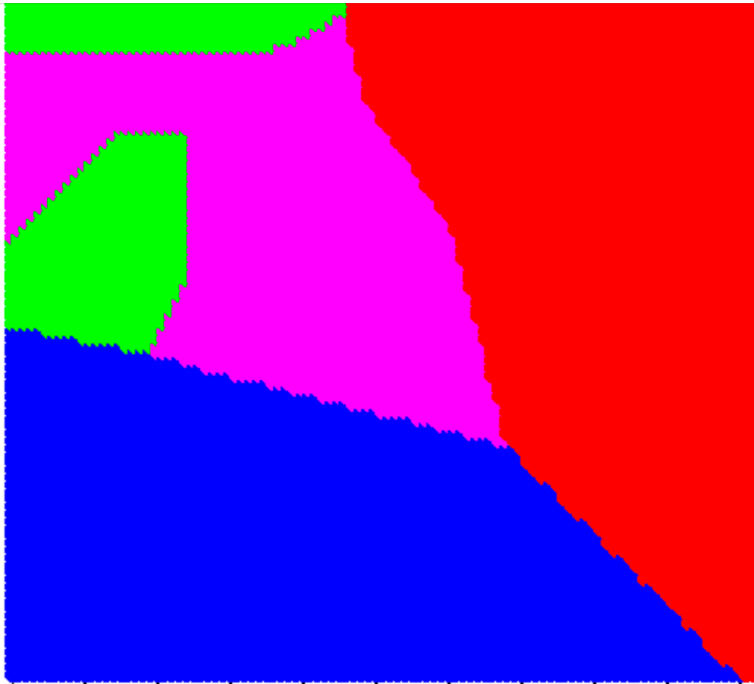# Decision Trees

# Announcements

HW6 and P6 will be released soon
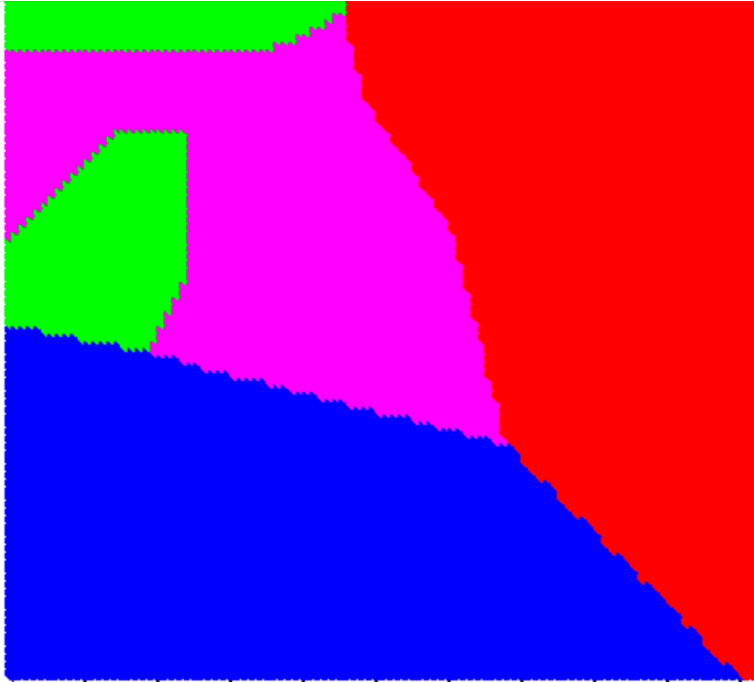
# Recap on the K-NN algorithm

K-NN can have complicated nonlinear decision boundaries



[1-NN decision boundary in prelim]

# Recap on the K-NN algorithm
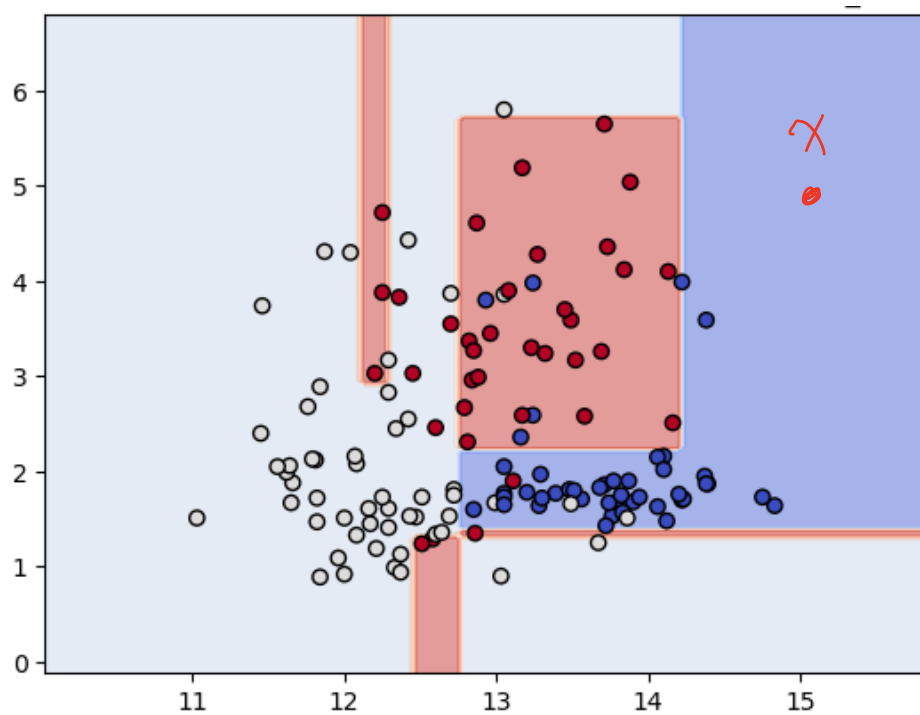
K-NN can have complicated nonlinear decision boundaries



k-NN is expensive in computation and memory

[1-NN decision boundary in prelim]

# Objective today

Decision tree — more efficient algorithm that
(1) splits space into regions with the same label, (2) is very fast in test time

# Objective today

Decision tree — more efficient algorithm that
(1) splits space into regions with the same label, (2) is very fast in test time

