# Small data technique I

# Lecture 13

Changho Suh

October 1, 2021

# Decision trees (DTs)

# Recap: DNNs

Work well with **enough data**.

Otherwise, we may face: <span style="color:red">Overfitting</span> problem

This motivates <span style="color:blue">simplifying DNNs</span>, being tailored for tasks of interest.

# Recap: CNNs

A model specialized for image data

Two key building blocks:

1. **Conv** layer (*mimicking* neurons in *visual cortex*)
2. **Pooling** layer (*mainly for reducing complexity*)

Design principles: As a network gets deeper:

1. Feature map size gets smaller;

2. # of feature maps gets bigger.

# Recap: RNNs

A model specialized for time-series data

Two key building blocks:
1. **Recurrent neurons**
2. **Memory cell**

**Basic RNNs**: Trained via truncated BTTP.

**LSTM:** Offers great performance and faster training.

# Recap: Tensorflow coding for RNNs

```python
from tensorflow.keras.datasets import imdb
from tensorflow.keras.preprocessing.sequence import pad_sequences
from keras.models import Sequential
from keras.layers import Dense, Embedding, SimpleRNN, LSTM

(X_train, y_train), (X_test, y_test) = imdb.load_data(num_words=10000)

# Preprocessing
X_train_pad = pad_sequences(X_train, value=0, padding='post', maxlen=256)

# Basic RNN
model = Sequential()
model.add(Embedding(num_words, 100, input_shape=(None,)))
model.add(SimpleRNN(128))
model.add(Dense(1, activation='sigmoid'))

# LSTM
model_LSTM = Sequential()
model_LSTM.add(Embedding(num_words, 100, input_shape=(None,)))
model_LSTM.add(LSTM(128))
model_LSTM.add(Dense(1, activation='sigmoid'))
```

# Questions

1. What if still <span style="color:red">unsatisfactory</span> performances?

   A better approach for the <span style="color:blue">small data</span> regime?


2. What about <span style="color:red">interpretability</span> of DNNs?

# Today's lectures

Will explore a technique that may enable a better performance for the small-data regime, as well as offer model interpretability:

## Random forests (RFs)

The **most powerful** ML algorithm in industry

# Outline of today's lectures

Specifically we will study:

1. **Decision trees (DTs)**:
   Fundamental components of RFs

2. **Ensemble learning**:
A generic technique that includes RFs as a special case.

3. **RFs** in depth

# Focus of Lecture 13

Specifically we will study:

1. **Decision trees (DTs)**:

   Fundamental components of RFs

2. **Ensemble learning**:

A generic technique that includes RFs as a special case.

3. **RFs** in depth

# A motivating example

Classification on **Iris** dataset:

