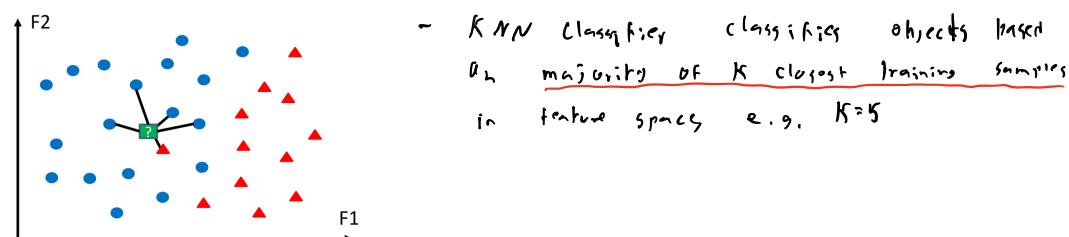
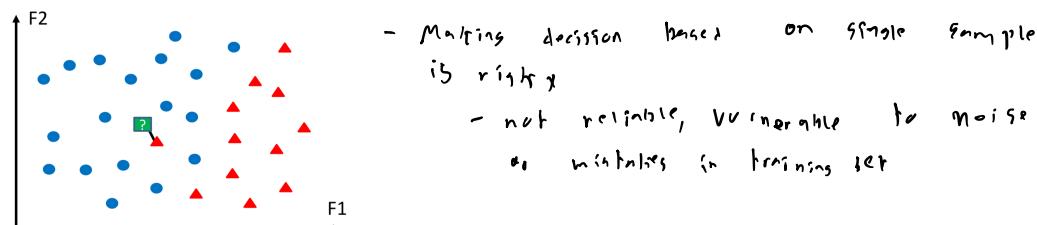
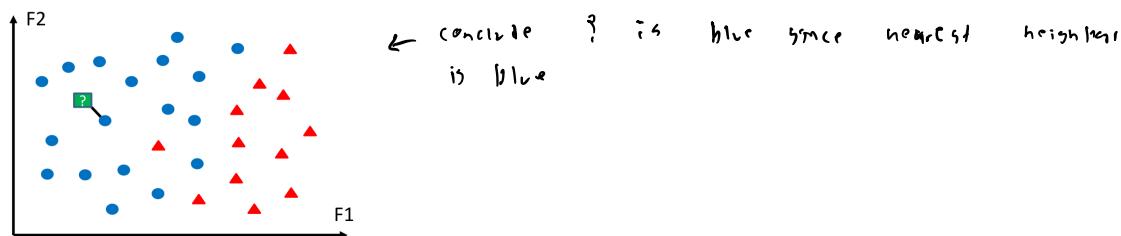


- K-Nearest Neighbor Classifier (KNN classifier)
- Classification also that classifies objects based on closest training samples
 - in the feature space



Feature Table

Training Feature Table:

- In this example we have 6 training data samples (6 different flowers), 4 features and label has 3 classes (3 different types of Iris).

	sepal length	sepal width	petal length	petal width	Label
i					y_i
1	5.3	3.7	1.5	0.2	setosa
2	5	3	2	0.2	setosa
3	7.0	3.2	4.7	1.4	versicolor
4	6.4	3.2	4.5	1.5	versicolor
5	6.3	2.7	4.9	1.8	virginica
6	7.9	3.8	6.4	2	virginica

Training dataset: $\{(x_i, y_i), \dots, (x_n, y_n)\}$:
n data samples used for training

	sepal length	sepal width	petal length	petal width	Label
x_1	5.3	3.7	1.5	0.2	setosa
x_2	5	3	2	0.2	setosa
x_3	7.0	3.2	4.7	1.4	versicolor
	6.4	3.2	4.5	1.5	versicolor
	6.3	2.7	4.9	1.8	virginica
	7.9	3.8	6.4	2	virginica

Defn

- Assuming we have N data samples (observations)
 - have D features
 - C classes

- Feature vector x (Input):** The vector of all features for one single observation: $x \in \mathbb{R}^D$
- Label y (Output):** Value of label for one observation $y \in \{0, 1, 2, \dots, C\}$
- Training dataset:** $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- * Special binary classification case:
 - Number of classes (> 2)
 - Labels: $(0, 1)$

1-Nearst Neighbor Classification

- training dataset: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ with known label
- new sample with unknown label: $(x, y=?)$

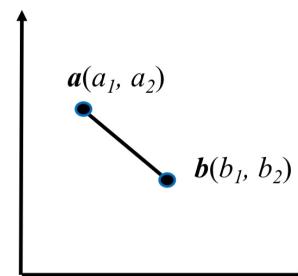
	sepal length	sepal width	petal length	petal width	Label	
x_1	5.3	3.7	1.5	0.2	setosa	y_1
x_2	5	3	2	0.2	setosa	y_2
:	7.0	3.2	4.7	1.4	versicolor	:
:	6.4	3.2	4.5	1.5	versicolor	:
	6.3	2.7	4.9	1.8	virginica	
x_N	7.9	3.8	6.4	2	virginica	y_N
x	7	3.9	5.9	1.3	???	$y=?$

