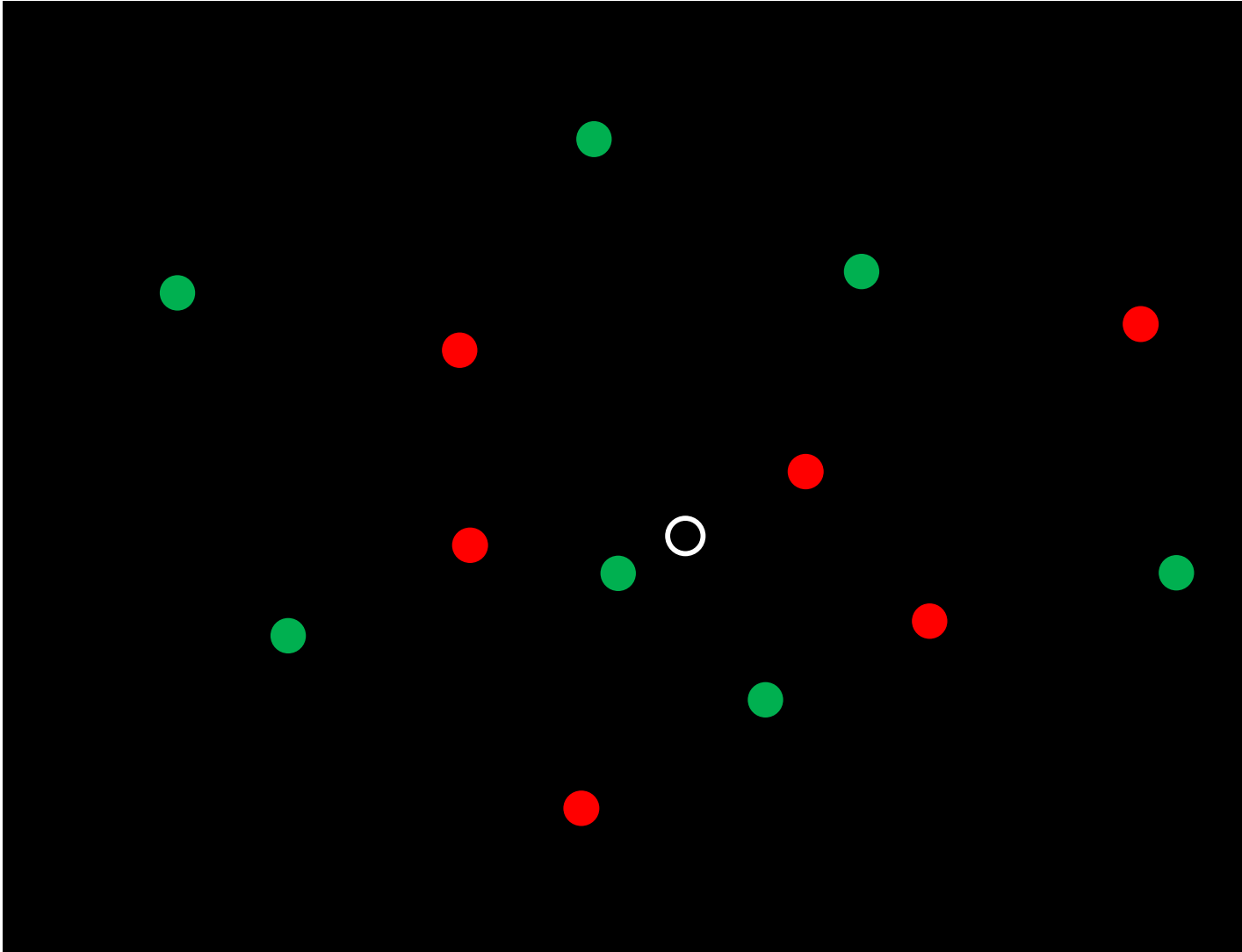# KNN CLASSIFICATION

Ratchainant Thammasudjarit, Ph.D.
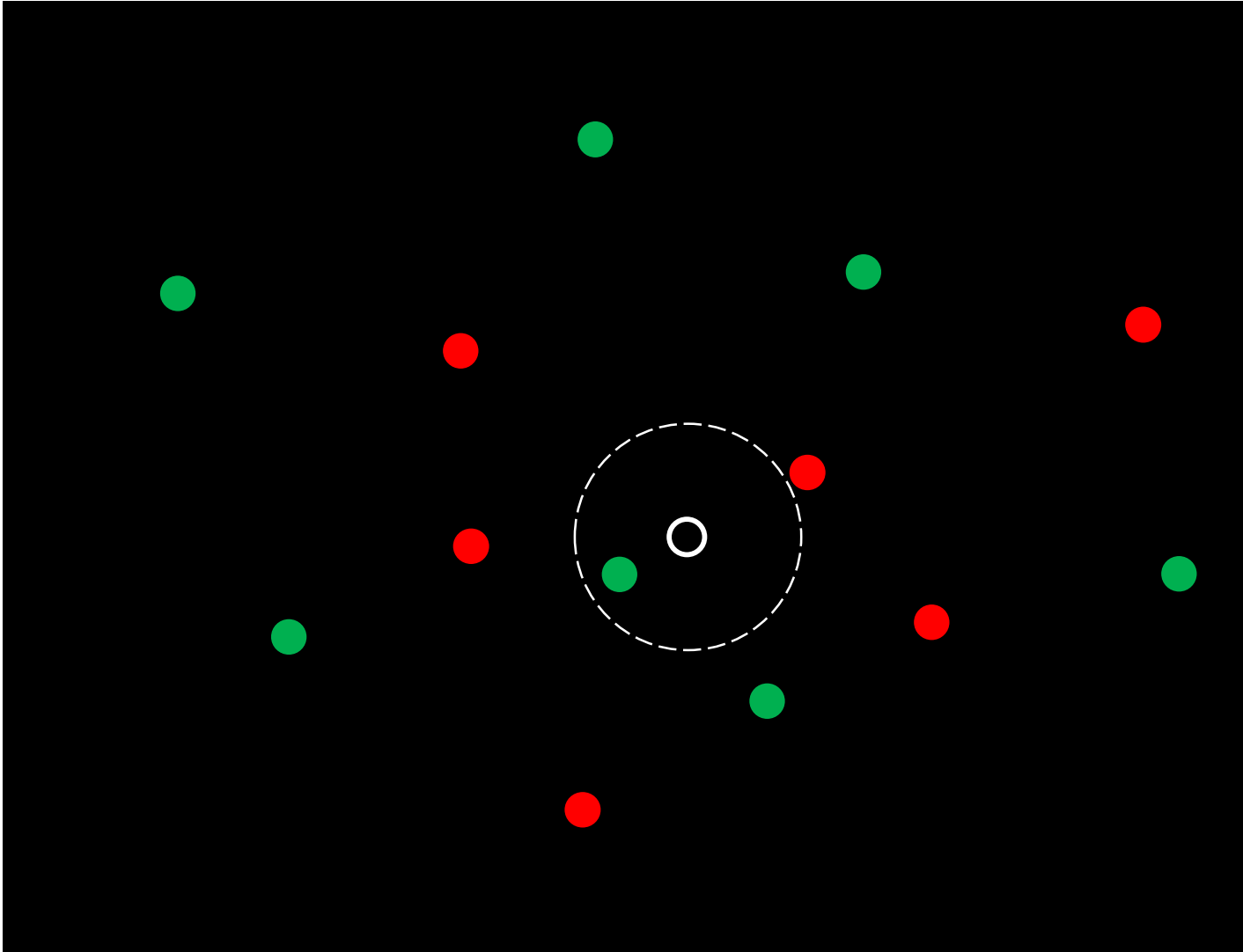
○ Peer Pressure