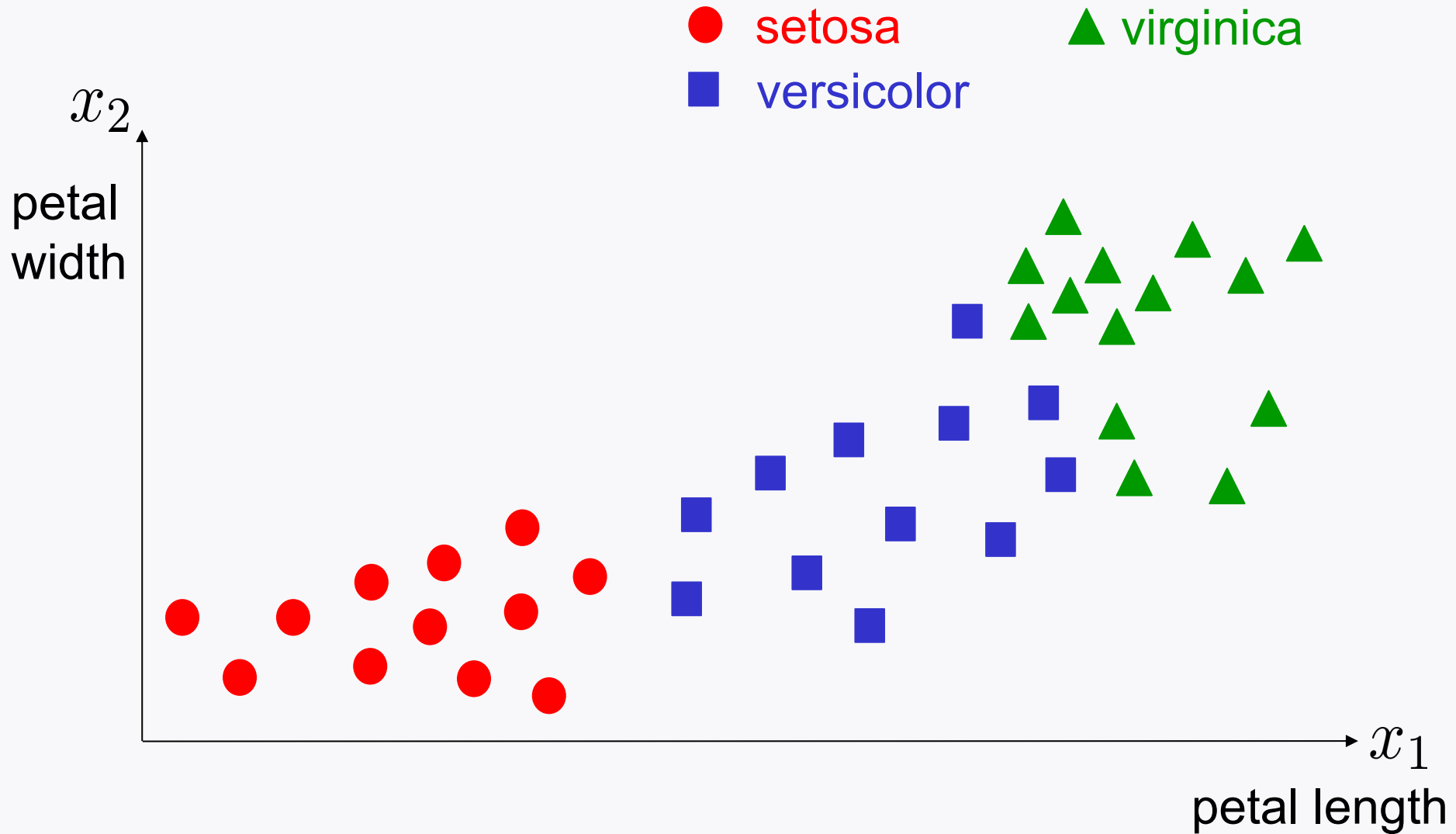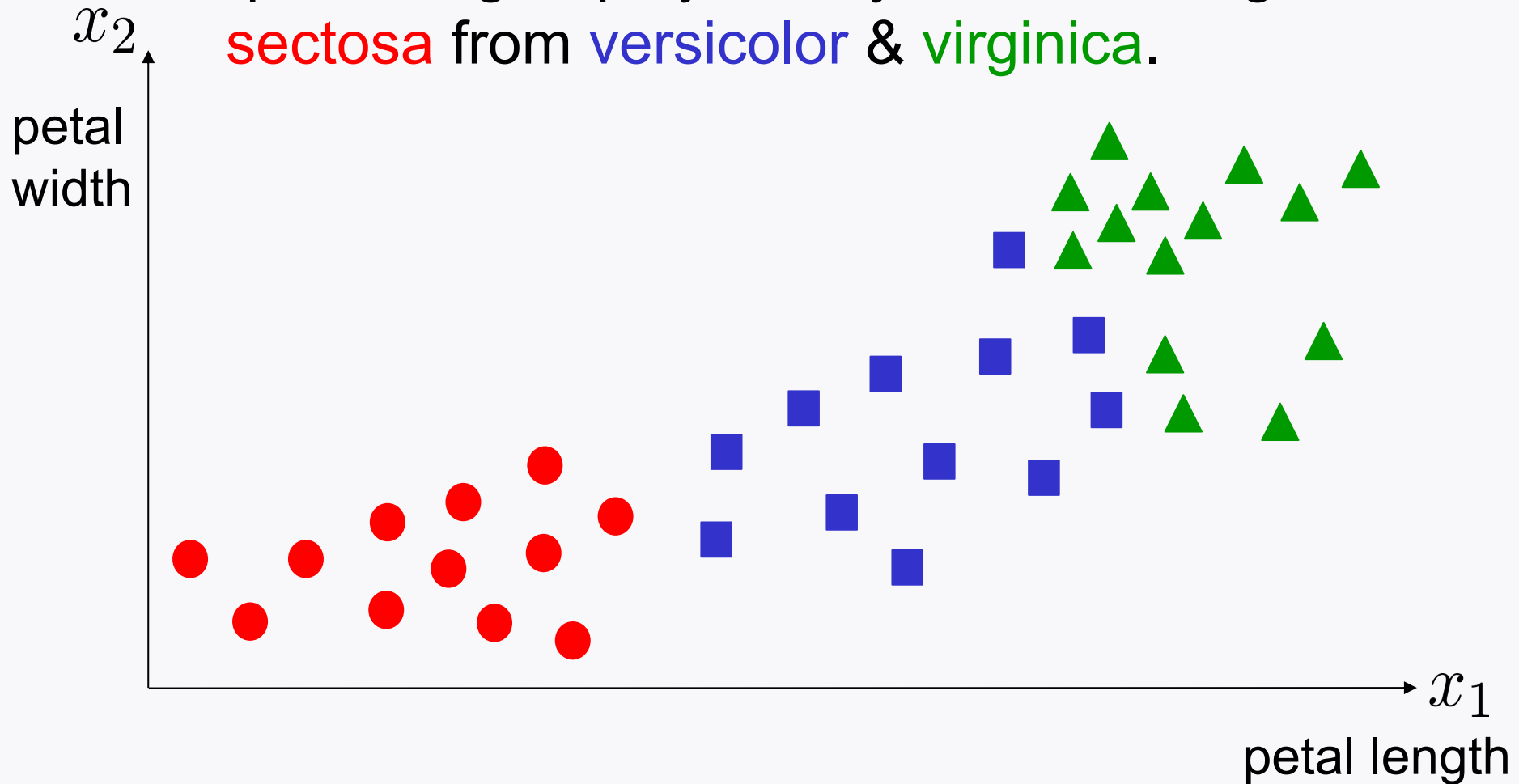Class:     <span style="color:red">setosa</span>         <span style="color:blue">versicolor</span>     <span style="color:green">virginica</span>
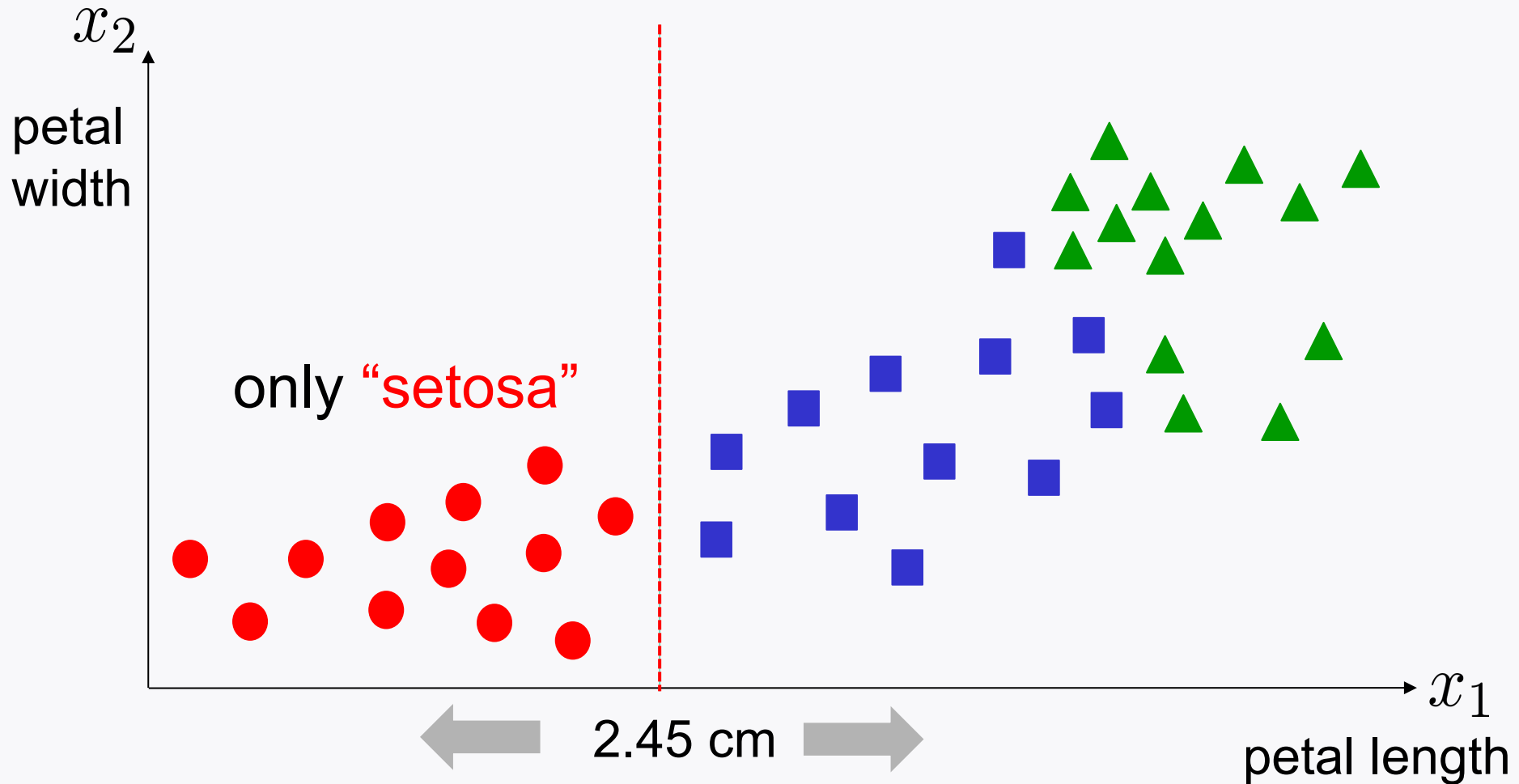