Euclidean Distance: Between 2 points in 2D Space

$$\text{Dis}(a, b) = \|a - b\| = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

Euclidean Distance: For n-dimensional Space:

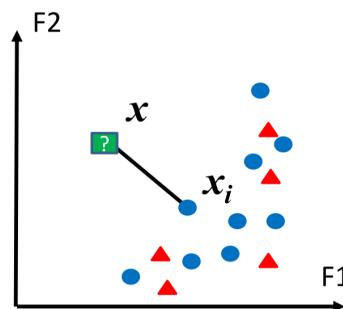
$$\text{Dis}(x_i, x) = \|x_i - x\| = \sqrt{\sum_{d=1}^n (x_{id} - x_d)^2}$$



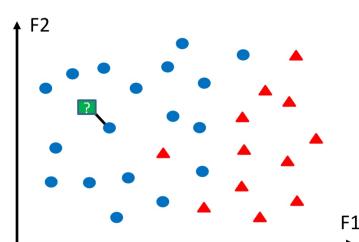
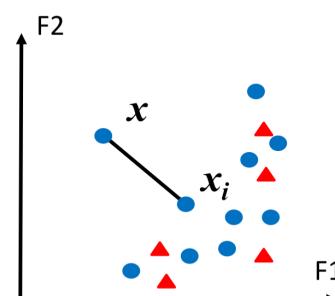
1-Nearst Neighbor Classification

- training dataset: $\{(x_1, y_1), \dots, (x_n, y_n)\}$
- new sample, unknown label: $(x, y=?)$
- Euclidean Distance between x and a known training point x_i :
- Nearest Neighbor to x ($NN(x)$)

$$\begin{aligned} NN(x) &= \arg \min_{i \in \{1, \dots, n\}} (\text{Dis}(x_i, x)) \\ &= \arg \min_{i \in \{1, \dots, n\}} (\|x_i - x\|^2) \\ &= \arg \min_{i \in \{1, \dots, n\}} \left(\sum_{d=1}^n (x_{id} - x_d)^2 \right) \end{aligned}$$

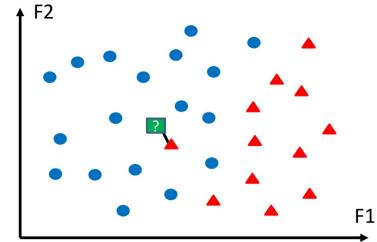


- After finding NN of x , label of x will be determined as label of its NN
- $y = y_{NN}(x)$

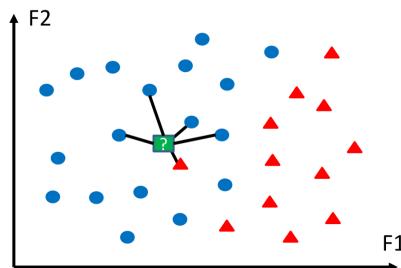


↪ NN works on this example, but will it always work?

- ← In this ex. NN classification would be wrong
- making decision on a single sample is risky



K-NN Classification



- classifies objects based on K closest training samples in the feature space e.g. $K=5$
- 1st NN: $NN_1(x)$
- 2nd NN: $NN_2(x)$
- ⋮
- K^{th} NN: $NN_K(x)$

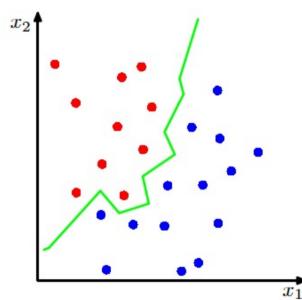
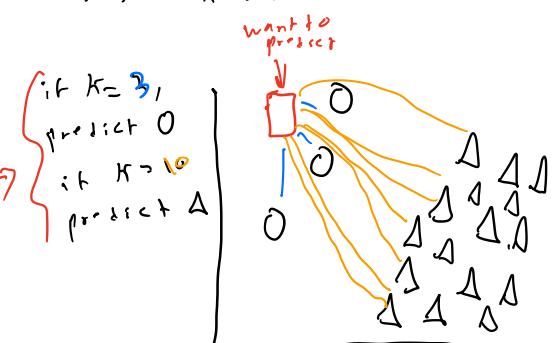
set of K -NNs: $KNN(x) = \{NN_1(x), \dots, NN_K(x)\}$

↑
Out of 5 NN,
4 blue, 1 red, ∴
prediction is blue

- **Classification Rule**: Select the label with the majority in $KNN(x)$

- Benefits large K
- ignoring the effect of outliers
- Better decision making

- Benefits of a small K
- avoid overfitting
 - simple model
 - low computational complexity
 - No bias towards popular labels



Decision Boundary

- At possible point in the space, can determine its label using KNN rule
- This gives us **decision boundary** that partitions the space into different regions