# Concepts

Which color should be for the white data point
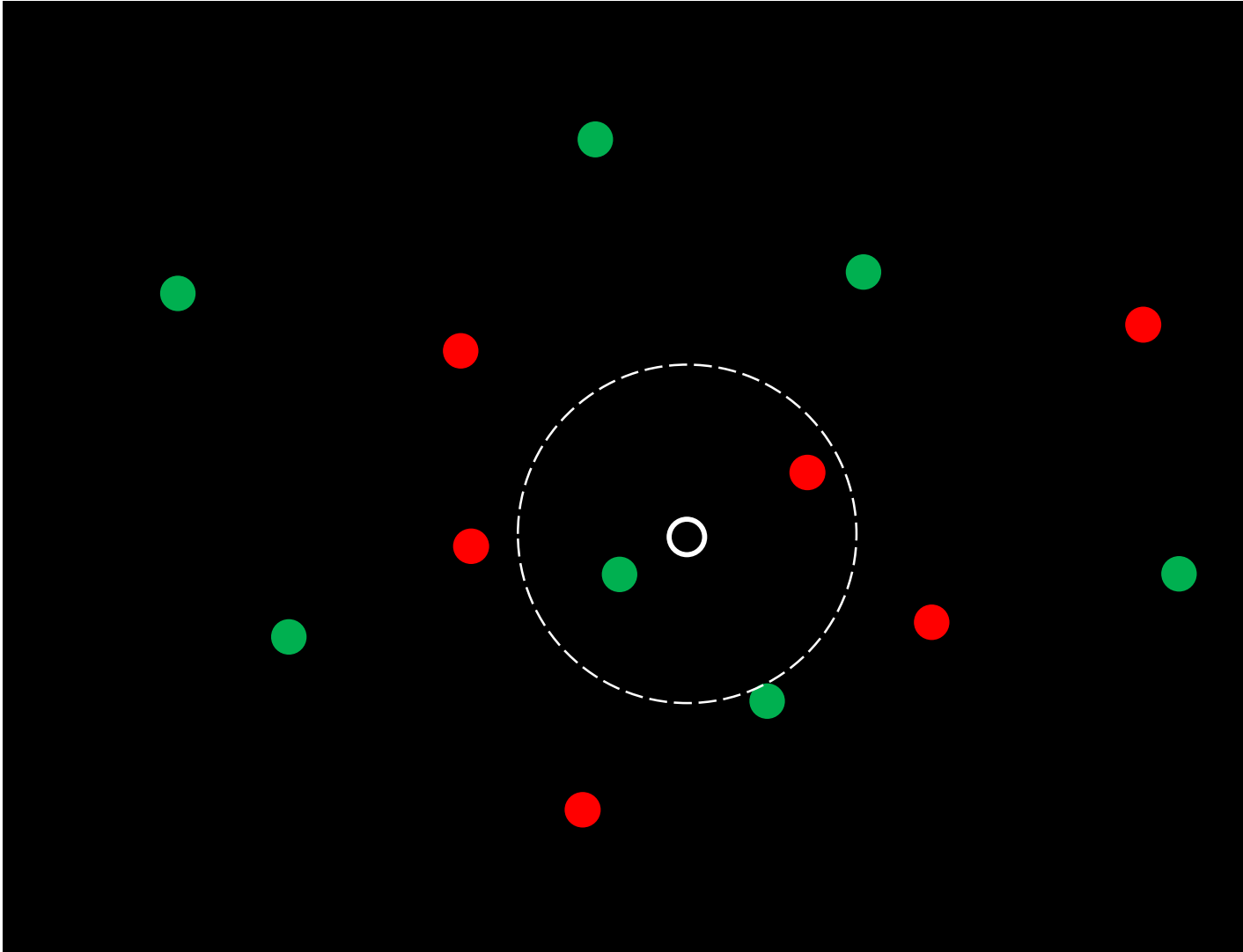
◦ 1 nearest data point

   ◦ Green is majority: Prediction is green



# Concepts

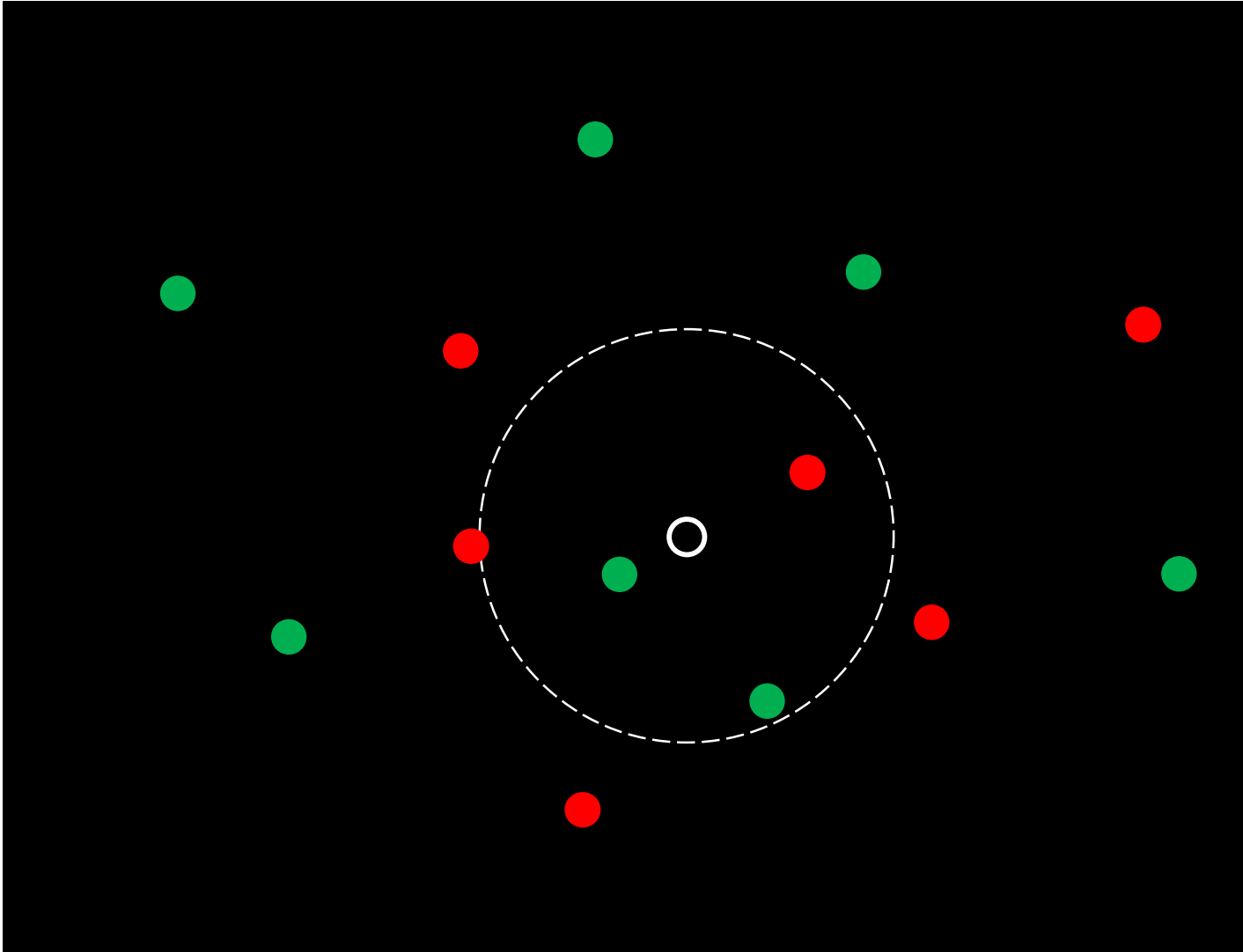Which color should be for the white data point