Features:     $x_1$ :   petal length

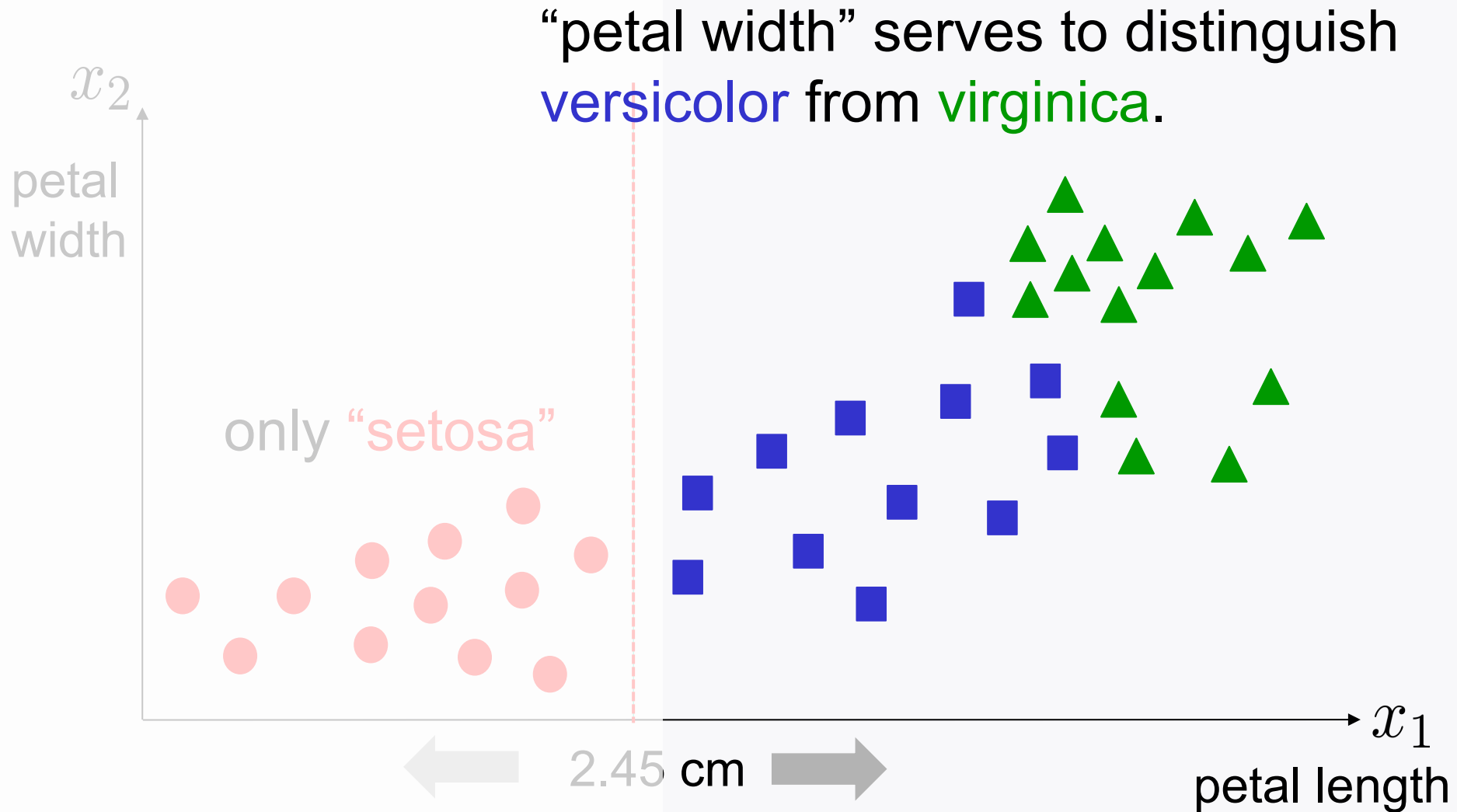                       $x_2$ :   petal width

