上海科技大学
ShanghaiTech University

# CS182: Introduction to Machine Learning – Decision Trees

Yujiao Shi

SIST, ShanghaiTech

Spring, 2025

# Q & A:

# How do these in-class polls work?

- Scan QR Code for accessing the polls

- Answer all poll questions **during lecture for full credit** or **within 24 hours for half credit**

- You have 8 free "poll points" for the semester that will excuse you from all polls from a single lecture; you cannot use more than 3 poll points consecutively.
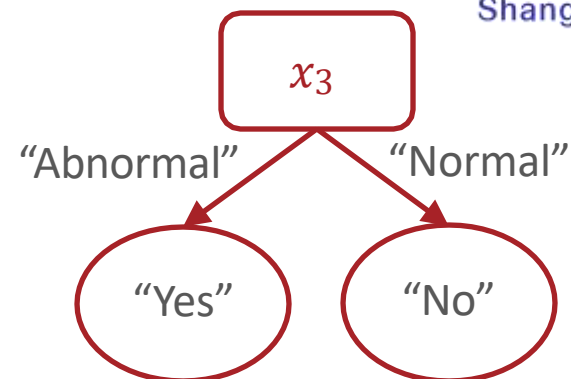
# Poll Question 1:

## Which of the following did you bring to class today? Select all that apply

A. A smartphone

B. A computer

C. A smart watch

D. No device

# Recall: Decision Stump

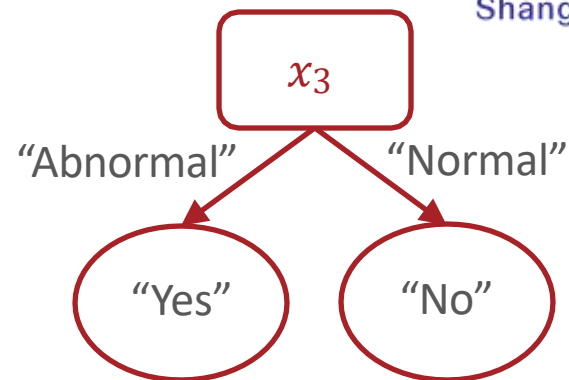| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |

$x_3$

"Abnormal"    "Normal"

"Yes"    "No"

# Recall: Decision Stump Questions

1. How can we pick which feature to split on?

2. Why stop at just one feature? **Don't!**

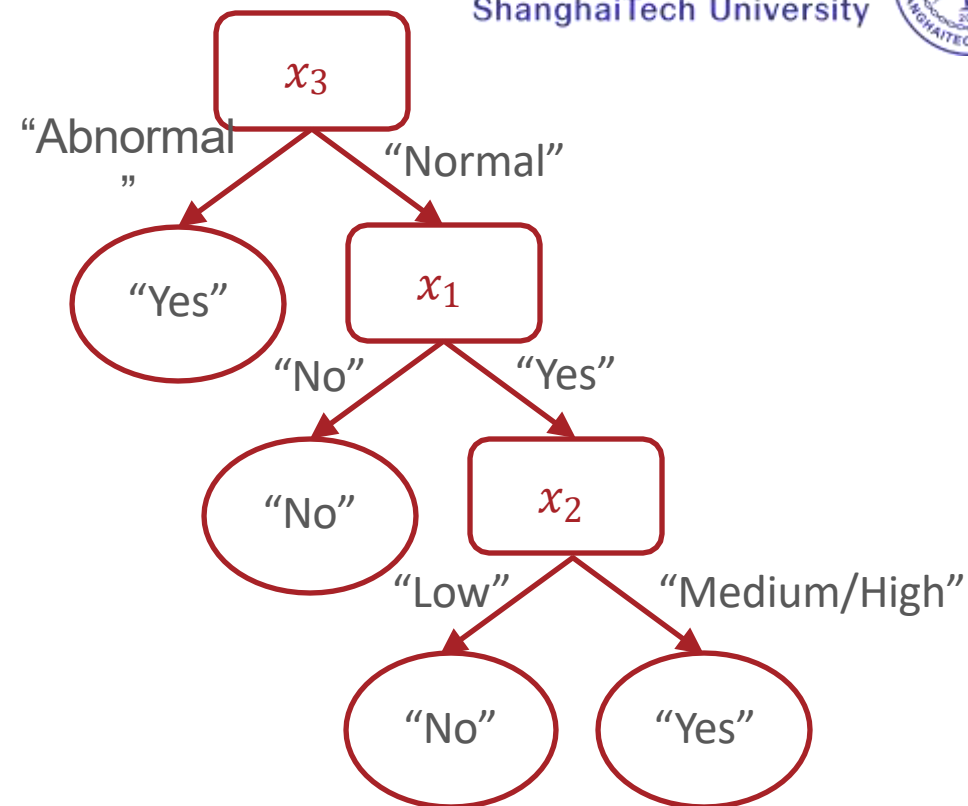   a) If we split on more than one feature, how do we decide the order to spilt on?