- 2 nearest data points
  - Cannot make decision



# Concepts

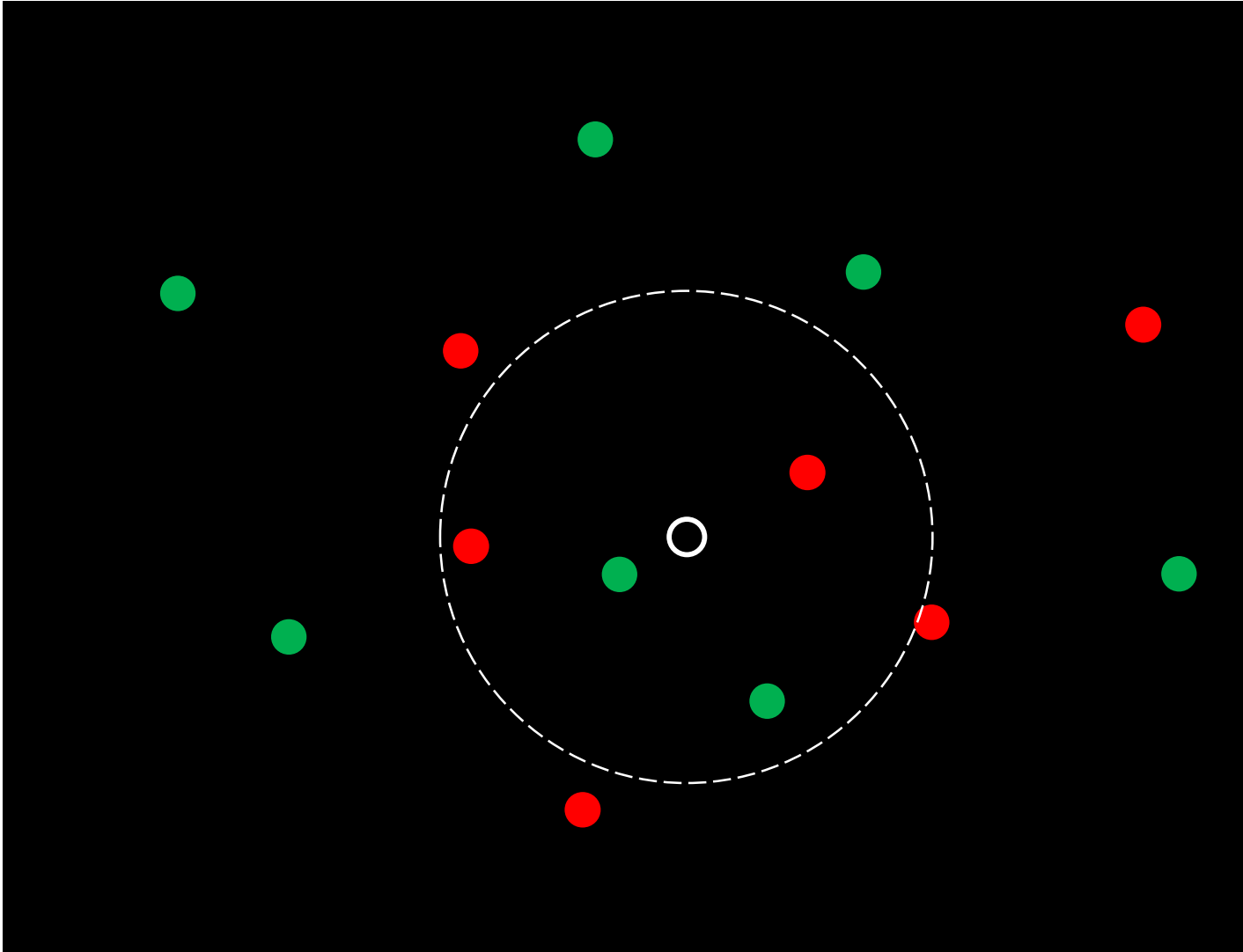Which color should be for the white data point

- 3 nearest data points
  - Green is majority: Prediction is green



# Concepts

Which color should be for the white data point
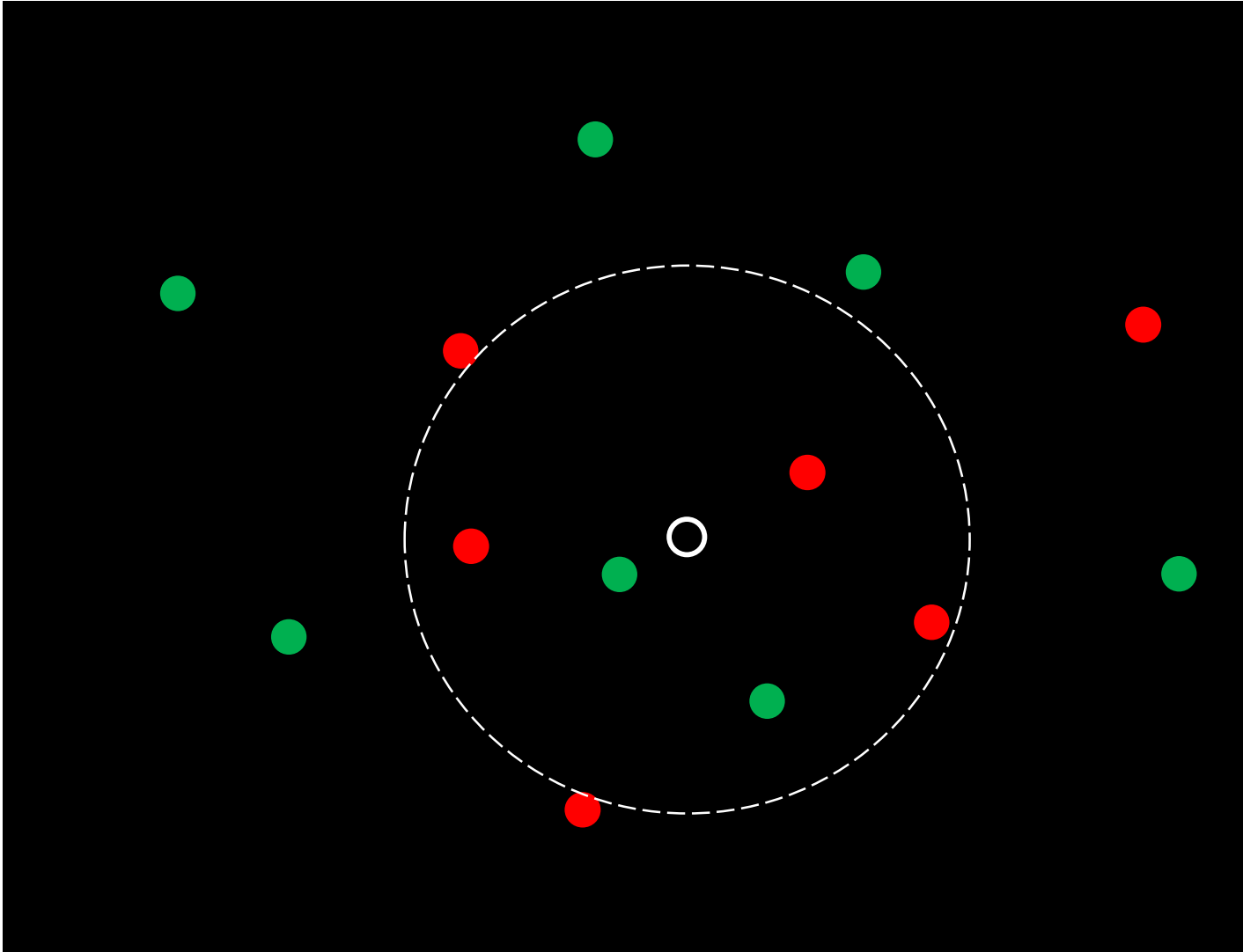
◦ 4 nearest data points
  ◦ Cannot make decision