Advantages KNN Classifier

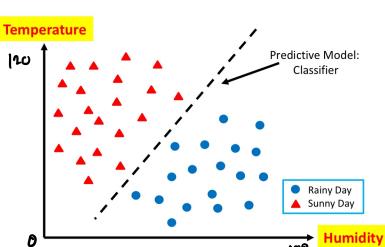
- easy & simple to implement
- low computational complexity
- Computationally intensive for large scale problems!
- inefficient for Big Data (needs to find millions of distances → find shortest each time)

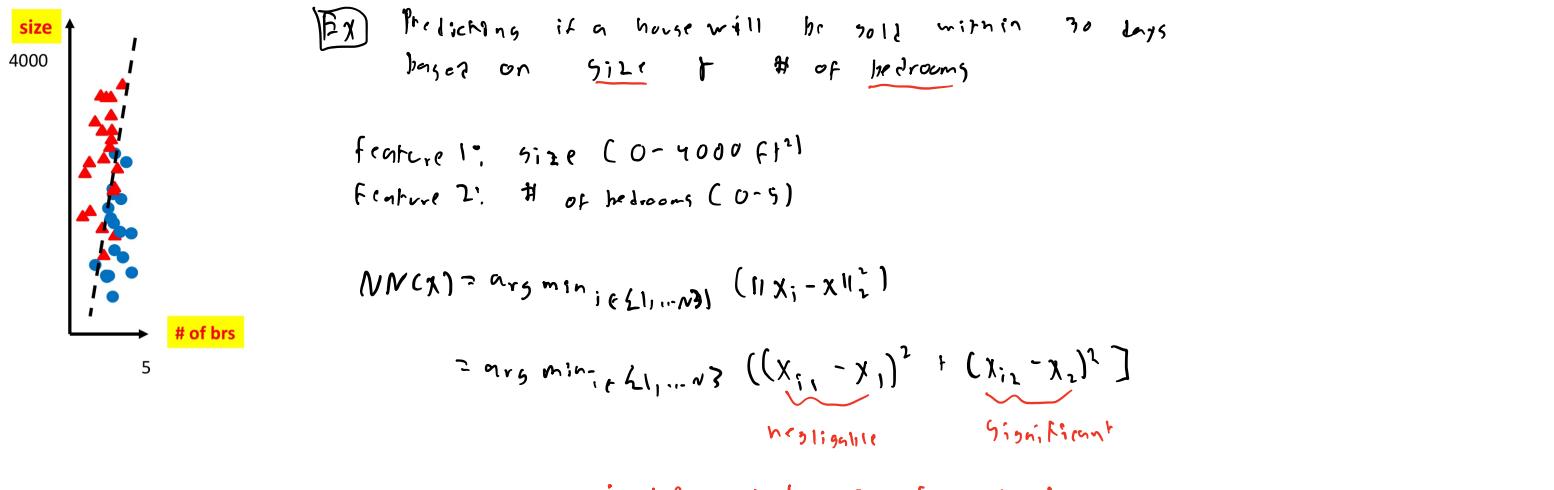
- choosing best K is challenging

Practical hint about KNN

- distances depend on unit of feature

← in ex. both features almost in same range so no problem





Easy way to scale features:

$$F1_{\text{Scaled}} = \text{size} / \max(\text{size})$$

$$F2_{\text{Scaled}} = \# \text{ of bedrooms} / \max \# \text{ of rooms}$$

After normalization:

$$(0 \leq F1_{\text{Scaled}} \leq 1)$$

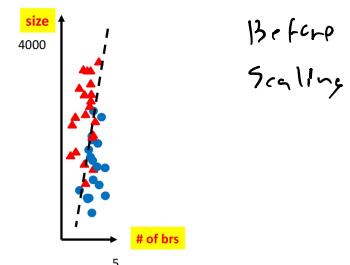
$$(0 \leq F2_{\text{Scaled}} \leq 1)$$

• Before Scaling:

	Number of Bedrooms	Size (sqr footage)	Label
1	1	1300	Sold
2	3	2400	Not-Sold
3	3	2270	Not-Sold
4	2	1450	Sold
5	2	1400	Sold
6	4	2900	Not-Sold

• After Scaling:

	Number of Bedrooms	Size (sqr footage)	Label
1	1/4	1300/2900	Sold
2	3/4	2400/2900	Not-Sold
3	3/4	2270/2900	Not-Sold
4	2/4	1450/2900	Sold
5	2/4	1400/2900	Sold
6	4/4	2900/2900	Not-Sold



More advanced technique for normalization:

- We can Normalize data samples to have zero mean and unit standard deviation in each dimension (i.e., for each feature)
 - ~ We should normalize for each column of the feature table individually
- Compute mean & std deviation for each feature
 $(x_{nd} \text{ is the } d\text{th feature of } n\text{th data sample})$:

$$\text{mean}(x_d) = \bar{x}_d = \frac{1}{N} \sum_{n=1}^N x_{nd} \quad \text{std}(x_d) = s_d = \sqrt{\frac{1}{N-1} \sum_{n=1}^N (x_{nd} - \bar{x}_d)^2}$$

- Scale the feature accordingly

$$x_{nd} \leftarrow \frac{x_{nd} - \bar{x}_d}{s_d}$$