# Outline of Today

1. Decision tree in classification

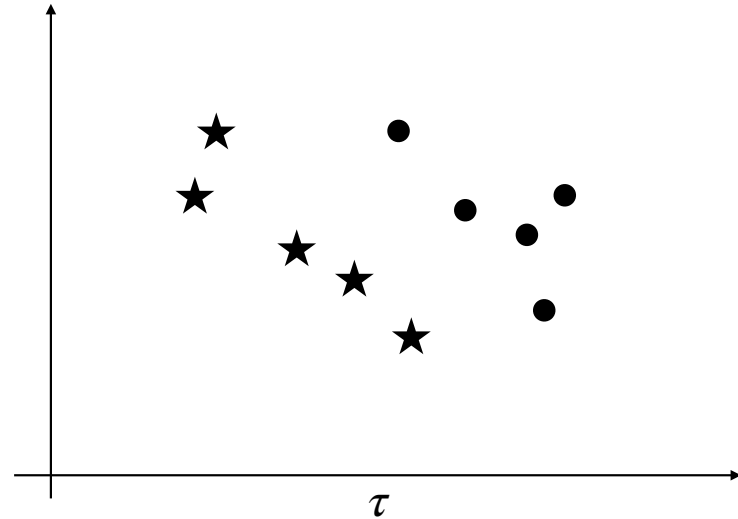2. Decision tree in regression

3. Demos of decision tree

# Overview of the Decision Tree algorithm

# How to split a tree node

Consider k-class classification, i.e., $y \in \{1, 2, \ldots, k\}$

$x \in R$

$S = \{x, y\}$



$\tau$

# How to split a tree node

Consider k-class classification, i.e., $y \in \{1, 2, \ldots, k\}$

$S = \{x, y\}$

# How to split a tree node
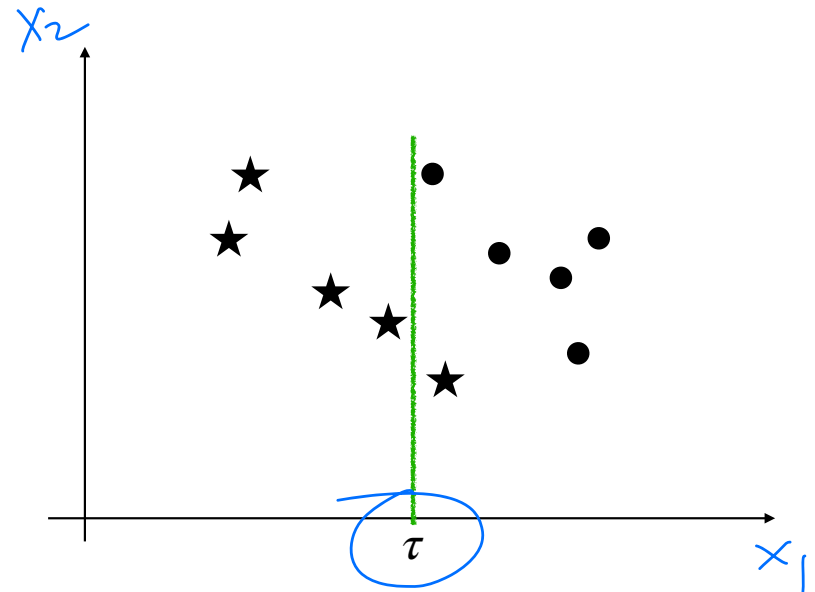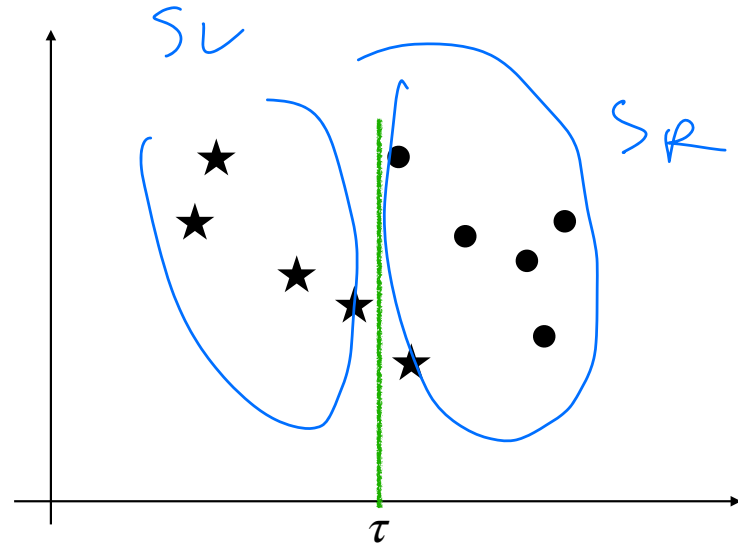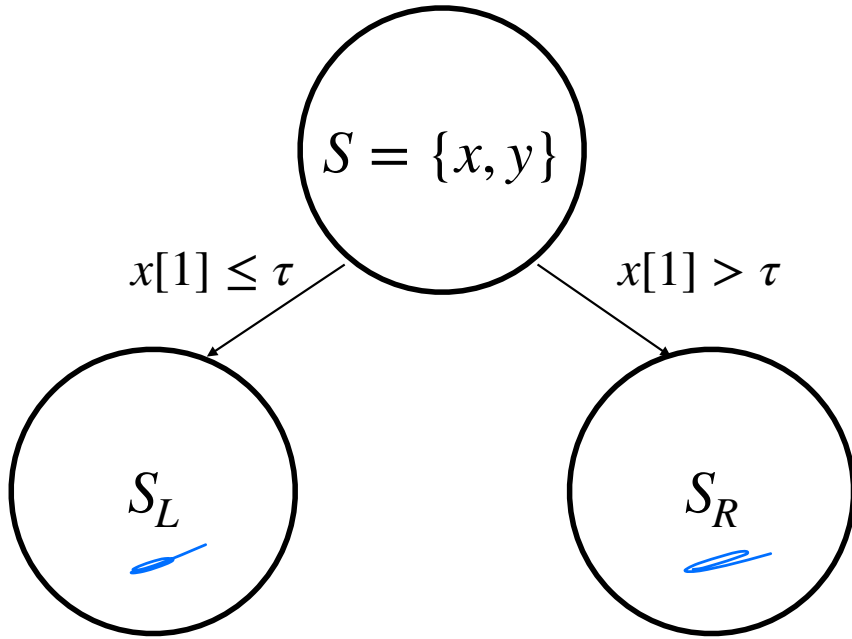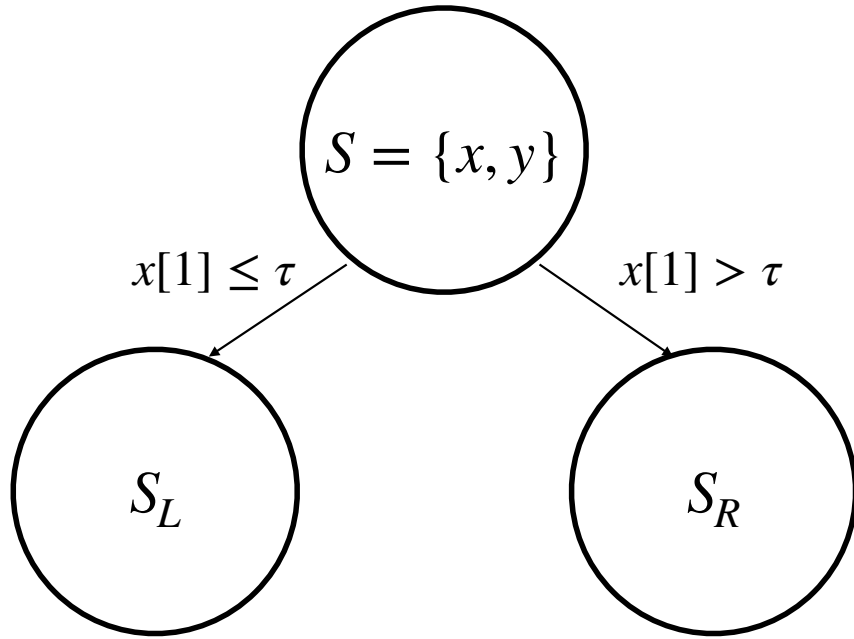
Consider k-class classification, i.e., $y \in \{1, 2, \ldots, k\}$