# Concepts

Which color should be for the white data point

◦ 5 nearest data points
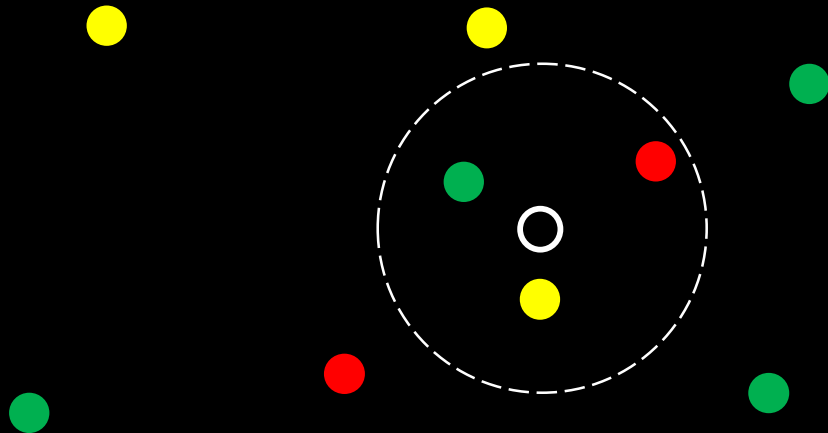  ◦ Red is majority: Prediction is red



# Concepts

Which color should be for the white data point

◦ For binary classification

  ◦ $k$ is recommended to be the odd number

◦ For multi-classification

  ◦ $k$ is recommended to be the odd number and at least $2C + 1$

**3 Classes:** If $k = 3$, this situation is undesirable



## Concepts

The number of nearest neighbor $k$ plays important role for prediction

In common practice, the $k$ is set to be some odd number

◦ For binary classification

  ◦ $k$ is recommended to be the odd number

◦ For multi-classification ($C$ classes)

  ◦ $k$ is recommended to be the odd number and at least $2C + 1$

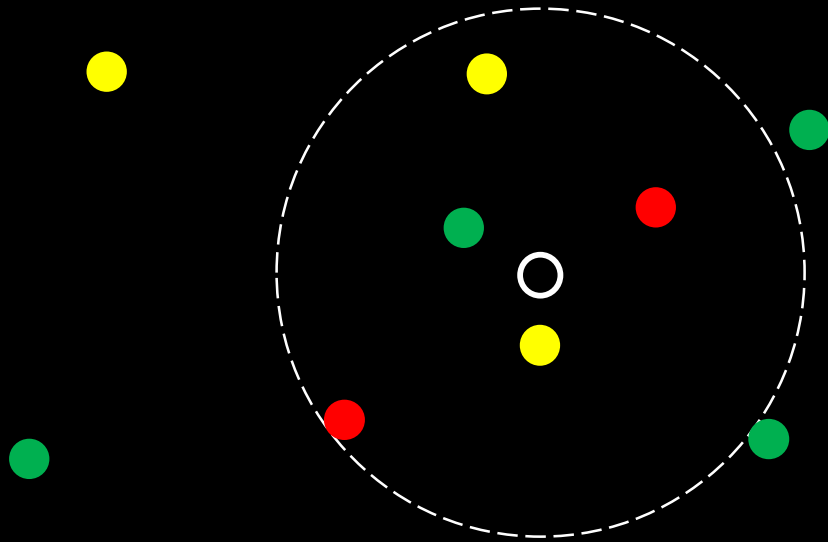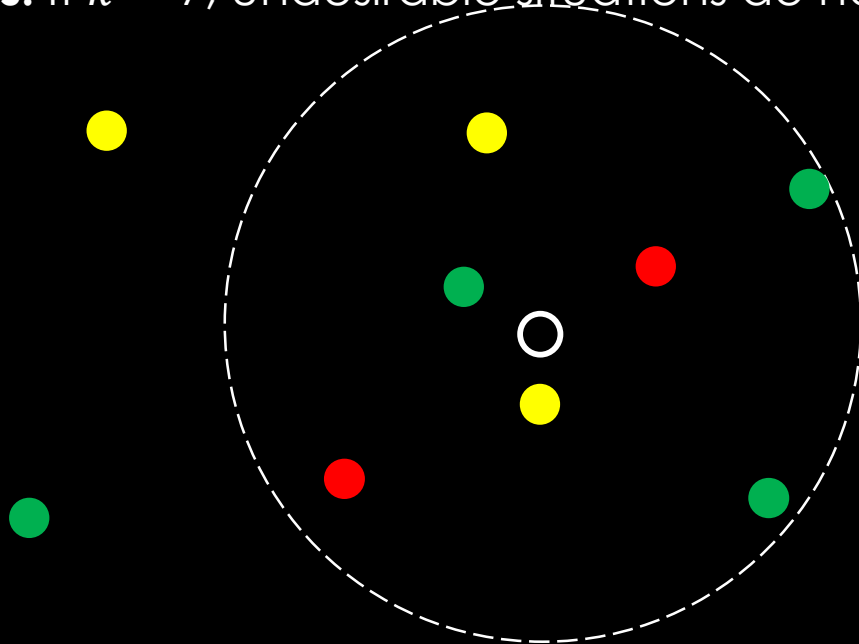**3 Classes:** If $k = 5$, this situation is undesirable

# Concepts

The number of nearest neighbor $k$ plays important role for prediction

In common practice, the $k$ is set to be some odd number

- For binary classification
  - $k$ is recommended to be the odd number

- For multi-classification
  - $k$ is recommended to be the odd number and at least $2C + 1$

