# Context

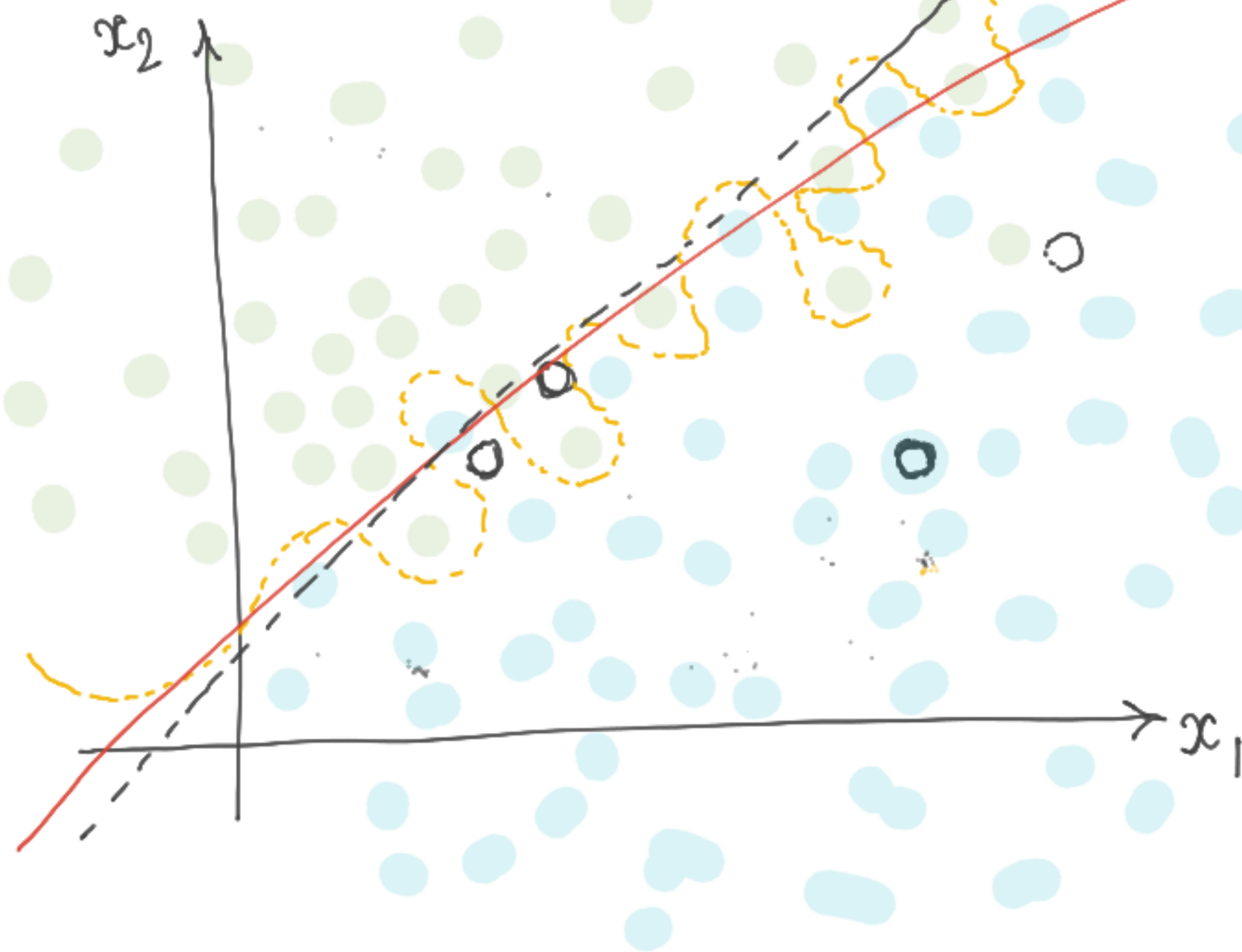- Linear Regression
- Logistic Regression

$\left.\right\}$ parametric

$$y = Xw + \epsilon$$

$$\ln\left(\frac{p}{1-p}\right) = Xw + \epsilon$$

→ parameters

- Decision Trees
- Random Forest
- K Nearest Neighbours

$\left.\right\}$ non-parametric → no statistical qty is being estimated

( Can be used for Regression & Classification )

# K-Nearest Neighbours

- non-parametric
- used for both regression & classification.

$x_2$

$k=10$

$k=1$

$k=6$

- if you want your pt to be influenced by only a small number of closest points; then value of $k$ should be small.

$x_1$

$k=1 \Rightarrow$ ●

$k=3 \Rightarrow$ ●

# KNN - classifier

- smaller value of $k$
  - $\Rightarrow$ complex decision boundary
  - $\Rightarrow$ overfit model
  - $\Rightarrow$ more susceptible to outliers.