# Data distribution

# Observation

"petal length" plays a key role to distinguish sectosa from versicolor & virginica.
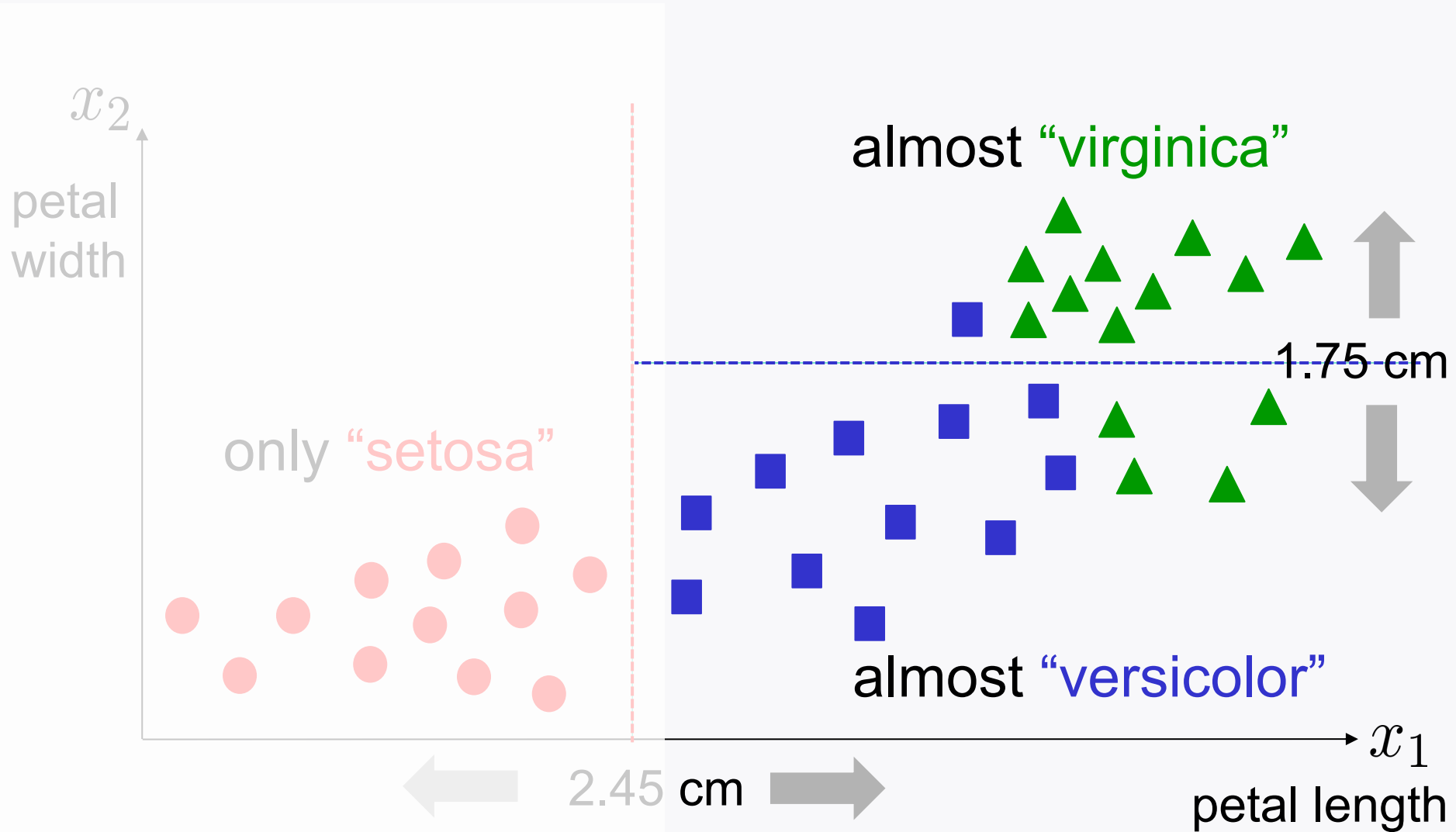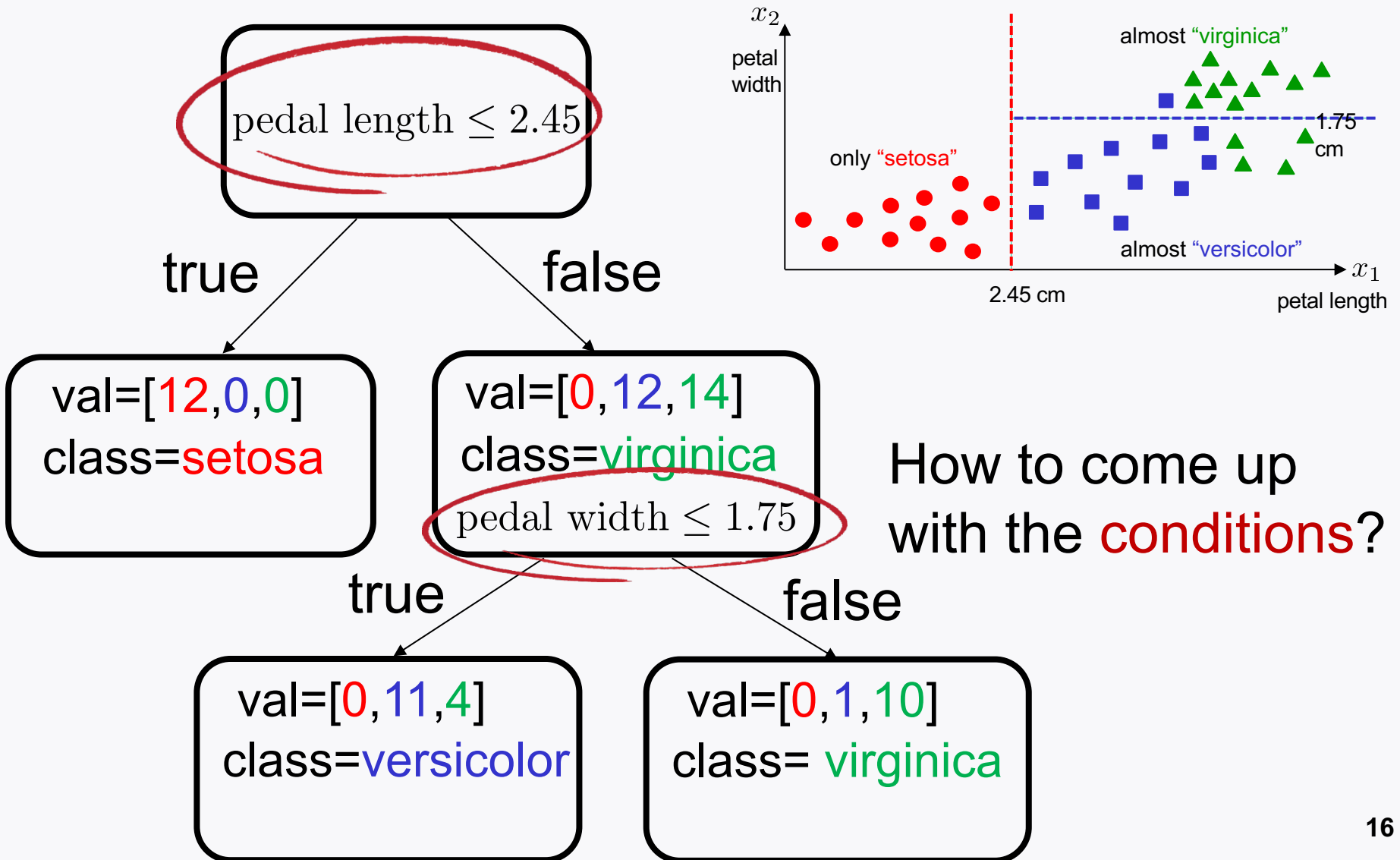
# A natural attempt for classification



only "setosa"

$x_2$
petal width

$x_1$
petal length

2.45 cm

# Another observation

"petal width" serves to distinguish versicolor from virginica.

$x_2$

petal width

only "setosa"

$x_1$

petal length

2.45 cm

# A follow-up natural attempt

# DT



pedal length $\leq 2.45$

true      false

val=[12,0,0]
class=setosa

val=[0,12,14]
class=virginica

pedal width $\leq 1.75$

true      false

val=[0,11,4]
class=versicolor

val=[0,1,10]
class= virginica

How to come up
with the conditions?

$x_2$
petal width
almost "virginica"
1.75 cm
only "setosa"
almost "versicolor"
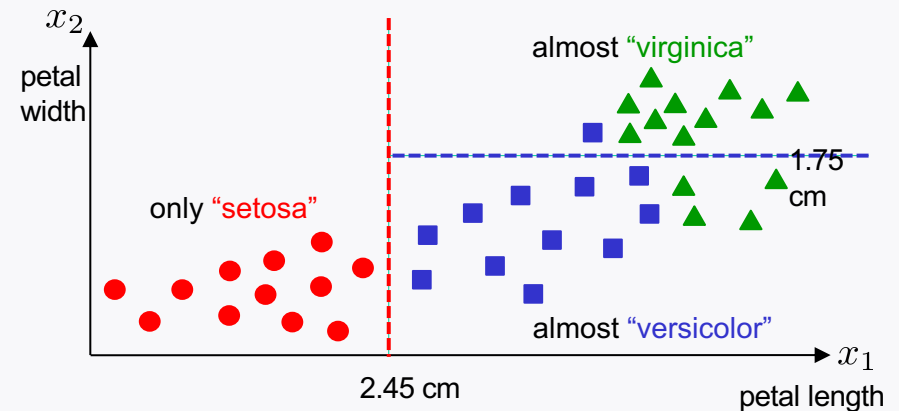$x_1$
petal length
2.45 cm