# How to split a tree node

Consider k-class classification, i.e., $y \in \{1, 2, \ldots, k\}$



$S = \{x, y\}$

$x[1] \leq \tau$      $x[1] > \tau$

$S_L$          $S_R$

$S_L + S_R = S, \; S_L \cap S_R = \varnothing$

$\tau$

# How to split a tree node

Consider k-class classification, i.e., $y \in \{1, 2, \ldots, k\}$



Goal: do an axis aligned split such that diversity of labels in leafs are reduced

$S = \{x, y\}$

$x[1] \leq \tau$

$x[1] > \tau$

$S_L$

$S_R$

$S_L + S_R = S, \; S_L \cap S_R = \varnothing$

# How to split a tree node

Consider k-class classification, i.e., $y \in \{1,2,\ldots,k\}$



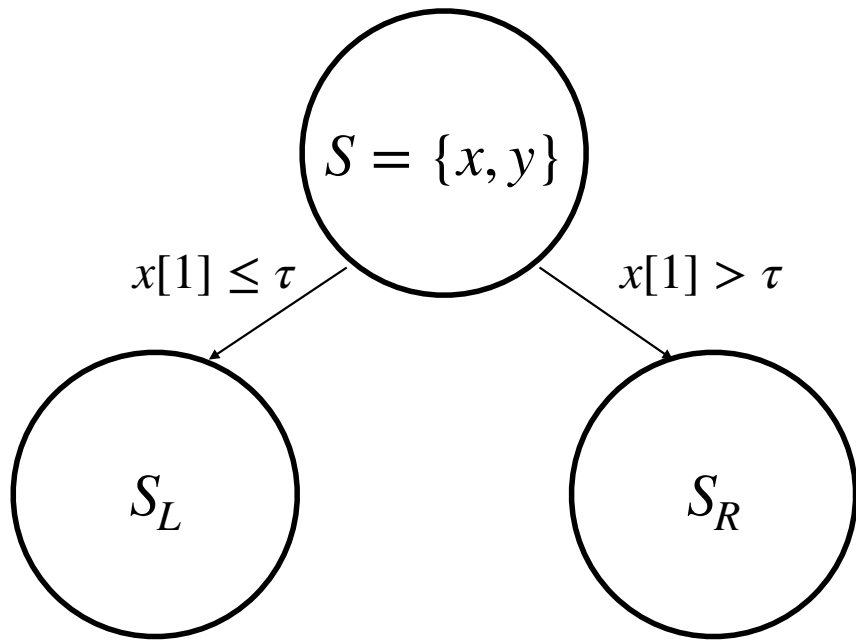Goal: do an axis aligned split such that diversity of labels in leafs are reduced

$S = \{x, y\}$

$x[1] \leq \tau$    $x[1] > \tau$

$S_L$    $S_R$

$S_L + S_R = S, \; S_L \cap S_R = \emptyset$

How to mathematically quantify "diversity"?

# Detour: Entropy

Given a set $S = \{x_i, y_i\}_{i=1}^n$, $y_i \in \{1, 2, \ldots, k\}$, measure the diversity of labels via entropy

# Detour: Entropy

Given a set $S = \{x_i, y_i\}_{i=1}^{n}$, $y_i \in \{1, 2, \ldots, k\}$, measure the diversity of labels via entropy

1. For each label i, Define $p_i = \dfrac{\text{number of label i}}{n} = \dfrac{\sum_{j=1}^{n} \mathbf{1}(y_j = i)}{n}$

# Detour: Entropy

Given a set $S = \{x_i, y_i\}_{i=1}^n$, $y_i \in \{1, 2, \ldots, k\}$, measure the diversity of labels via entropy

1. For each label i, Define $p_i = \dfrac{\text{number of label i}}{n} = \dfrac{\sum_{j=1}^n \mathbf{1}(y_j = i)}{n}$

(Probability of y being label i)

# Detour: Entropy

Given a set $S = \{x_i, y_i\}_{i=1}^n$, $y_i \in \{1,2,\ldots,k\}$, measure the diversity of labels via entropy

1. For each label i, Define $p_i = \dfrac{\text{number of label i}}{n} = \dfrac{\sum_{j=1}^n \mathbf{1}(y_j = i)}{n}$

<span style="color:green">(Probability of y being label i)</span>

2. Entropy: $H(S) = \sum_{i=1}^k -p_i \ln(p_i)$

# Detour: Entropy

Given a set $S = \{x_i, y_i\}_{i=1}^n$, $y_i \in \{1, 2, \ldots, k\}$, measure the diversity of labels via entropy