- larger value of $k$
  - $\Rightarrow$ simpler (linear-ish) decision boundary
  - $\Rightarrow$ underfit model
  - $\Rightarrow$ ignores subtle patterns present in data.

$\Rightarrow$ How do I select $k$?
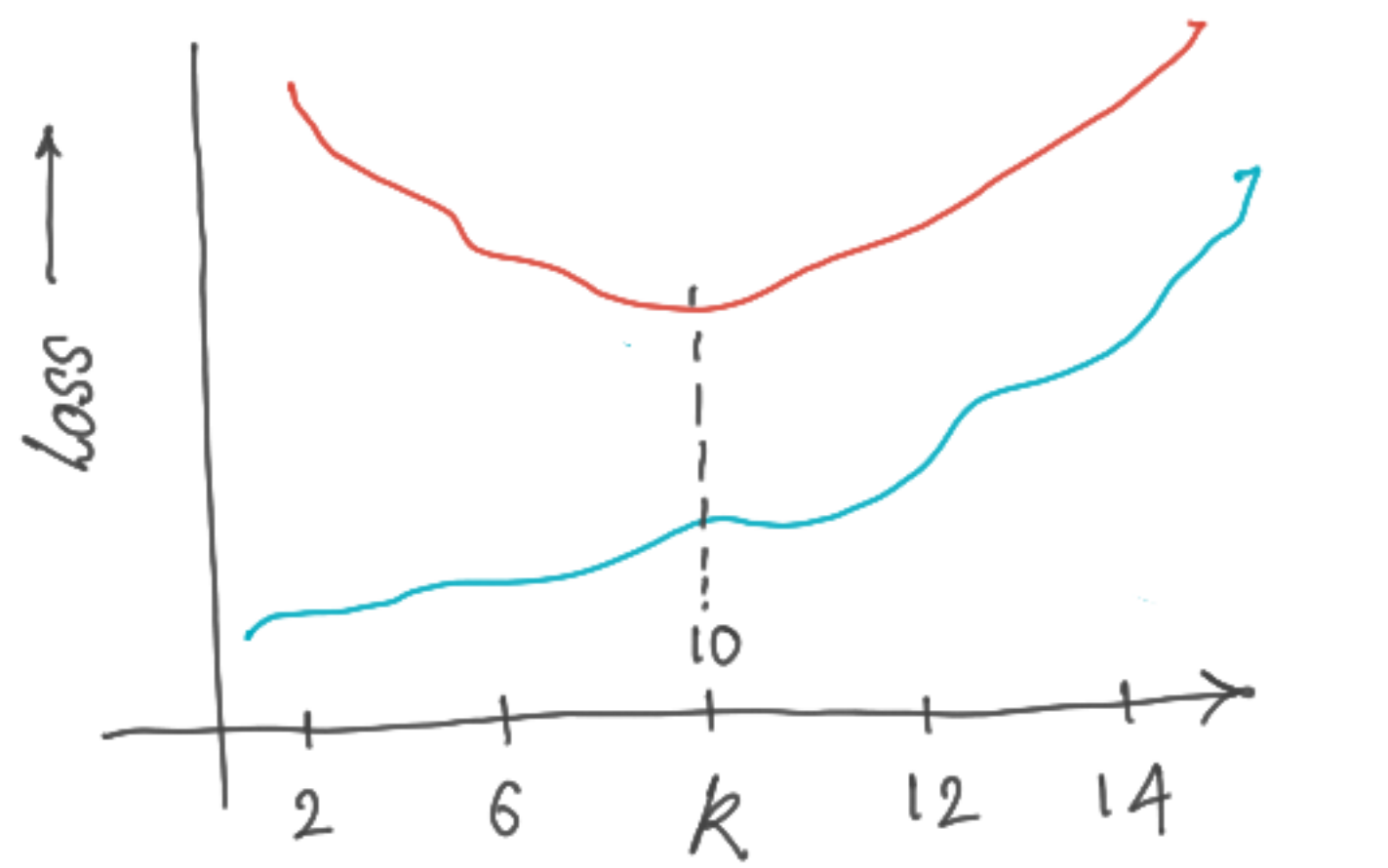
# How to select k ?

- Divide DS → Train, Test
- For k=1 → k=20 ;

$\Rightarrow$ Calculate evaluation metrics for each value of k
  - ↳ accuracy, recall, precision, misclassif. rate
  - ↳ mape, mae, rmse

for train & test separately.