# From Decision Stump ...

| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |



$x_3$

"Abnormal"  "Normal"

"Yes"  "No"

上 海 科 技 大 学
ShanghaiTech University
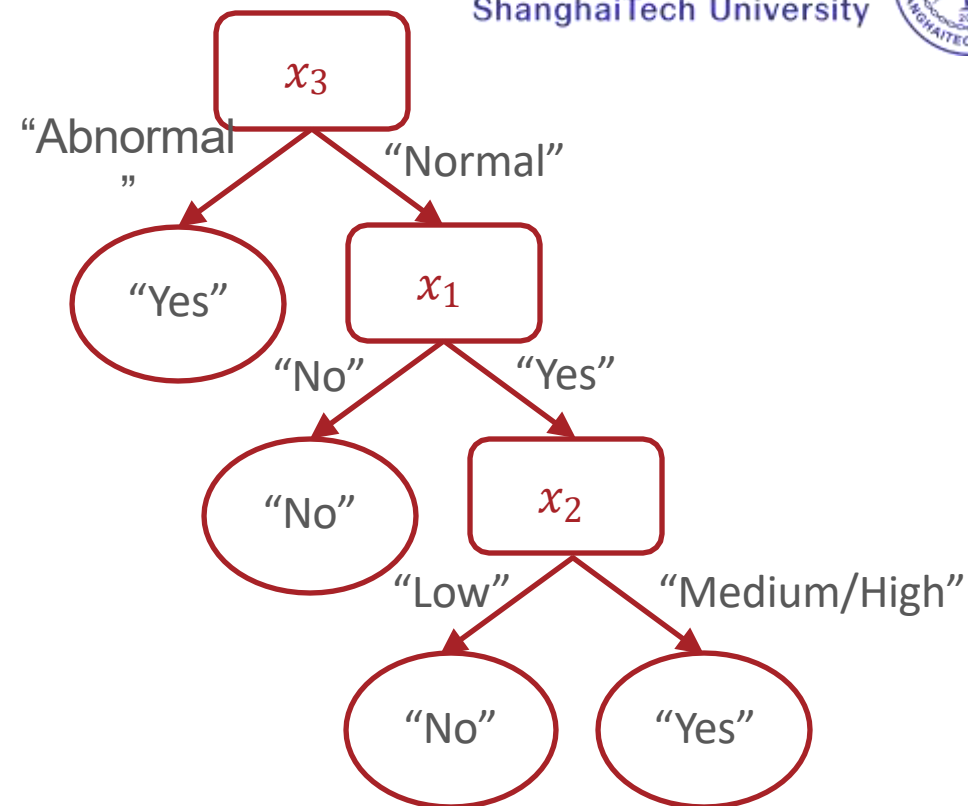
# From Decision Stump to Decision Tree

| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |

# From Decision Stump to Decision Tree

| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |
| | | | |
| No | High | Normal | No] |