# CART (Classification And Regression Tree) algorithm

$k$ : feature index

$t_k$ : threshold



**Step 1:** Find $(k, t_k)$ such that $J(k, t_k)$ is minimized.

$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

smaller → more pure

impurity of the left split: Gini index (0~1)

$$G_{\text{left}} := 1 - \sum_{c=1}^{3} r_{\text{left},c}^2 = 1 - (1^2 + 0^2 + 0^2) = 0$$

$$G_{\text{right}} = 1 - \left(0^2 + \left(\frac{12}{26}\right)^2 + \left(\frac{14}{26}\right)^2\right) = 0.497$$

**17**

# CART (Classification And Regression Tree) algorithm

**Step 1:** Find $(k, t_k)$ such that $J(k, t_k)$ is minimized.

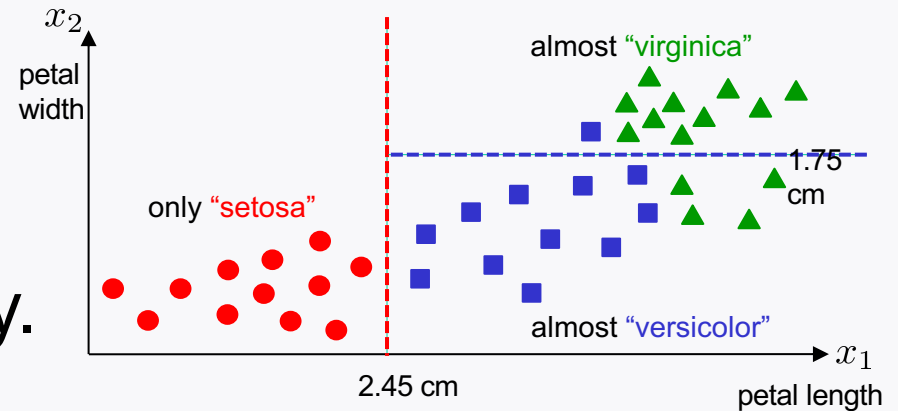$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

**Step 2:** Repeat Step 1 for each split:

$$G_{\text{left}} \begin{cases} G_{\text{left,left}} \\ G_{\text{left,right}} \end{cases} \qquad G_{\text{right}} \begin{cases} G_{\text{right,left}} \\ G_{\text{right,right}} \end{cases}$$

Stopping criteria?

# Stopping criteria

1. Cannot find a split that further reduces impurity.

OR

2. Reach "**max_depth**".

hyperparameter

max_depth=2 (in the example)

$x_2$

petal width

only "setosa"

almost "virginica"

almost "versicolor"

1.75 cm

2.45 cm

$x_1$

petal length

# Hyperparameters

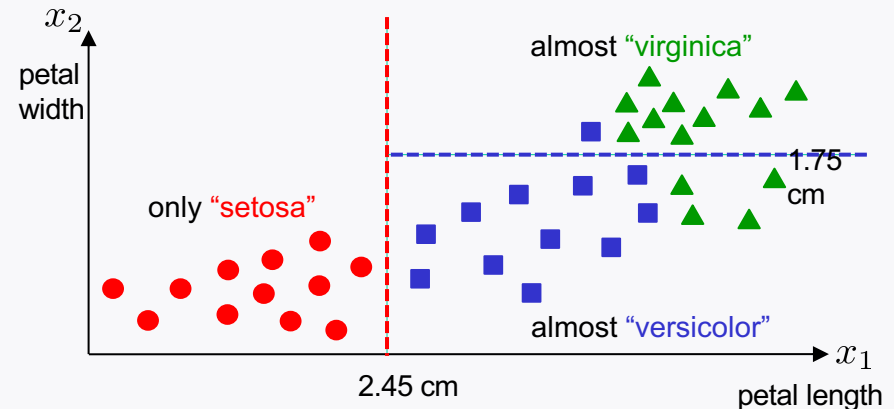1. "**max_depth**"

2. "**min_samples_split**"

   Min # of samples a node must have prior to splitting.

3. "**min_samples_leaf**"

   Min # of samples a leaf must have.

4. "**max_leaf_nodes**"

   Max # of leaf nodes



almost "virginica"

only "setosa"

almost "versicolor"

petal width

$x_2$

1.75 cm

2.45 cm

$x_1$

petal length

# Hyperparameters vs. regularization

1. "**max_depth**"    →    More regularized

2. "**min_samples_split**" ↑    More regularized

   Min # of samples a node must have prior to splitting.