1. For each label i, Define $p_i = \dfrac{\text{number of label i}}{n} = \dfrac{\sum_{j=1}^n \mathbf{1}(y_j = i)}{n}$

<span style="color:green">(Probability of y being label i)</span>

2. Entropy: $H(S) = \sum_{i=1}^k - p_i \ln(p_i)$

<span style="color:red">High entropy means "diverse, chaos, uncertain"</span>

# Entropy

Consider a Bernoulli distribution

$$P(y = 1) = p, P(y = 0) = 1 - p$$

**Entropy** $H(y)$:

# Entropy

Consider a Bernoulli distribution

$$P(y = 1) = p, \; P(y = 0) = 1 - p$$

**Entropy $H(y)$:**

$$-P(y = 1) \cdot \ln P(y = 1) - P(y = 0) \cdot \ln P(y = 0)$$

$p$

$1-p$

# Entropy

Consider a Bernoulli distribution

$$P(y = 1) = p, \ P(y = 0) = 1 - p$$

**Entropy $H(y)$:**

$$-P(y = 1) \cdot \ln P(y = 1) - P(y = 0) \cdot \ln P(y = 0)$$

$$= -p \ln p - (1 - p)\ln(1 - p)$$

# Entropy

Consider a Bernoulli distribution

$$P(y = 1) = p, \; P(y = 0) = 1 - p$$

**Entropy $H(y)$:**

$$-P(y = 1) \cdot \ln P(y = 1) - P(y = 0) \cdot \ln P(y = 0)$$
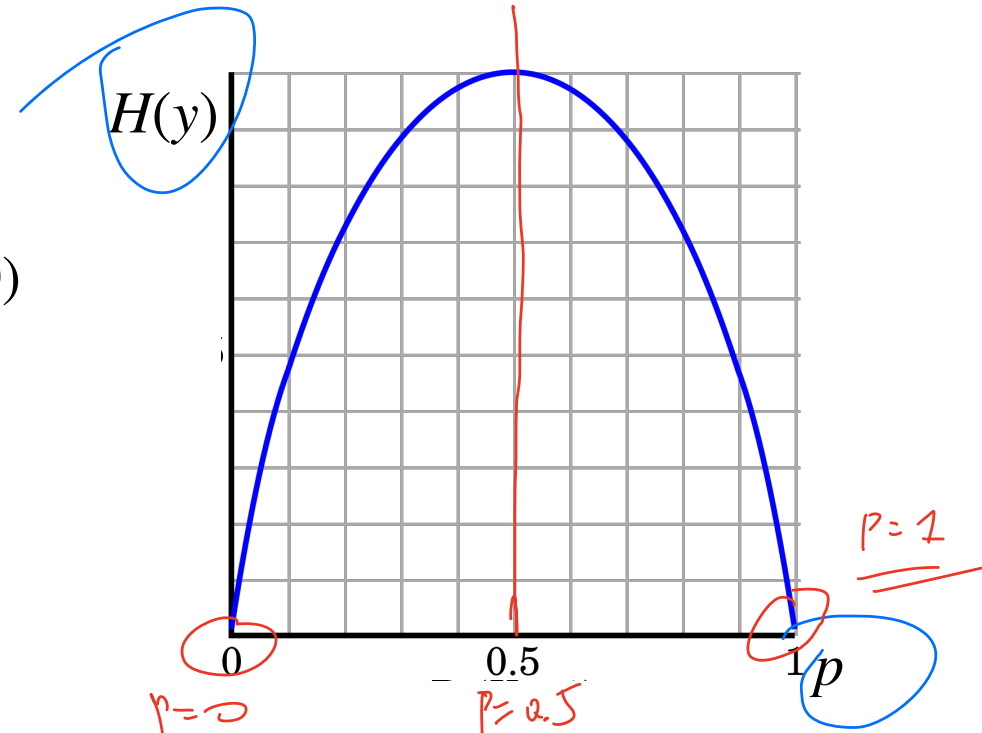
$$= - p \ln p - (1 - p)\ln(1 - p)$$

# Entropy

Consider a categorical distribution

$$y \in \{1, 2, \ldots, k\}, \ P(y = i) = p_i \geq 0, \ \sum_{i=1}^{k} p_i = 1$$

$$\sum_{i=1}^{k} -\left(p_i \ln p_i\right)$$

Q: when is entropy maximized?

$$p_1 = p_2 \ldots = p_k = \frac{1}{k}$$

# Back to tree node split…

Consider a split, i.e, dim i and threshold $\tau$,



$S = \{x, y\}$

$x[i] \leq \tau$        $x[i] > \tau$

$S_L$            $S_R$

Optimization:

# Back to tree node split…

Consider a split, i.e, dim i and threshold $\tau$,

Optimization:



$S = \{x, y\}$

$x[i] \leq \tau$     $x[i] > \tau$

$S_L$       $S_R$

$$H(S) = \sum_{i=1}^{k} - p_i \ln p_i$$

$$p_i = \frac{\# \text{ of } \cancel{\text{example}} \text{ ins with label } i}{|S|}$$

# Back to tree node split...

Consider a split, i.e, dim i and threshold $\tau$,



$$H(S) = \sum_{i=1}^{k} - p_i \ln p_i$$

$$\left(\frac{|S_L|}{|S|}\right) H(S_L) + \left(\frac{|S_R|}{|S|}\right) H(S_R)$$
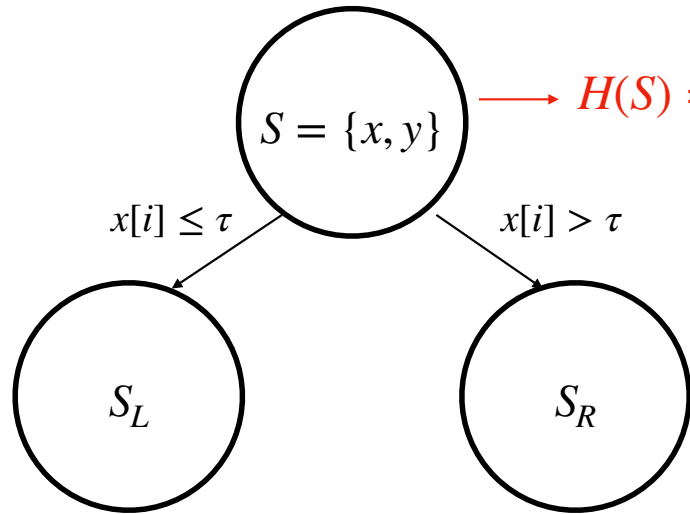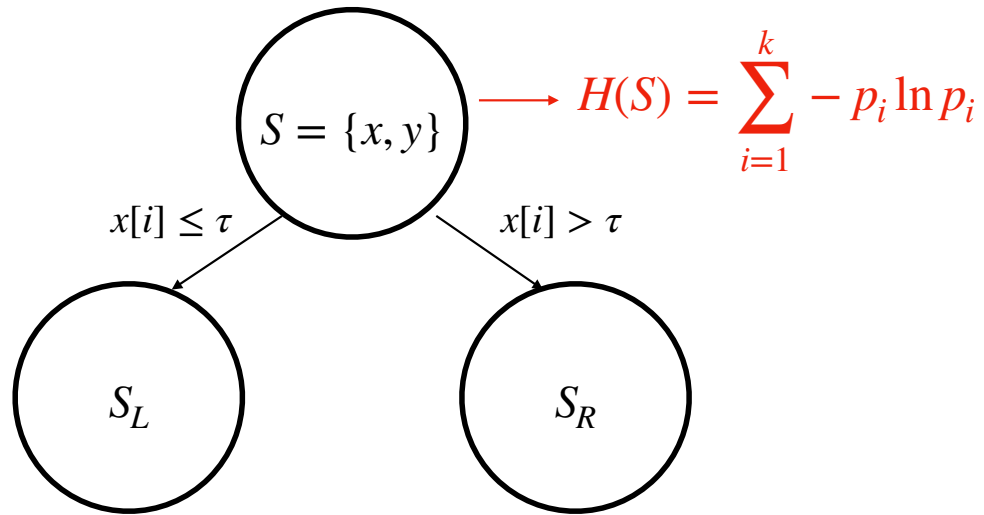
Optimization:

# Back to tree node split…

Consider a split, i.e, dim i and threshold $\tau$,



$$H(S) = \sum_{i=1}^{k} -p_i \ln p_i$$

$x[i] \le \tau$     $x[i] > \tau$

$S = \{x, y\}$

$S_L$       $S_R$
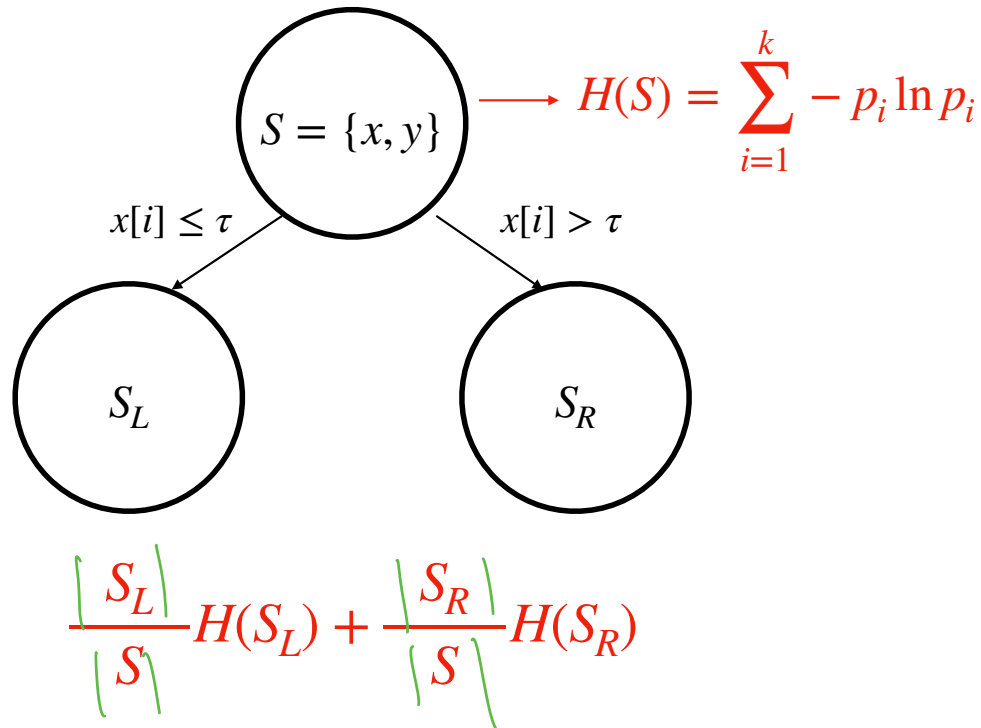
$$\frac{|S_L|}{|S|}H(S_L) + \frac{|S_R|}{|S|}H(S_R)$$

Optimization:

Find a split $(i, \tau)$ such that

$$\frac{|S_L|}{|S|}H(S_L) + \frac{|S_R|}{|S|}H(S_R)$$   is the smallest

# Back to tree node split...

Consider a split, i.e, dim i and threshold $\tau$,

Optimization: $\tau \in R$ threshold



$$H(S) = \sum_{i=1}^{k} -p_i \ln p_i$$

$S = \{x, y\}$

$x[i] \leq \tau$    $x[i] > \tau$

$S_L$    $S_R$

$$\frac{|S_L|}{|S|} H(S_L) + \frac{|S_R|}{|S|} H(S_R)$$

Find a split $(i, \tau)$ such that

$$\frac{|S_L|}{|S|} H(S_L) + \frac{|S_R|}{|S|} H(S_R)$$ is the smallest

# Back to tree node split…
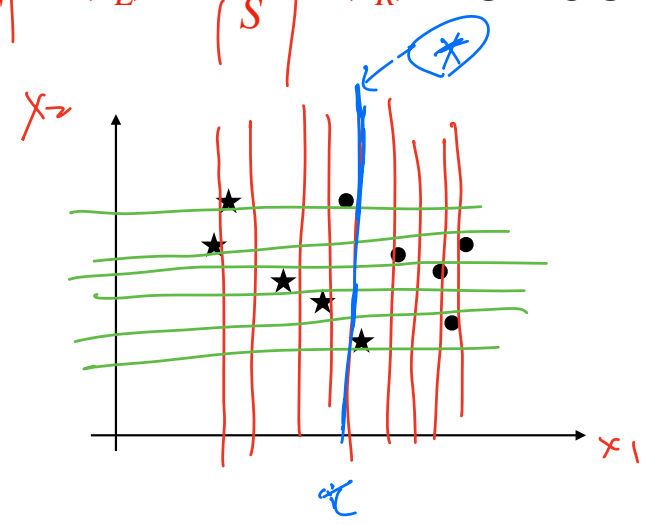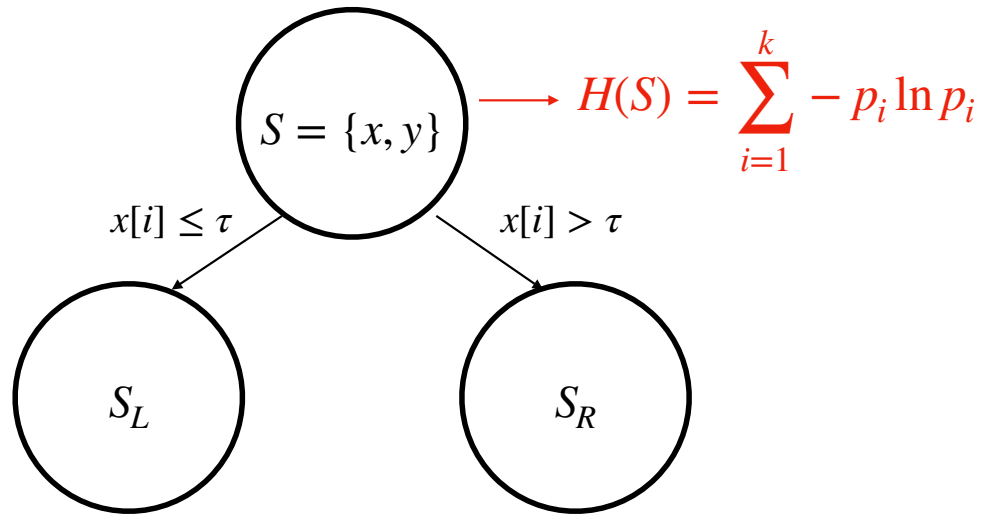
Consider a split, i.e, dim i and threshold $\tau$,