11
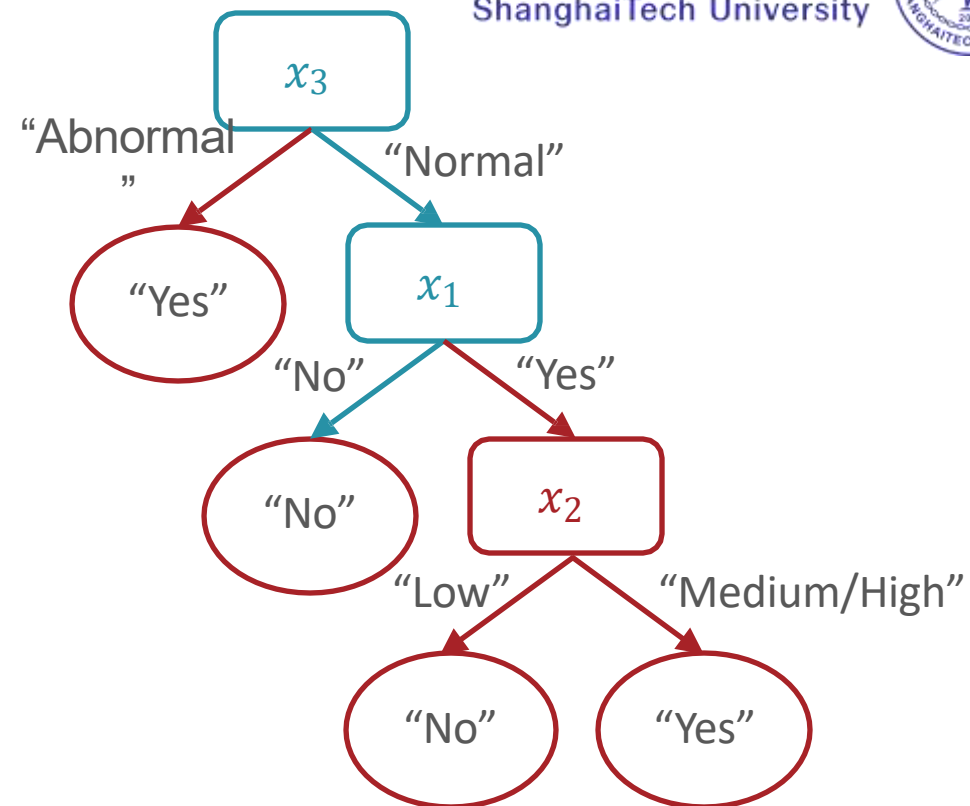
# From Decision Stump to Decision Tree

| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |
|  |  |  |  |
| No | High | Normal | No |

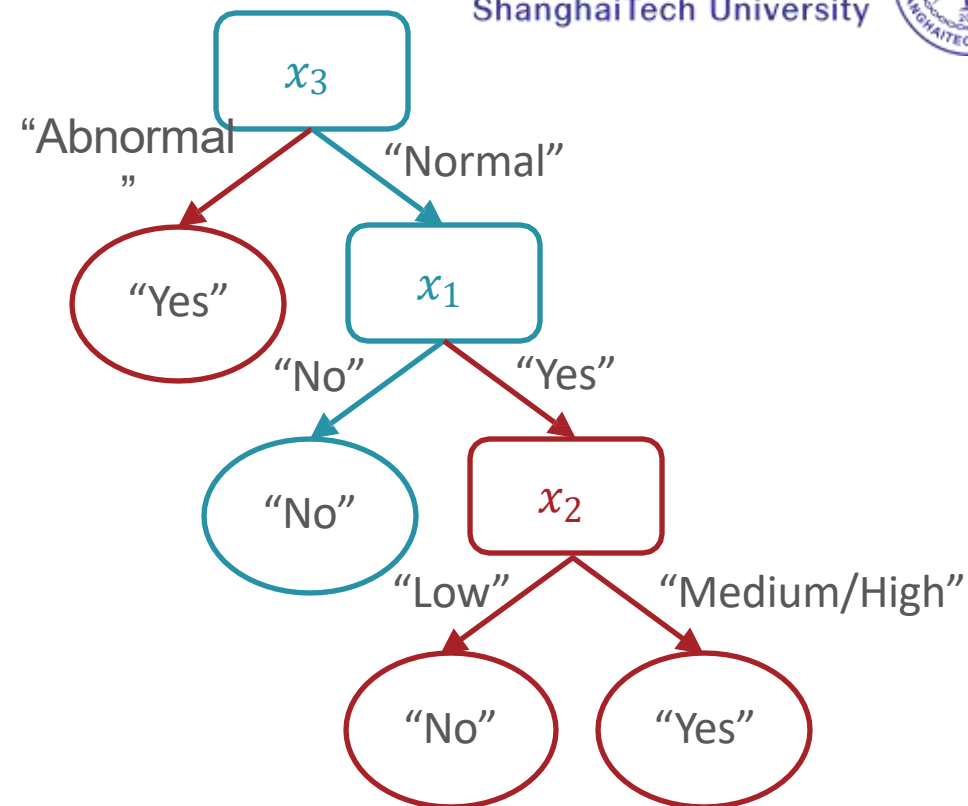# From Decision Stump to Decision Tree

| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |
| | | | |
| No | High | Normal | No |



Decision tree:
- $x_3$
  - "Abnormal" → "Yes"
  - "Normal" → $x_1$
    - "No" → "No"
    - "Yes" → $x_2$
      - "Low" → "No"
      - "Medium/High" → "Yes"