3. "**min_samples_leaf**" ↑    More regularized

   Min # of samples a leaf must have.
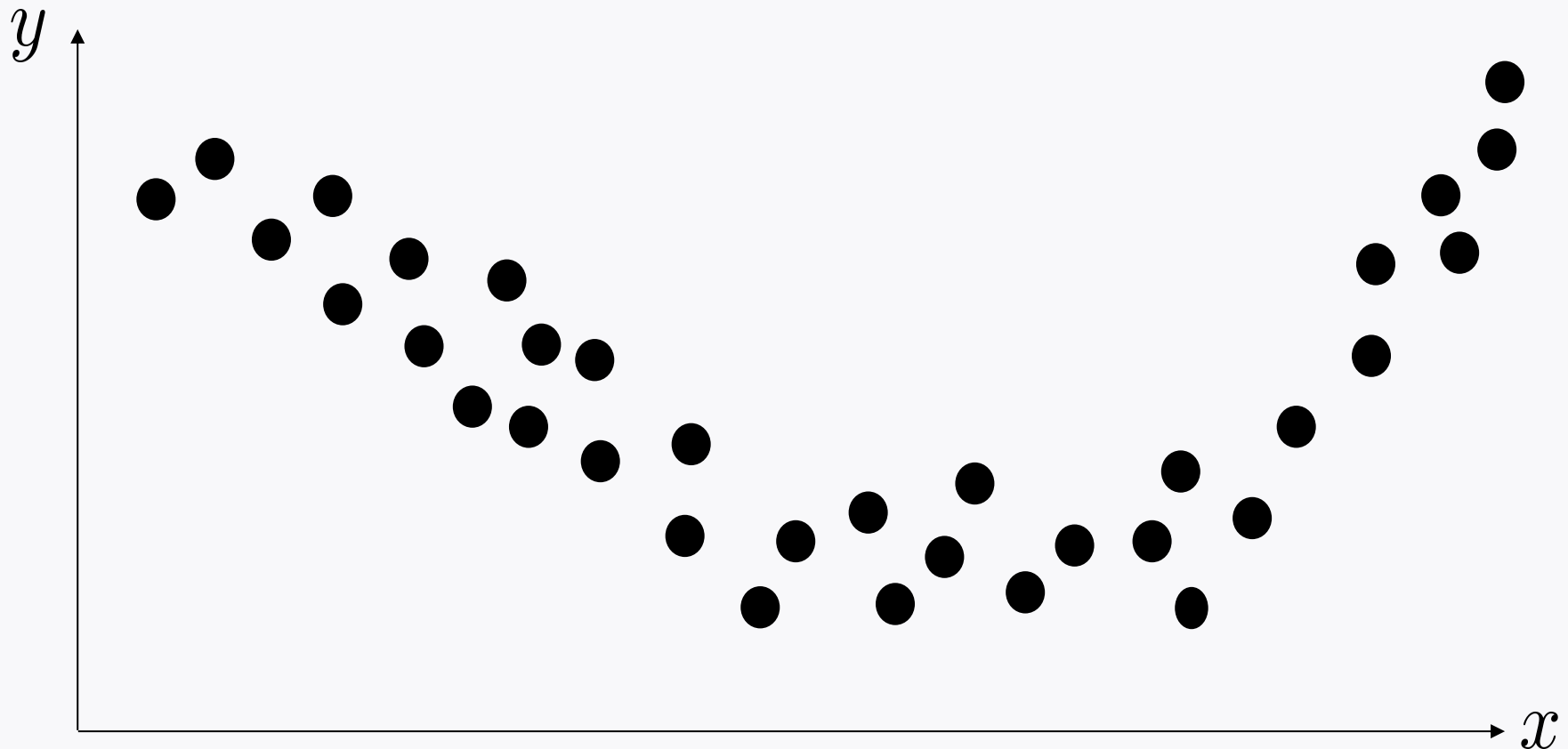
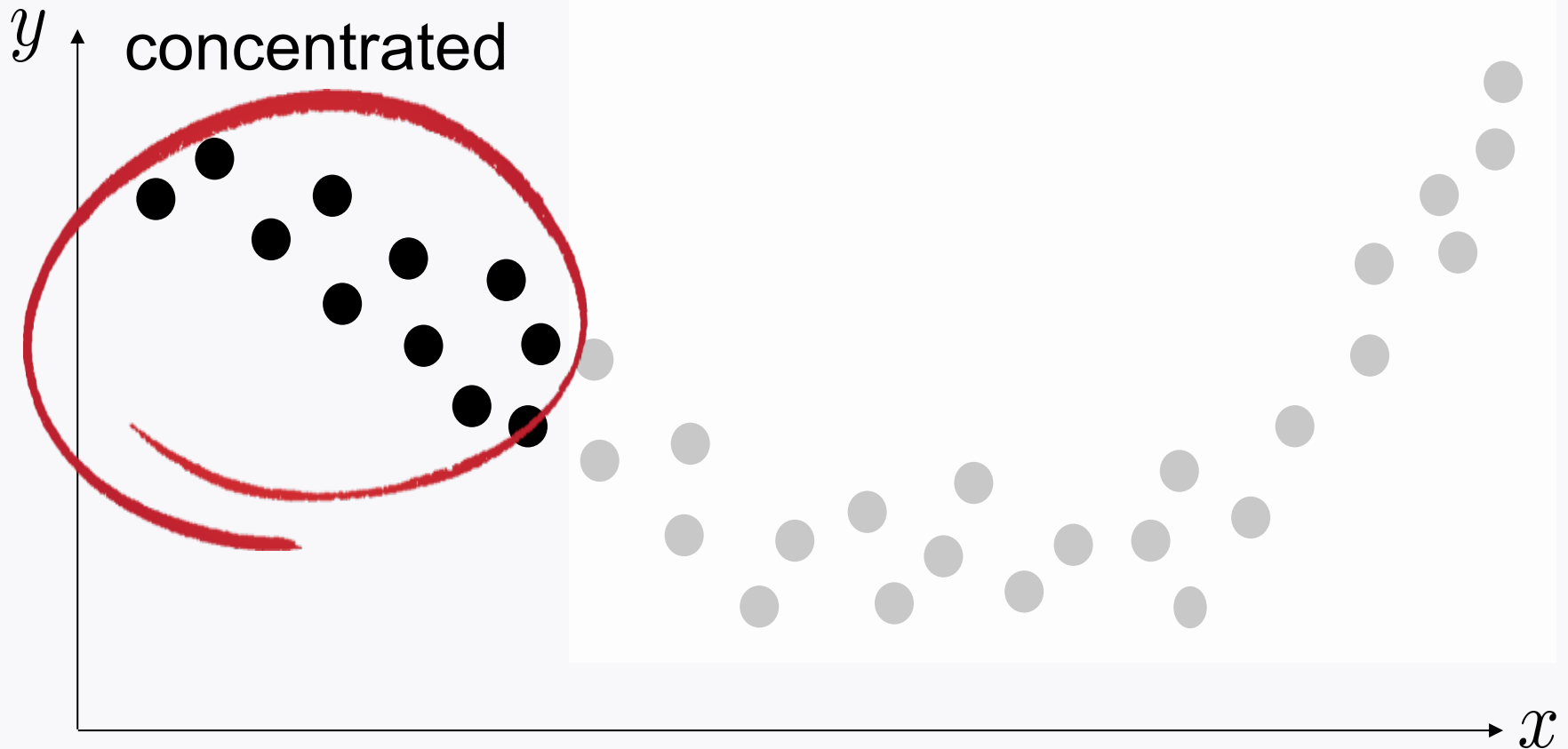4. "**max_leaf_nodes**"    More regularized

