store all data points

testing: given test point
 \vec{x}

$$(\vec{x}_1 = \begin{bmatrix} 2 \\ 10 \end{bmatrix}, y_1 = \bullet) \in D$$

$$D \subseteq X \times Y$$

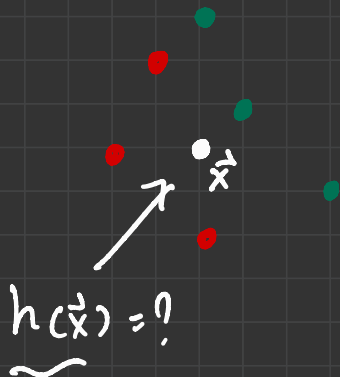
$$X = \mathbb{R}^2$$

$$Y = \{ \bullet, \bullet \}$$

Binary classification

k-NN's assumption: "Close (similar) points
should have similar labels."

Closeness \Rightarrow similarity

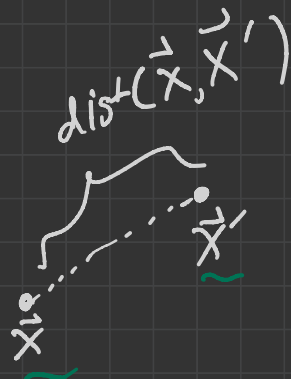


Classification Rule:

The label of the test point
is determined by the majority of k-NN

k-NN:

- Store all data points in the data set D .
- Given the test point \vec{x}
- Denote $S_{\vec{x}}$ the set of k -NN of \vec{x}



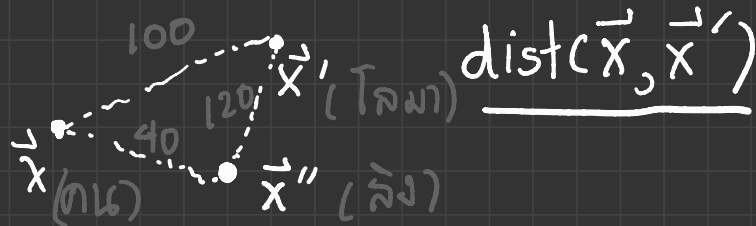
$$\underline{S_{\vec{x}}} \subseteq D \text{ s.t. } |S_{\vec{x}}| = k \text{ and}$$

$$\forall (\vec{x}', y') \in D - S_x \quad \boxed{\text{dist}(\vec{x}, \vec{x}')} \geq \max_{(\vec{x}'', y'') \in S_x} \text{dist}(\vec{x}, \vec{x}'')$$

$$\Rightarrow \underline{h(x) = \text{mode}(\{y'' : (\vec{x}'', y'') \in S_x\})}$$

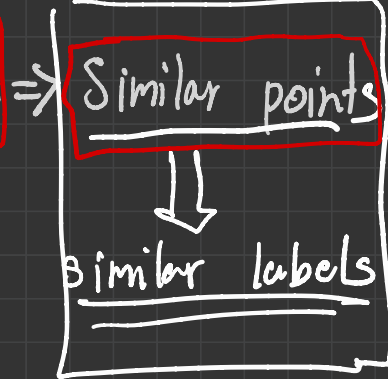
Distance function / metric

- Distance metrics measure levels of similarity among data points.



- The most common choice is "Minkowski distance"

Close points



Minkowski distance:

$$\text{dist}(\vec{x}, \vec{x}') = \left(\sum_{i=1}^d |x_i - x'_i|^p \right)^{\frac{1}{p}}$$

- $p=1$ (Manhattan)

- $p=2$ (Euclidean)

- $p=+\infty$ (Max)

- $p=-\infty$ (??)

$$\begin{aligned} \text{dist}(\vec{x}, \vec{x}') &= \left(\sum_{i=1}^d |x_i - x'_i|^p \right)^{\frac{1}{p}} \\ &= \max_{i=1}^d |x_i - x'_i| \end{aligned}$$

$$\begin{aligned} \text{dist}(\vec{x}, \vec{x}') &= \sqrt{|0-2|^2 + |0-1|^2} \\ &= \sqrt{2^2 + 1^2} \approx \sqrt{2^2} \\ &= 2 \end{aligned}$$