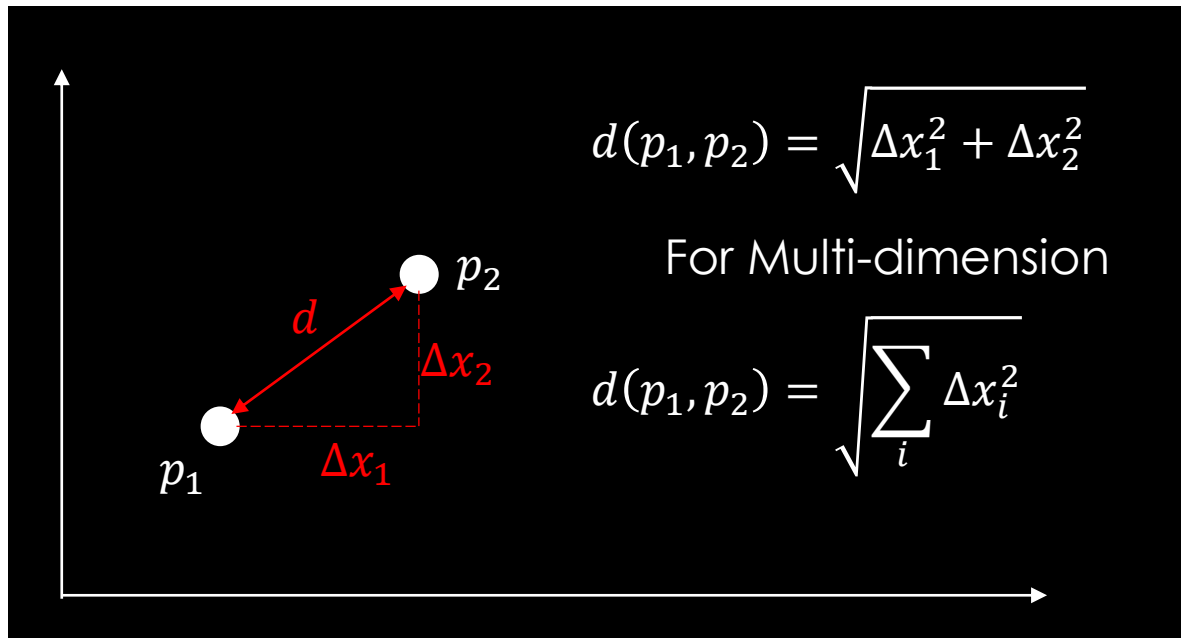**3 Classes:** If $k = 7$, undesirable situations do not exist

## Concepts

The number of nearest neighbor $k$ plays important role for prediction

In common practice, the $k$ is set to be some odd number

◦ Euclidean Distance (0 to ∞)