   Max # of leaf nodes

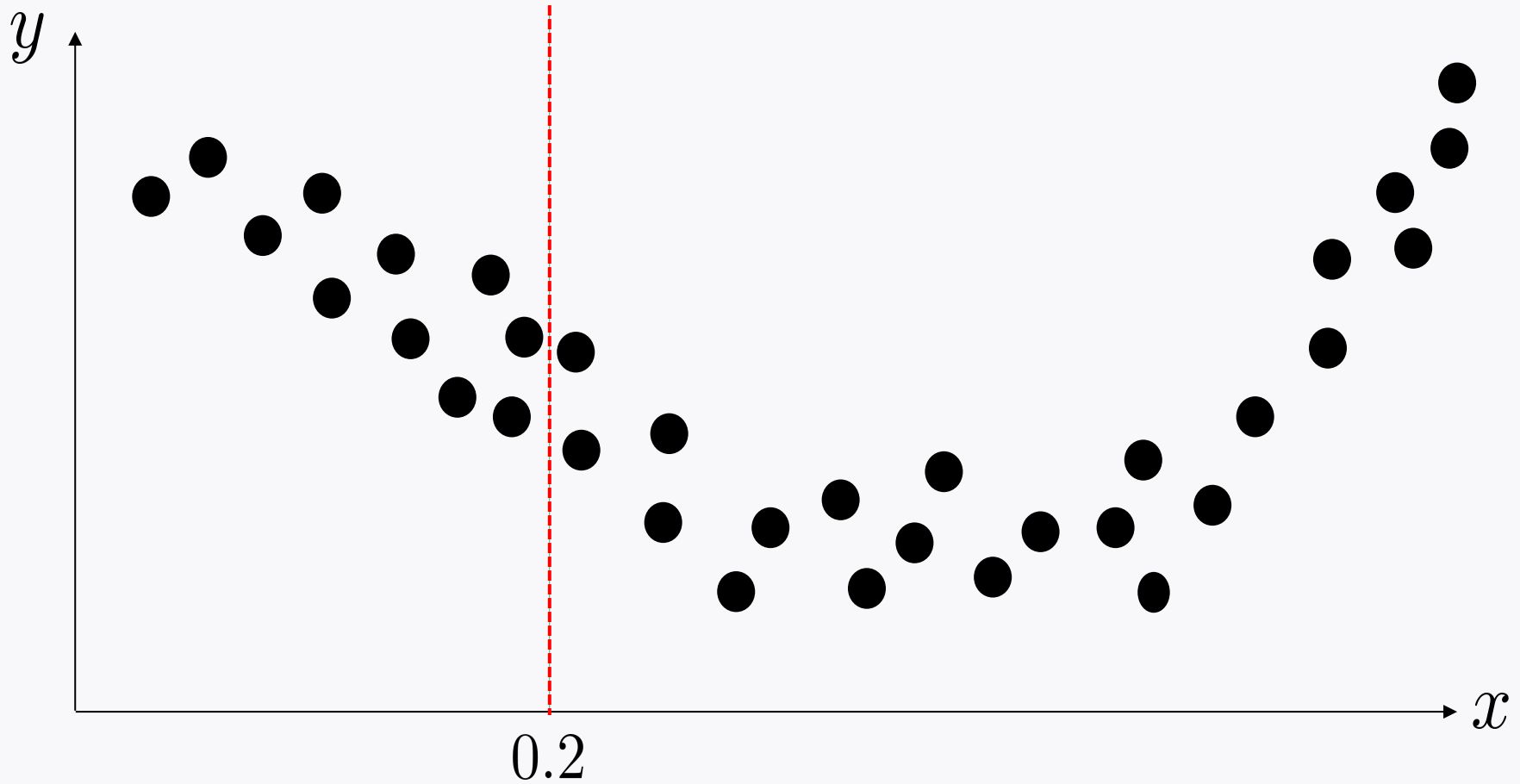# DTs for regression

# A motivating example

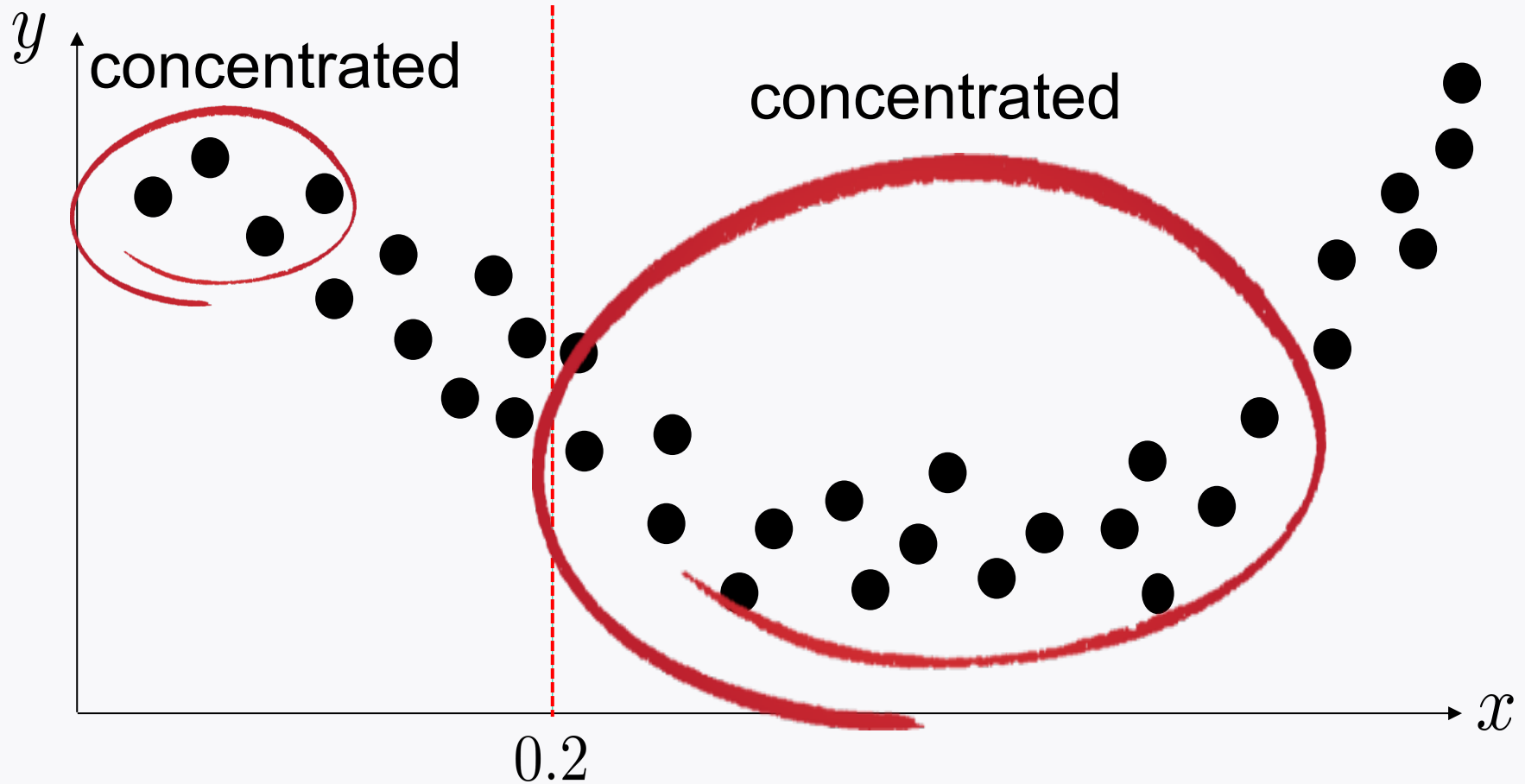$x \in \mathbf{R} \quad y \in \mathbf{R}$
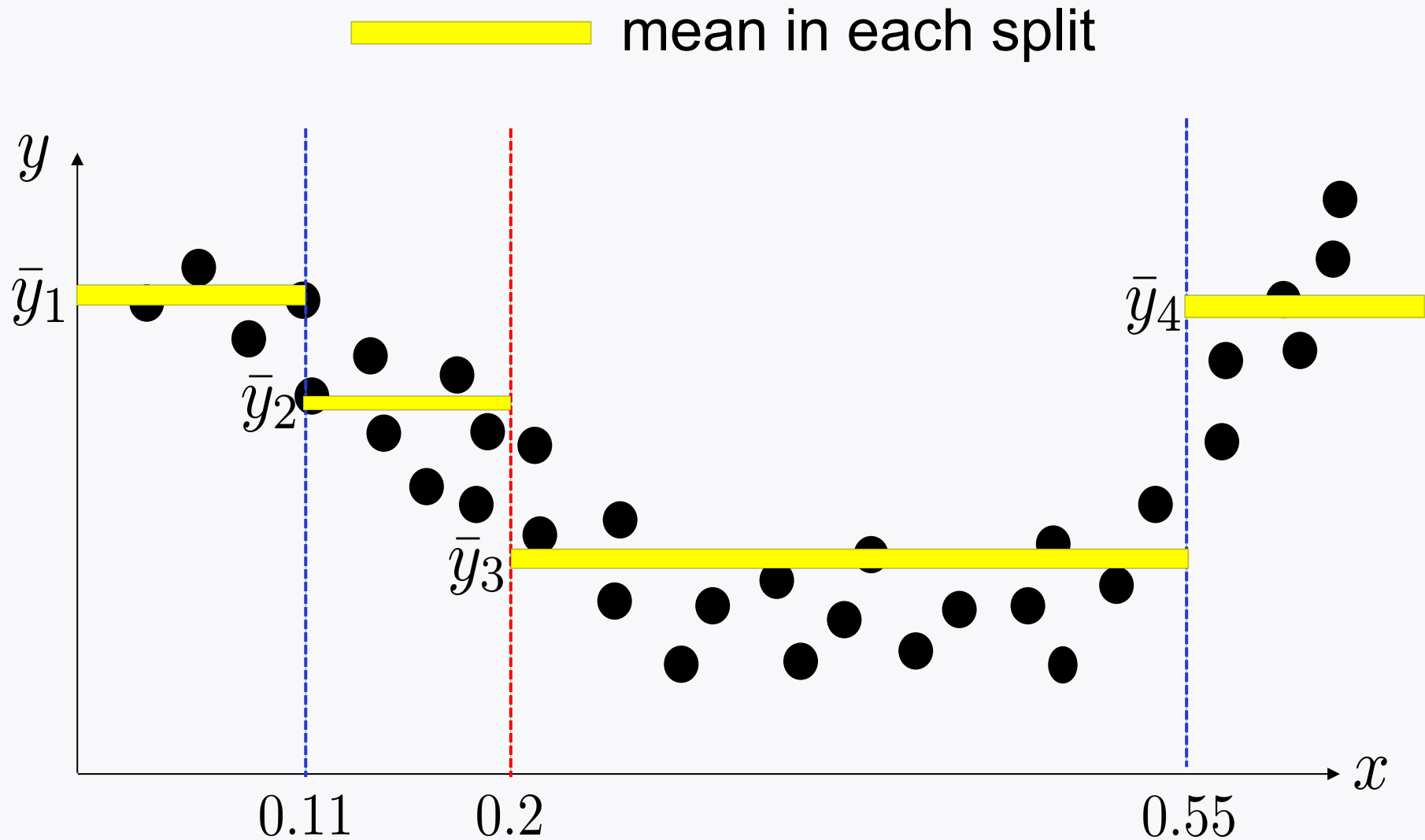
# Observation