$$H(S) = \sum_{i=1}^{k} - p_i \ln p_i$$

$x[i] \leq \tau$      $x[i] > \tau$
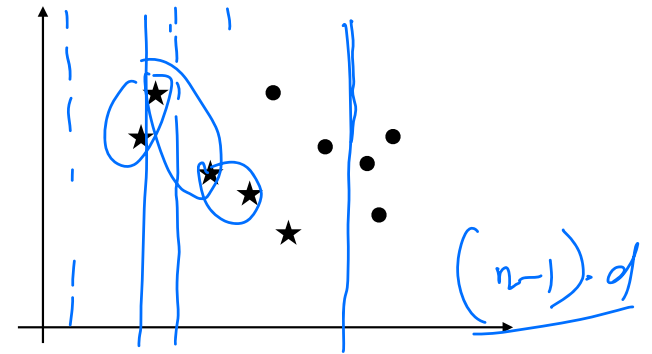
$S = \{x, y\}$

$S_L$         $S_R$

$$\frac{|S_L|}{|S|} H(S_L) + \frac{|S_R|}{|S|} H(S_R)$$

Optimization:

Find a split $(i, \tau)$ such that

$$\frac{|S_L|}{|S|} H(S_L) + \frac{|S_R|}{|S|} H(S_R) \quad \text{is the smallest}$$



$(n-1) \cdot d$

Q: how many splits we need to check?

$n$ points

# Put everything together — ID3 algorithm

Input: training set $S = \{x, y\}$

**Decision_tree($S$):**

# Put everything together — ID3 algorithm

Input: training set $S = \{x, y\}$

**Decision_tree($S$):**

- If all $y$ in $S$ are the same

# Put everything together — ID3 algorithm

Input: training set $S = \{x, y\}$

**Decision_tree($S$):**

- If all $y$ in $S$ are the same

    Done, and return this label
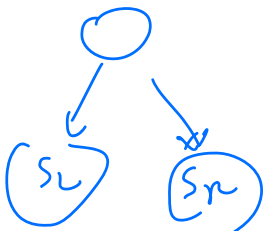
# Put everything together — ID3 algorithm

Input: training set $S = \{x, y\}$

**Decision_tree($S$):**

- If all $y$ in $S$ are the same

    Done, and return this label

- Else:

    Find a split $(i, \tau)$ that minimizes
    weighted entropy $\longrightarrow$ $\dfrac{|S_L|}{|S|} \cdot H(S_L) + \dfrac{|S_R|}{|S|} \cdot H(S_R)$

# Put everything together — ID3 algorithm

Input: training set $S = \{x, y\}$

**Decision_tree($S$):**

- If all $y$ in $S$ are the same

    Done, and return this label

- Else:

    Find a split $(i, \tau)$ that minimizes weighted entropy

    Call **Decision_tree($S_L$)** & **Decision_tree($S_R$)**

# Put everything together — ID3 algorithm

Input: training set $S = \{x, y\}$

**Decision_tree($S$):**

- If all $y$ in $S$ are the same

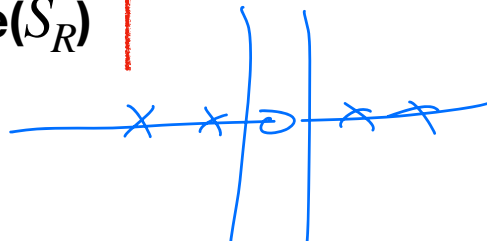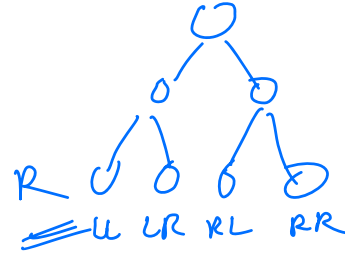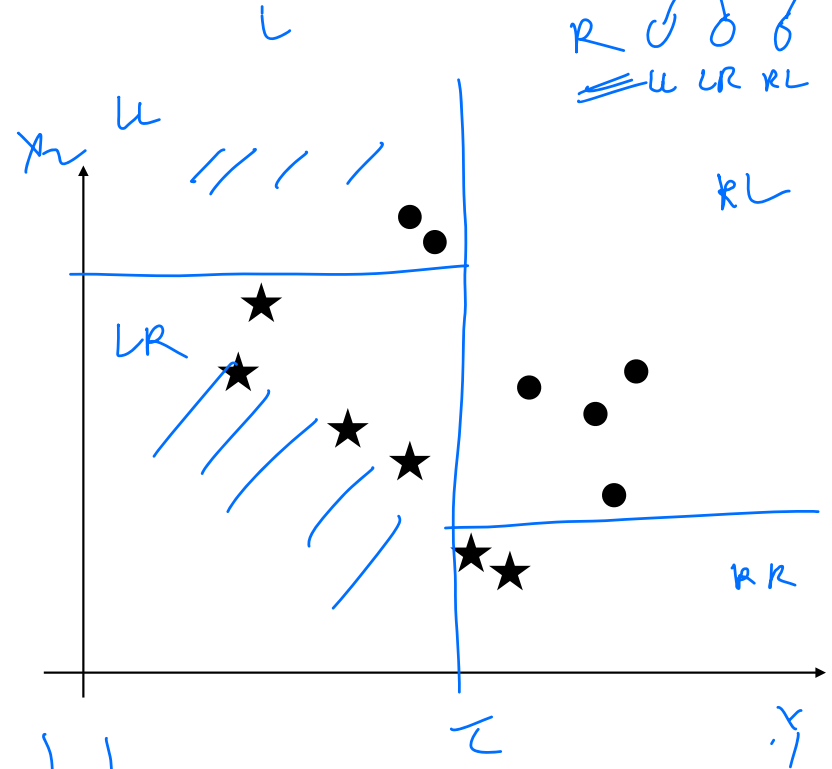  Done, and return this label

- Else:

  Find a split $(i, \tau)$ that minimizes weighted entropy

  Call **Decision_tree($S_L$)** & **Decision_tree($S_R$)**

$$x_i \quad x_j$$

$$\begin{pmatrix} x_i = x_j \\ y_i \pm y_j \end{pmatrix}$$

# Outline of Today

1. Decision tree in classification

2. Decision tree in regression

3. Demos of decision tree

# Regression

How to split the note, i.e., what is the diversity measure?

$$y \in \mathbb{R}$$

# Regression

How to split the note, i.e., what is the diversity measure?