  ◦ 0: Exactly the same

  ◦ ∞: Completely different

◦ In practice, ∞ of Euclidean distance is relatively impossible

◦ Magnitude does matter

◦ All features are recommended to have the same scale

$$d(p_1, p_2) = \sqrt{\Delta x_1^2 + \Delta x_2^2}$$

For Multi-dimension

$$d(p_1, p_2) = \sqrt{\sum_i \Delta x_i^2}$$

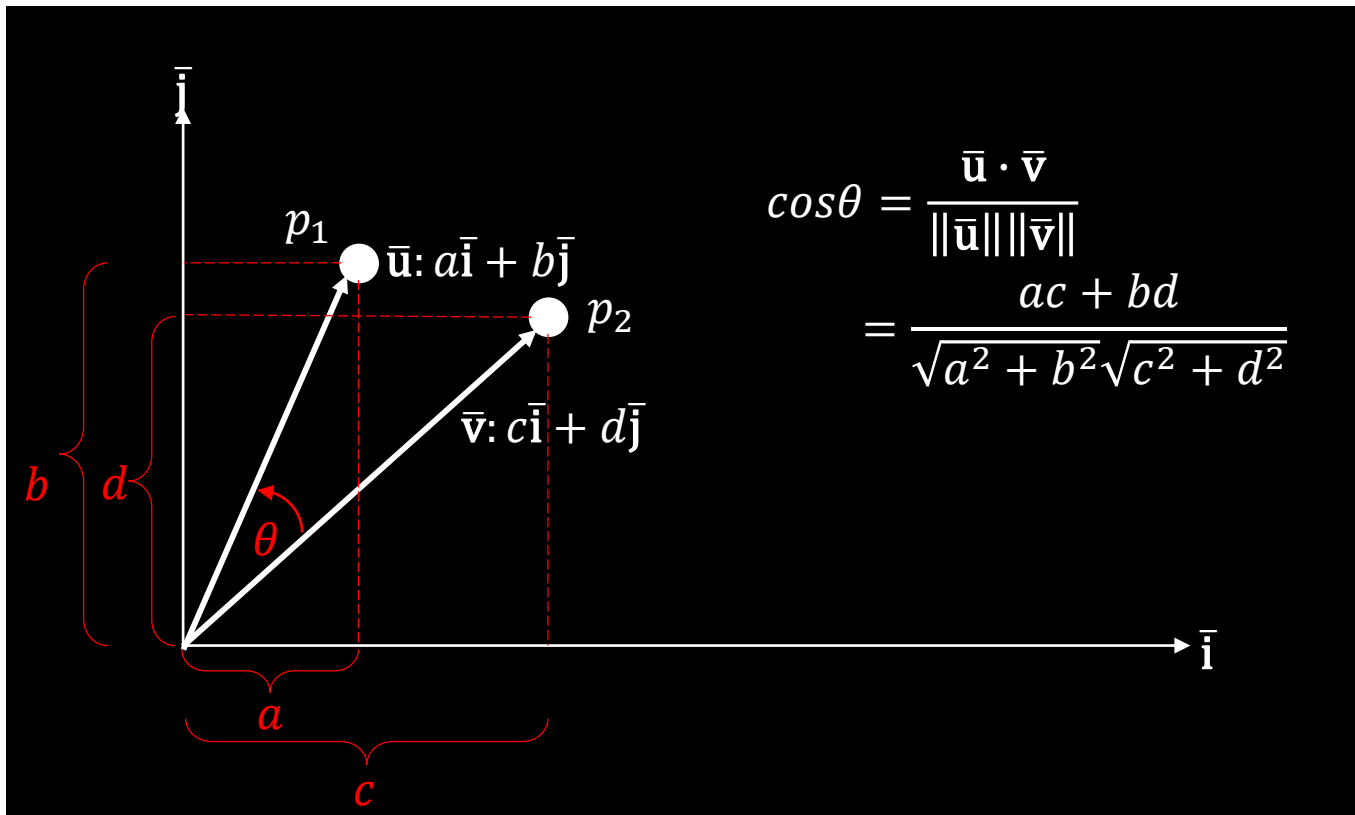$p_2$

$d$

$\Delta x_2$

$p_1$

$\Delta x_1$

# Nearest Neighbors

Determining the $k$ nearest neighbors relies on the **similarity measure**

There are several similarity measures

- Cosine similarity (-1 to 1)
  - 1: Similar
  - 0: Unable to detect similarity
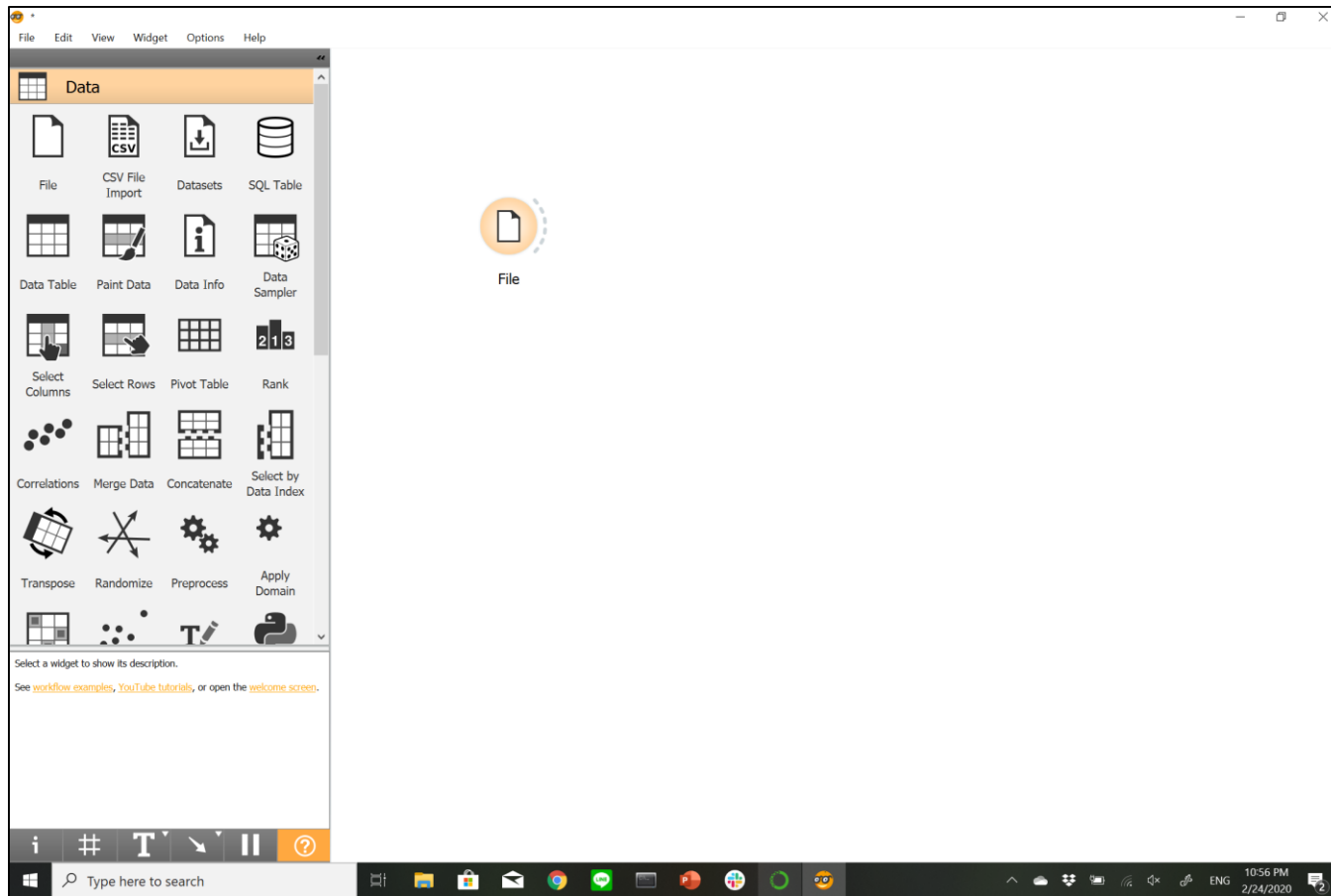  - -1: Unsimilar
  - Magnitude does not matter

$$cos\theta = \frac{\bar{\mathbf{u}} \cdot \bar{\mathbf{v}}}{\|\bar{\mathbf{u}}\|\|\bar{\mathbf{v}}\|}$$

$$= \frac{ac + bd}{\sqrt{a^2 + b^2}\sqrt{c^2 + d^2}}$$



# Nearest Neighbors

Determining the $k$ nearest neighbors relies on the **similarity measure**
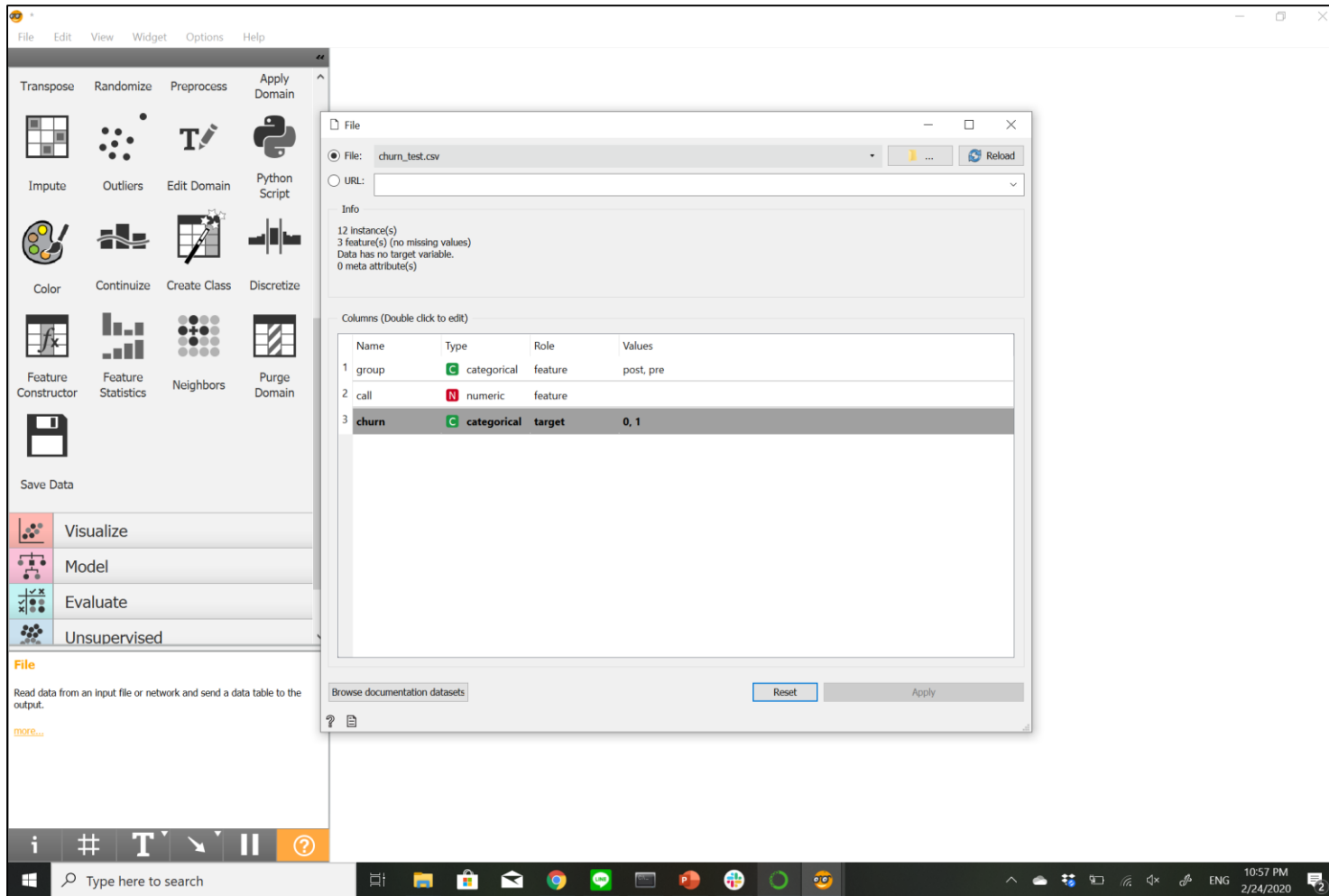
There are several similarity measures