$y$ concentrated

$x$

# A natural attempt for separation

# Observation in each split

# A follow-up natural attempt

# Decision tree

$$x \leq 0.2$$

true          false

$$x \leq 0.11$$          $$x \leq 0.55$$

true          false          true          false

$$\text{val} = \bar{y}_1$$          $$\text{val} = \bar{y}_2$$          $$\text{val} = \bar{y}_3$$          $$\text{val} = \bar{y}_4$$

$y$

$\bar{y}_1$

$\bar{y}_2$

$\bar{y}_3$

$\bar{y}_4$

$x$

0.11          0.2          0.55

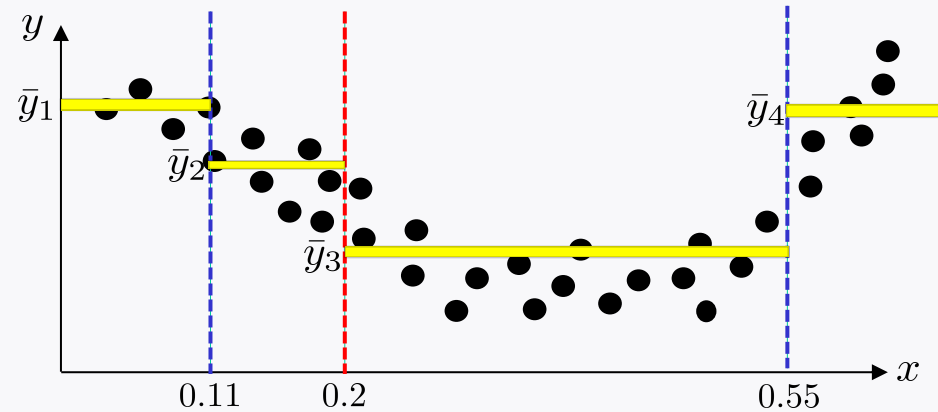How to come up
with the conditions?

# CART algorithm

$k$ : feature index

$t_k$ : threshold



**Step 1:** Find $(k, t_k)$ such that $J(k, t_k)$ is minimized.

$$J(k, t_k) = \frac{m_{\text{left}}}{m} \text{MSE}_{\text{left}} + \frac{m_{\text{right}}}{m} \text{MSE}_{\text{right}}$$

$$\text{MSE}_{\text{left}} := \sum_{i \in \text{left}} \left( y^{(i)} - \bar{y}_{\text{left}} \right)^2 \qquad \bar{y}_{\text{left}} = \frac{1}{m_{\text{left}}} \sum_{i \in \text{left}} y^{(i)}$$

# CART algorithm

**Step 1:** Find $(k, t_k)$ such that $J(k, t_k)$ is minimized.

$$J(k, t_k) = \frac{m_{\text{left}}}{m} \text{MSE}_{\text{left}} + \frac{m_{\text{right}}}{m} \text{MSE}_{\text{right}}$$
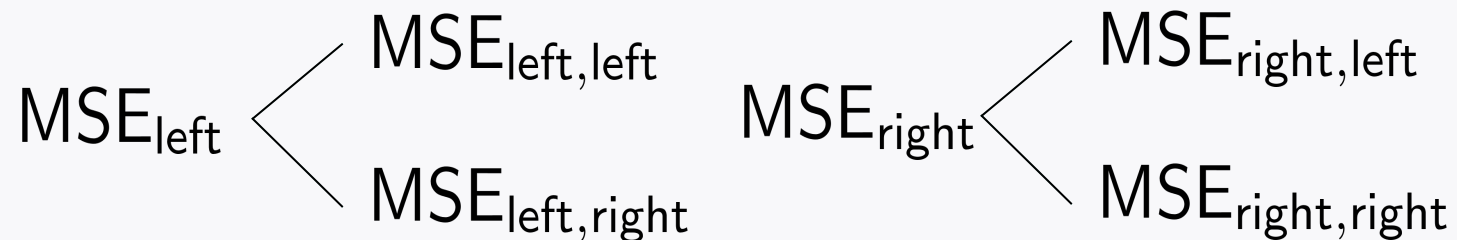
**Step 2:** Repeat Step 1 for each split:

$$\text{MSE}_{\text{left}} \begin{cases} \text{MSE}_{\text{left,left}} \\ \text{MSE}_{\text{left,right}} \end{cases} \qquad \text{MSE}_{\text{right}} \begin{cases} \text{MSE}_{\text{right,left}} \\ \text{MSE}_{\text{right,right}} \end{cases}$$

Stopping criteria & hyperparameters are the same as those of classification.

# Look ahead

1. Investigate a challenge that arises in DTs.

2. Explore a way to address the challenge:

**Ensemble learning**