Consider a set of training points $S = \{x_i, y_i\}_{i=1}^m, y_i \in \mathbb{R}$

# Regression

How to split the note, i.e., what is the diversity measure?

Consider a set of training points $S = \{x_i, y_i\}_{i=1}^{m}, y_i \in \mathbb{R}$

Define the sample mean $\hat{y}_S = \sum_{i=1}^{m} y_i / m$

# Regression

How to split the note, i.e., what is the diversity measure?

Consider a set of training points $S = \{x_i, y_i\}_{i=1}^m, y_i \in \mathbb{R}$

Define the sample mean $\hat{y}_S = \sum_{i=1}^m y_i/m$

Impurity: sample variance $\widehat{Var}(S) = \sum_{i=1}^m (y_i - \bar{y}_S)^2/m$

$\Longrightarrow Var(y)$ when $m \to \infty$

# Regression Tree

Regression_Tree( $S$ ):

# Regression Tree

Regression_Tree( $S$ ):

- IF $|S| \leq k$:

  Set leaf value to be $\bar{y}_S$

Threshod $K \in \mathbb{N}^+$

$(e.g \; k=3)$

# Regression Tree

Regression_Tree( $S$ ):

- IF $|S| \leq k$:

    Set leaf value to be $\bar{y}_S$

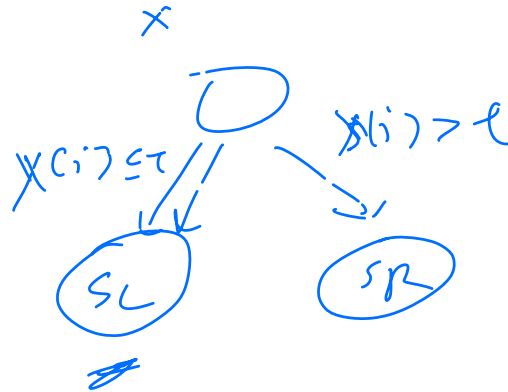- ELSE:

# Regression Tree

Regression_Tree( $S$ ):

- IF $|S| \leq k$:

  Set leaf value to be $\bar{y}_S$

- ELSE:

  For all $(i, \tau)$, find the split such that minimizes $\dfrac{|S_L|}{|S|} \widehat{Var}(S_L) + \dfrac{|S_R|}{|S|} \widehat{Var}(S_R)$

# Regression Tree

$x$

$x(i) \leq \tau$   $x(i) > \ell$

$S_L$   $S_R$

Regression_Tree( $S$ ):

- IF $\lfloor S \rfloor \leq k$:

    Set leaf value to be $\bar{y}_S$

- ELSE:

    For all $(i, \tau)$, find the split such that minimizes $\dfrac{\lfloor S_L \rfloor}{\lfloor S \rfloor} \widehat{Var}\,(S_L) + \dfrac{\lfloor S_R \rfloor}{\lfloor S \rfloor} \widehat{Var}\,(S_R)$

    Call Regression_Tree( $S_L$ ) & Regression_Tree( $S_R$ )