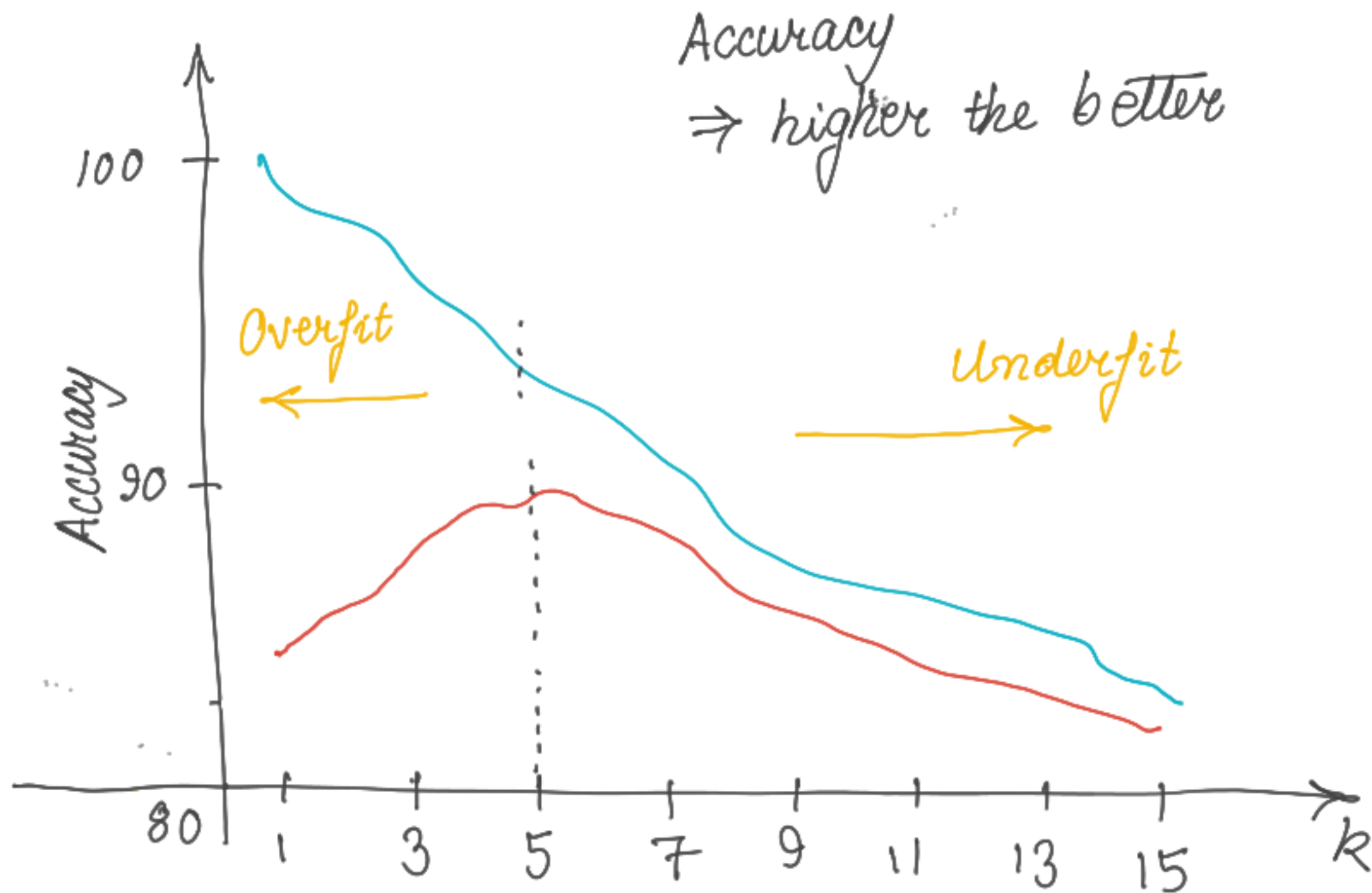- choose that value of k that gives you best eval. metric for test.