○ Import data



KNN Classification in Orange
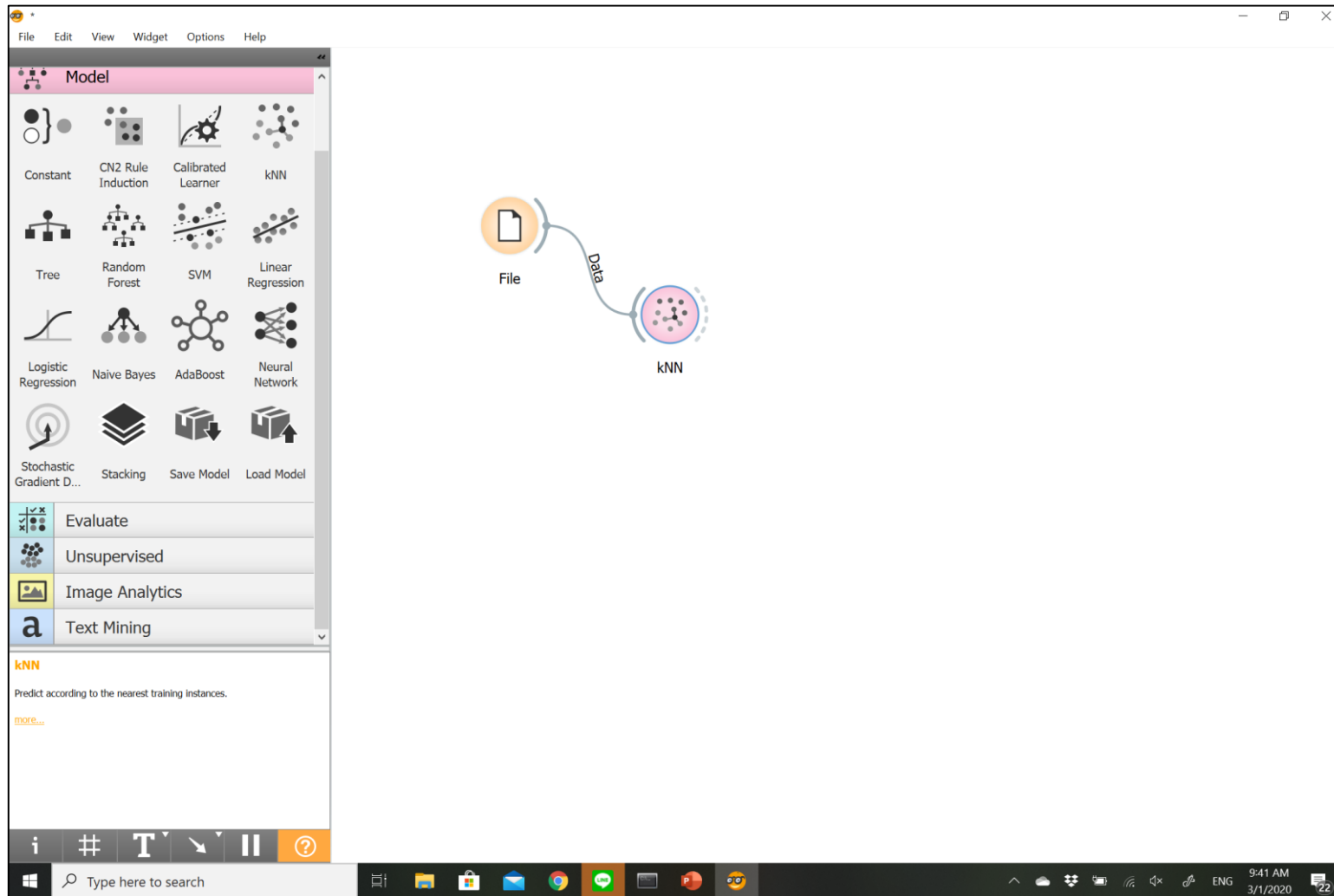
Build your model in 10 seconds

○ Identify features and target



KNN Classification in Orange
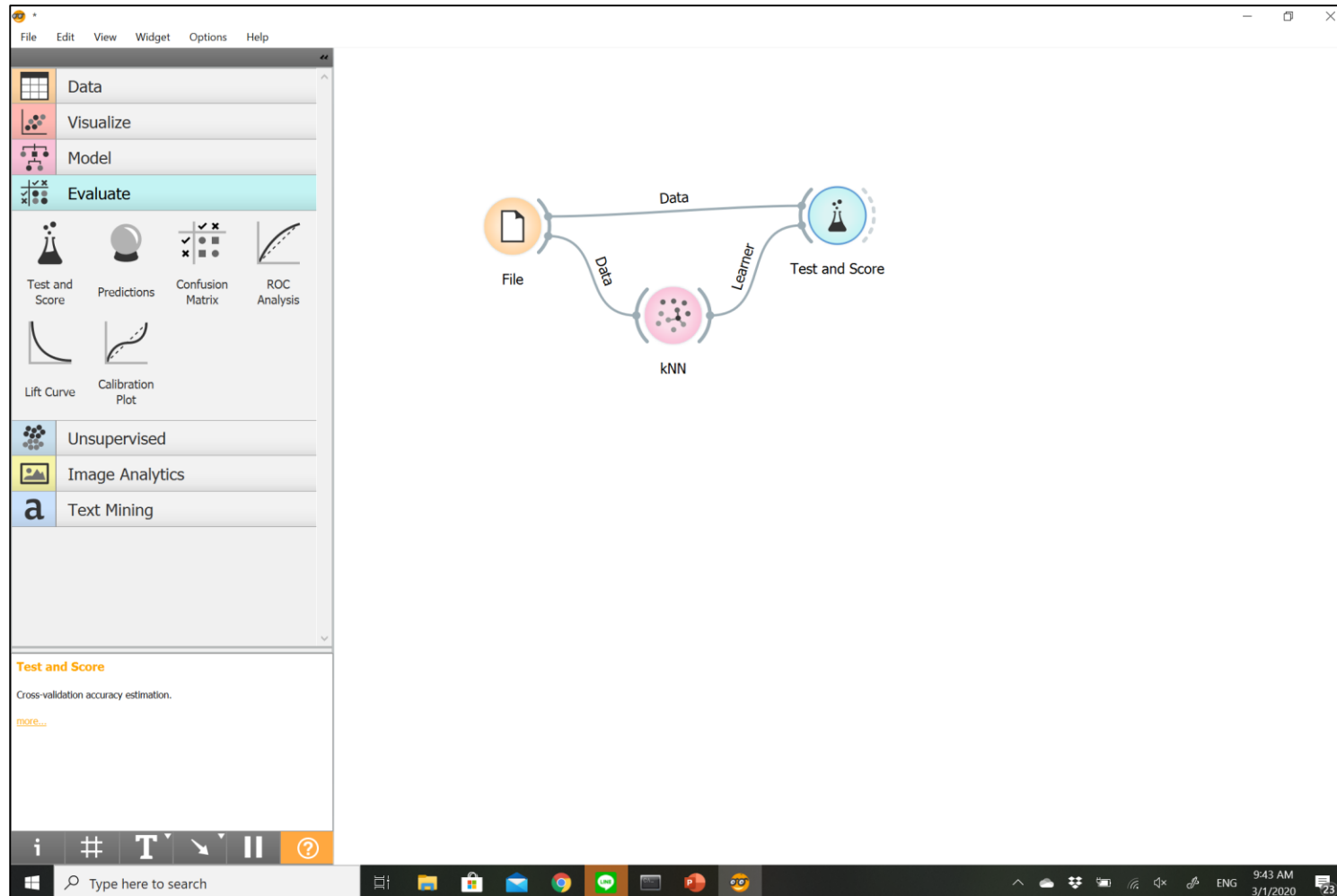
Build your model in 10 seconds