# From Decision Stump to Decision Tree

| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |
| | | | |
| No | High | Normal | No |

$x_3$

"Abnormal"     "Normal"

"Yes"     $x_1$

"No"     "Yes"

"No"     $x_2$

"Low"     "Medium/High"

"No"     "Yes"

## Decision Tree: Pseudocode

```
def h(𝒙'):

   - walk from root node to a leaf node

      while(true):

         if current node is internal (non-leaf):

            check the associated attribute, $x_d$

            go down branch according to $x'_d$

         if current node is a leaf node:

            return label stored at that leaf
```

## Decision Tree Questions

1. How can we pick which feature to split on?

2. Why stop at just one feature?

   a) If we split on more than one feature, how do we decide the order to spilt on?
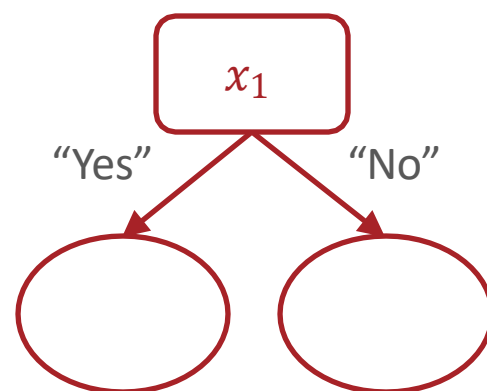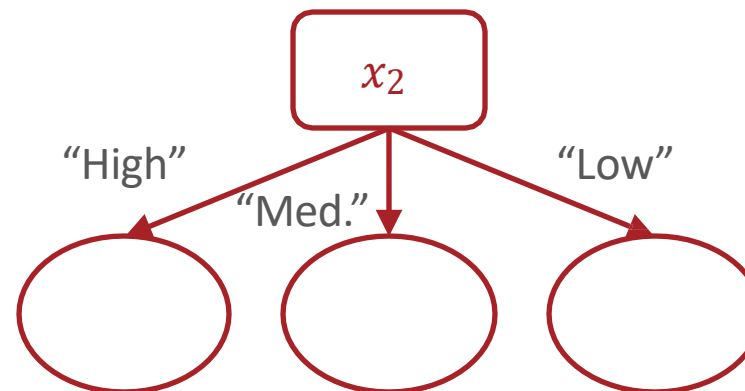
# Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*

- Idea: when deciding which feature to split on, use the one that optimizes the splitting criterion

| $x_1$ Family History | $x_2$ Resting Blood Pressure | $x_3$ Cholesterol | $y$ Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |

Training Error Rate as a Splitting Criterion



$x_1$
"Yes"   "No"

Training error rate:

$x_2$
"High"   "Med."   "Low"

Training error rate:

$x_3$
"Abnormal"   "Normal"

Training error rate:

Poll Question 2:

■ Which feature would you split on using training error rate as the splitting criterion?

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |

A. $x_1$

B. $x_2$

C. Either $x_1$ or $x_2$

D. Neither $x_1$ nor $x_2$

上 海 科 技 大 学
ShanghaiTech University

Poll Question 2:

■ Which feature would you split on using training error rate as the splitting criterion?

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |



Training error rate:

## Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*

- Idea: when deciding which feature to split on, use the one that optimizes the splitting criterion

- Potential splitting criteria:
  - Training error rate (minimize)
  - Gini impurity (minimize) → CART algorithm
  - Mutual information (maximize) → ID3 algorithm

# Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*

- Idea: when deciding which feature to split on, use the one that optimizes the splitting criterion

- Potential splitting criteria:
    - Training error rate (minimize)
    - Gini impurity (minimize) → CART algorithm
    - **Mutual information** (maximize) → ID3 algorithm

# Entropy

- The **entropy** of a *random variable* describes the uncertainty of its outcome: the higher the entropy, the less certain we are about what the outcome will be.