**Left figure (Loss vs k):**
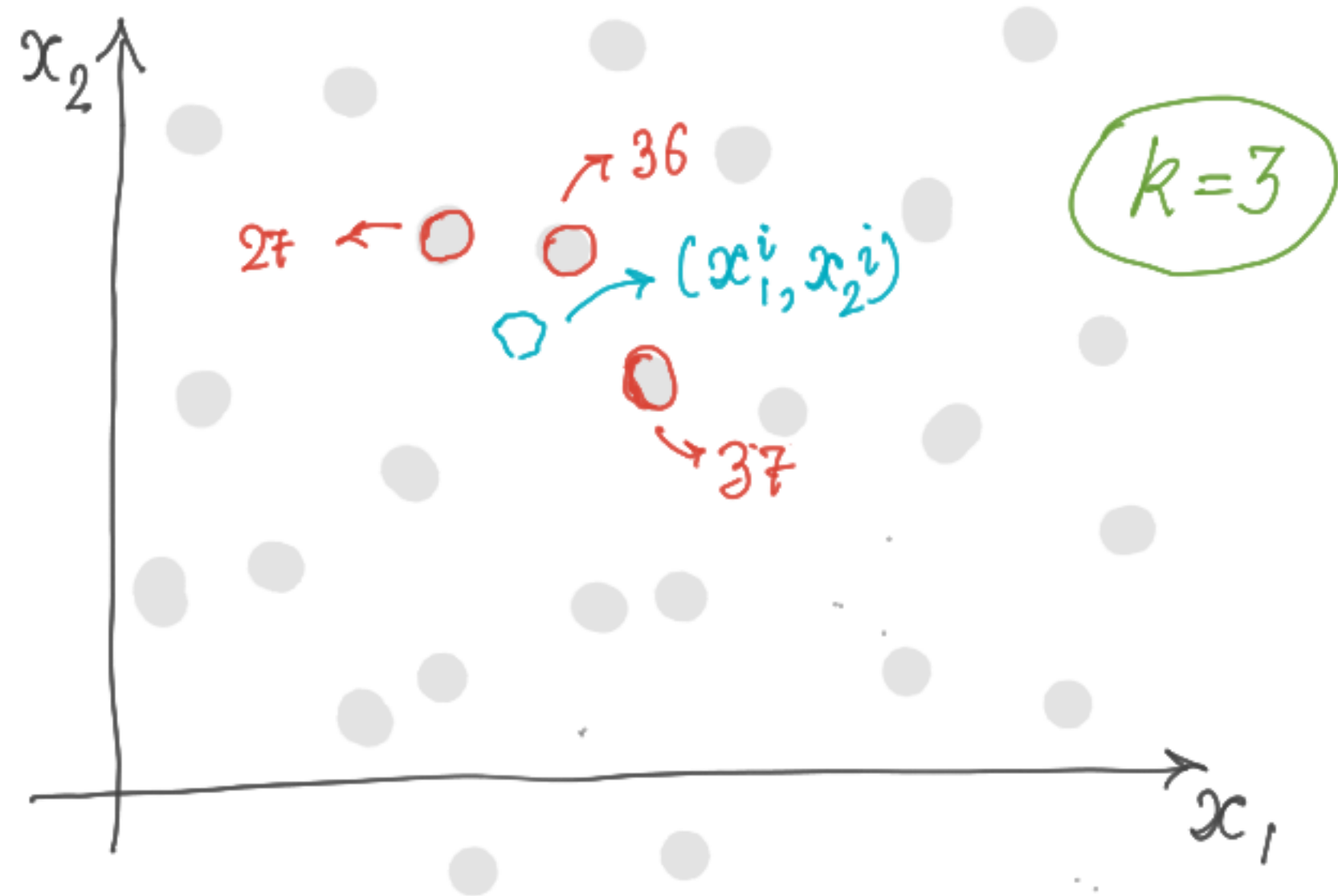
loss

Test (red) — Train (blue)

10

2    6    k    12    14

← Overfit        Underfit →

loss
⇒ Lower the better

**Right figure (Accuracy vs k):**

— Test
— Train.

Accuracy
⇒ higher the better

Accuracy

100

90

80

1    3    5    7    9    11    13    15    k

Overfit ←        → Underfit

# KNN for Regression



$x_1, x_2 \quad y$

1
2
3
$\vdots$
$n$

$\underbrace{\qquad}_{X} \quad \underbrace{\qquad}_{y}$

$(x_1^i, x_2^i) \to y^i\,?$

$k = 3$

$x_2$

$36$

$27$

$(x_1^i, x_2^i)$

$37$

$x_1$

$$y^i = \frac{37 + 36 + 27}{3} = 33.\overline{3}$$

# KNN for Regression

- smaller $k$

$\Rightarrow$ predictions influenced by outliers

$\Rightarrow$ overfit model

- larger $k$

$\Rightarrow$ predictions insensitive to subtle patterns in dataset

$\Rightarrow$ underfit model

We choose $k$ as before $\Rightarrow$ best eval. metric value of test dataset

# Misc

$$\text{Minkowski distance} = \left\{ \sum_{i=1}^{n} |x_i - y_i|^P \right\}^{1/P}$$

$D_P(x, y)$

$$\downarrow$$

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

$$D_2(x, y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

$$= \text{Euclidean distance}$$

$$D_1(x, y) = \sum_{i=1}^{n} |x_i - y_i|$$

$$= \text{Manhattan distance}$$