○ Add model



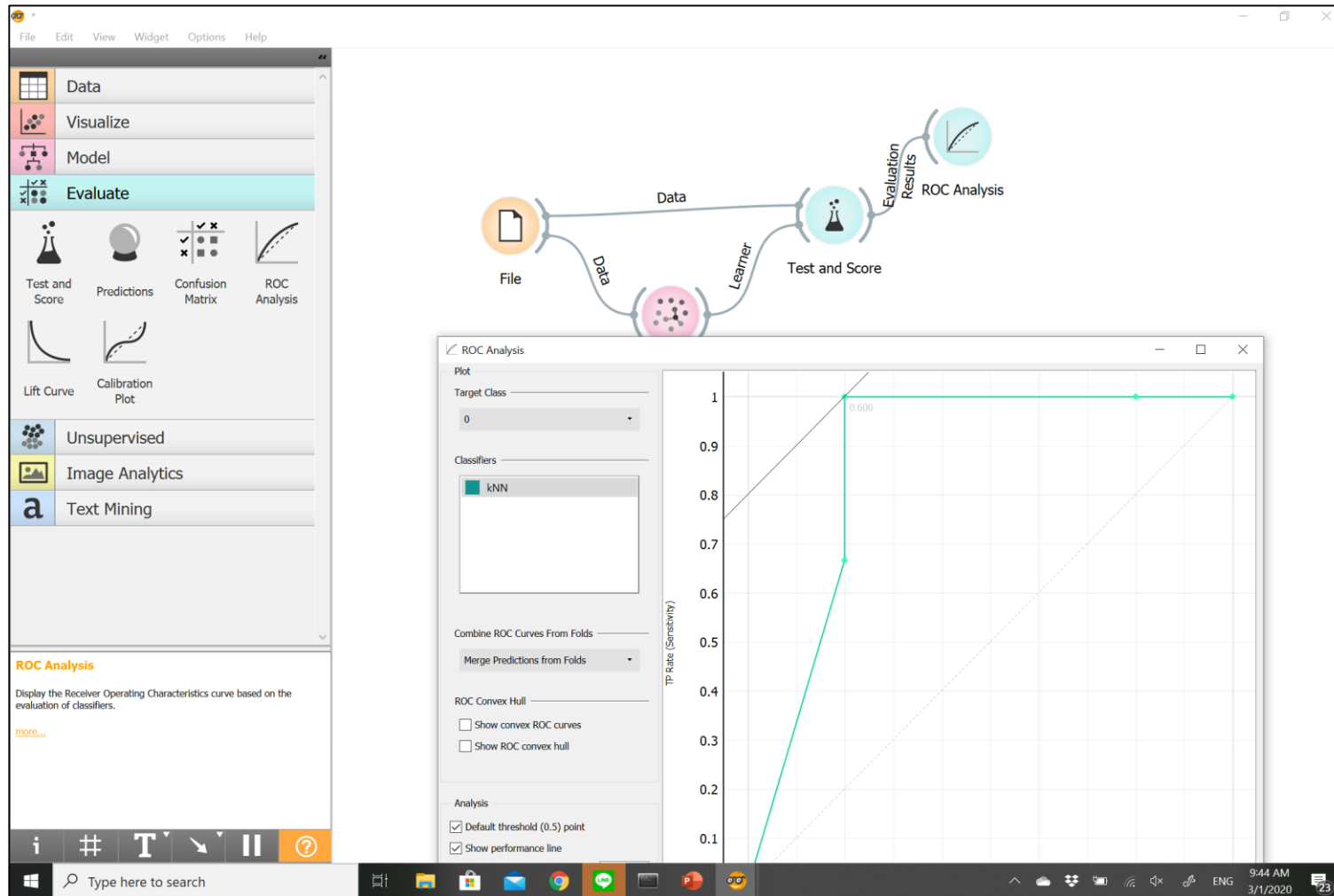# KNN Classification in Orange

Build your model in 10 seconds

○ Evaluate Model



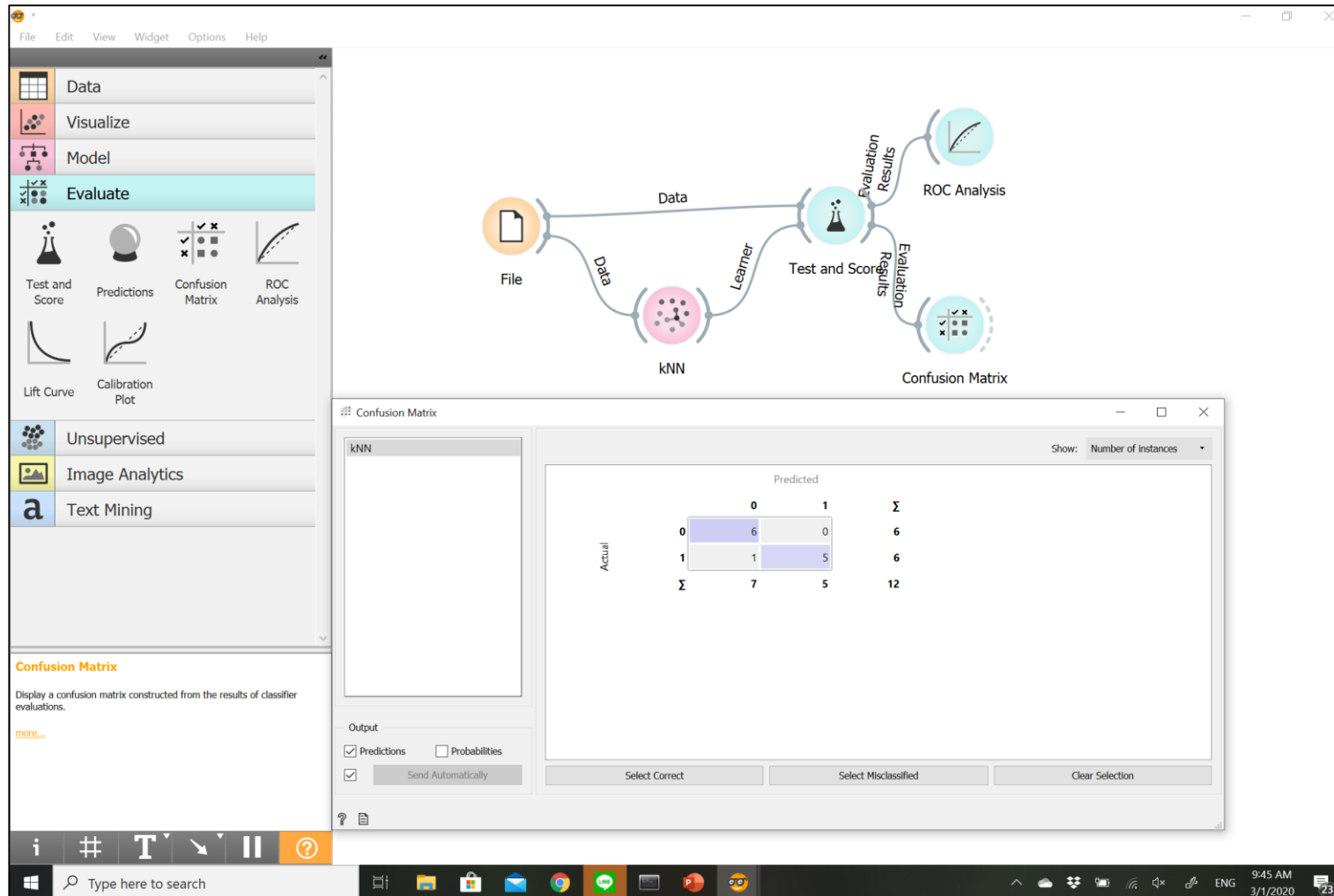# KNN Classification in Orange

Build your model in 10 seconds

○ ROC Plot



KNN Classification in Orange

Build your model in 10 seconds

○ Confusion Matrix



KNN Classification in Orange

Build your model in 10 seconds

◦ Build your model using Kaggle dataset

# Exercise

Our sample data is very small for demonstration purpose

Now it is the time to work with the dataset churn prediction of telecom in Kaggle Competition