$$H(X) = - \sum_{v \in V(X)} P(X = v) \log_2\big(P(X = v)\big)$$

where $X$ is a (discrete) random variable

$V(X)$ is the set of possible values $X$ can take on

# Entropy

- The **entropy** of a *set* describes how uniform or pure it is: the higher the entropy, the more impure or "mixed-up" the set is

$$H(S) = - \sum_{v \in V(S)} \frac{|S_v|}{|S|} \log_2 \left( \frac{|S_v|}{|S|} \right)$$

where $S$ is a collection of values,

$V(S)$ is the set of unique values in $S$

$S_v$ is the collection of elements in $S$ with value $v$

- If all the elements in $S$ are the same, then

$$\mathrm{H}(S) = -1 \log_2 (1) = 0$$

# Entropy

- The **entropy** of a *set* describes how uniform or pure it is: the higher the entropy, the more impure or "mixed-up" the set is

$$H(S) = - \sum_{v \in V(S)} \frac{|S_v|}{|S|} \log_2 \left( \frac{|S_v|}{|S|} \right)$$

where $S$ is a collection of values,

$V(S)$ is the set of unique values in $S$

$S_v$ is the collection of elements in $S$ with value $v$

- If $S$ is split fifty-fifty between two values, then

$$\text{H}(S) = -\frac{1}{2} \log_2 \left( \frac{1}{2} \right) - \frac{1}{2} \log_2 \left( \frac{1}{2} \right) = -\log_2 \left( \frac{1}{2} \right) = 1$$

# Mutual Information

- The **mutual information** between *two random variables* describes how much clarity knowing the value of one random variables provides about the other

$$I(Y; X) = H(Y) - H(Y|X)$$

$$= H(Y) - \sum_{v \in V(X)} P(X = v)H(Y|X = v)$$

where $X$ and $Y$ are random variables

$V(X)$ is the set of possible values $X$ can take on

$H(Y|X = v)$ is the conditional entropy of $Y$ given $X = v$

## Mutual Information

- The **mutual information** between *a feature and the label* describes how much clarity knowing the feature provides about the label

$$I(y; x_d) = H(y) - H(y|x_d)$$

$$= H(y) - \sum_{v \in V(x_d)} f_v \left( H(Y_{x_d=v}) \right)$$

where $x_d$ is a feature and $y$ is the set of all labels

$V(x_d)$ is the set of possible values $x_d$ can take on

$f_v$ is the fraction of data points where $x_d = v$

$Y_{x_d=v}$ is the set of all labels where $x_d = v$

# Mutual Information Example

| $x_d$ | $y$ |
|:-:|:-:|
| 1 | 1 |
| 1 | 1 |
| 0 | 0 |
| 0 | 0 |

$$I(x_d, Y) = H(Y) - \sum_{v \in V(x_d)} (f_v)\left(H(Y_{x_d=v})\right)$$

$$= 1 - \frac{1}{2}H(Y_{x_d=0}) - \frac{1}{2}H(Y_{x_d=1})$$

$$= 1 - \frac{1}{2}(0) - \frac{1}{2}(0) = 1$$

# Mutual Information Example

| $x_d$ | $y$ |
|:---:|:---:|
| 1 | 1 |
| 0 | 1 |
| 1 | 0 |
| 0 | 0 |

$$I(x_d, Y) = H(Y) - \sum_{v \in V(x_d)} (f_v)\left(H(Y_{x_d=v})\right)$$

$$= 1 - \frac{1}{2}H(Y_{x_d=0}) - \frac{1}{2}H(Y_{x_d=1})$$

$$= 1 - \frac{1}{2}(1) - \frac{1}{2}(1) = 0$$

## Poll Question 3:

Which feature would you split on using mutual information as the splitting criterion?

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |
| 1 | 1 | 1 |

A. $x_1$

B. $x_2$

C. Either $x_1$ or $x_2$

D. Neither $x_1$ nor $x_2$

# Decision Tree: Questions

1. How can we pick which feature to split on?

2. Why stop at just one feature?

   a) If we split on more than one feature, how do we decide the order to spilt on?

## Decision Tree: Pseudocode

```
def train(𝒟):
    store root = tree_recurse(𝒟)

def tree_recurse(𝒟′):

    q = new node()

    base case – if (SOME CONDITION):

    recursion – else:
        find best attribute to split on, x_d

        q.split = x_d

        for v in V(x_d), all possible values of x_d:
```

$$\mathcal{D}_v = \left\{ \left( x^{(n)}, y^{(n)} \right) \in \mathcal{D} \mid x_d^{(n)} = v \right\}$$

```
            q.children(v) = tree_recurse(𝒟_v)

    return q
```

## Decision Tree: Pseudocode

```
def train(𝒟):
    store root = tree_recurse(𝒟)

def tree_recurse(𝒟′):
    q = new node()
    base case – if (𝒟′ is empty OR
        all labels in 𝒟′ are the same OR
        all features in 𝒟′ are identical OR
        some other stopping criterion):
            q.label = majority_vote(𝒟′)

    recursion – else:
    return q
```