# Outline of Today

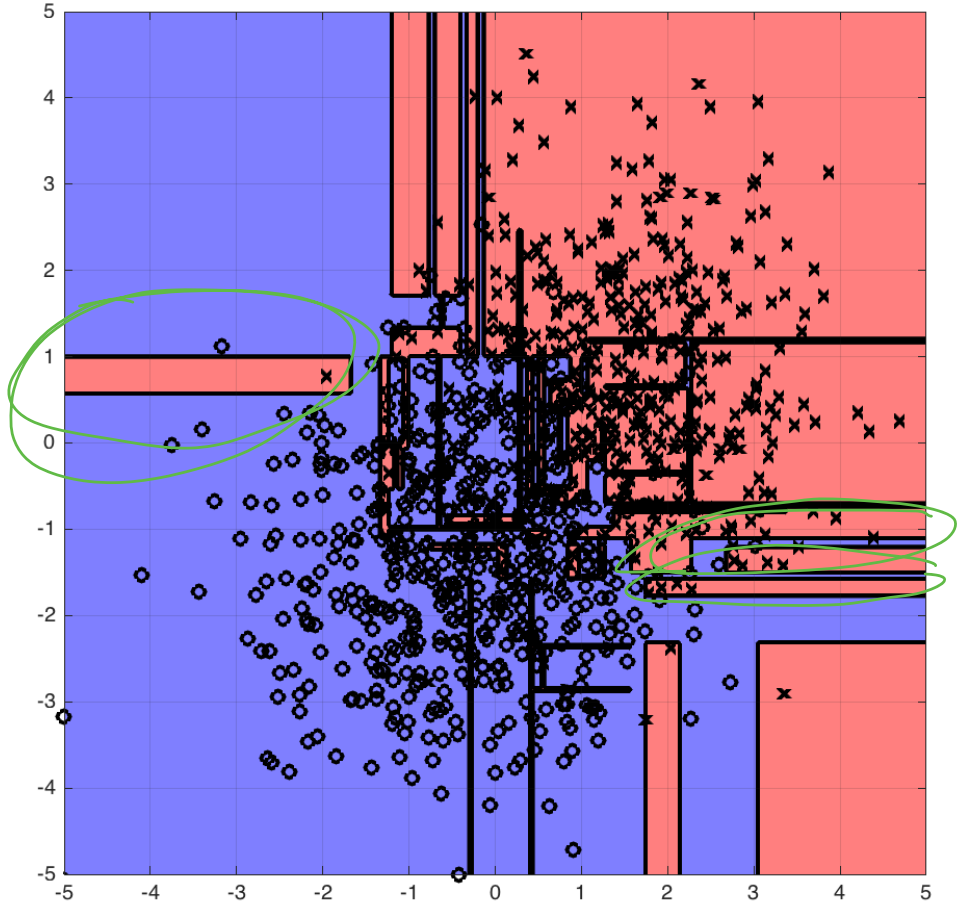1. Decision tree in classification

2. Decision tree in regression
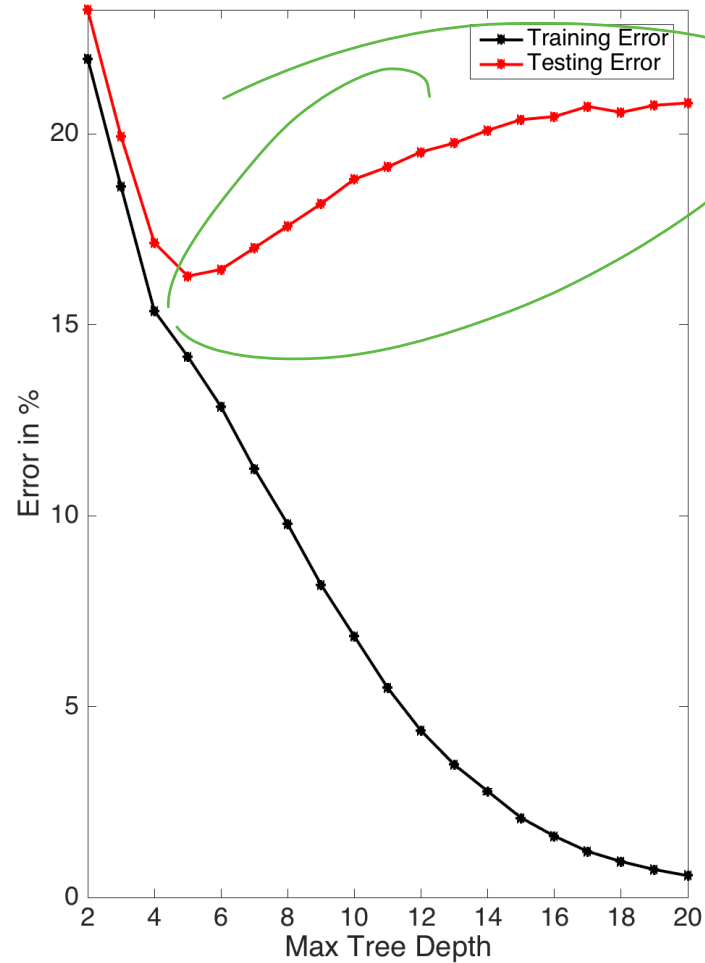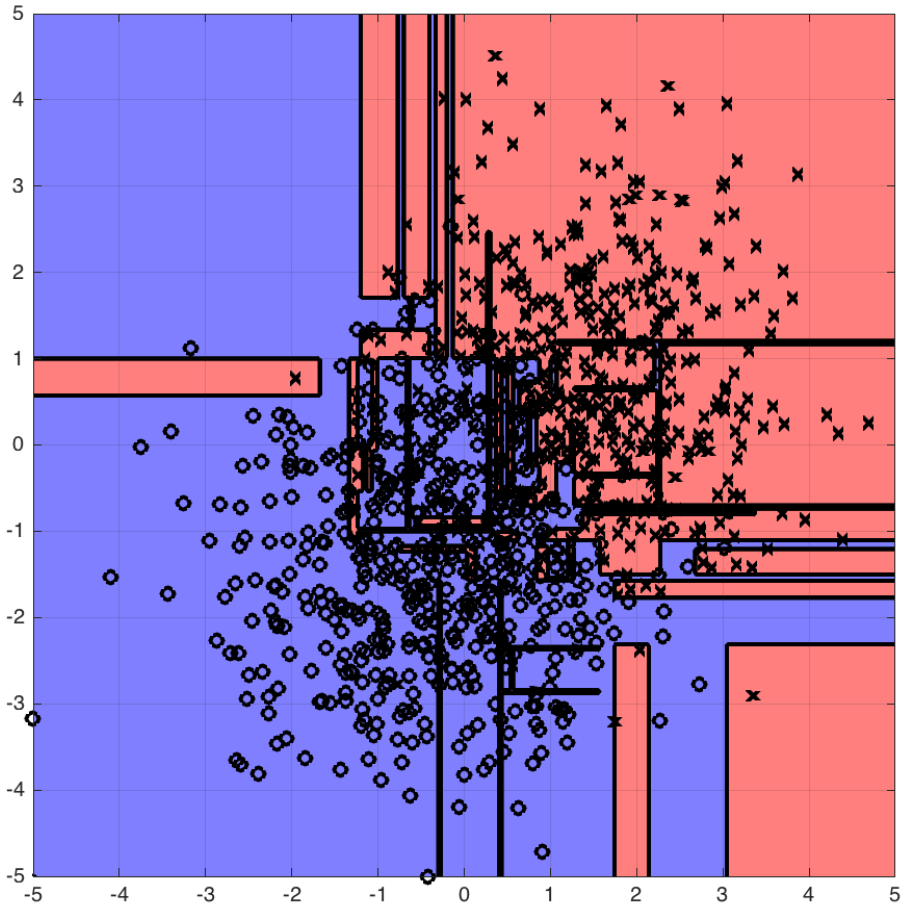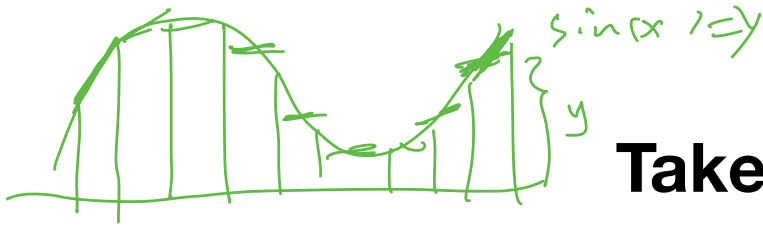
3. Demos of decision tree

# Issue of Decision Trees

Decision Tree can have high variance, i.e., overfilling!

# Issue of Decision Trees



Decision Tree can have high variance, i.e., overfilling!

# Issue of Decision Trees



Decision Tree can have high variance, i.e., overfilling!

$sin(x) = y$

$y$

# Take-home messages

1 Decision tree algorithms splits space into axis-aligned regions

Each region ideally should only contain one unique label

2: Split a node such that the entropy of labels in the leafs are minimized

$\sum_{y \in S_L} \frac{y}{|S_L|}$

$x < $

$(x < y)$

3: Can easily overfit as the depth of the tree increases
(limiting the depth of the tree is a good regularization)

$S_L$

$S_R$