# Decision Trees
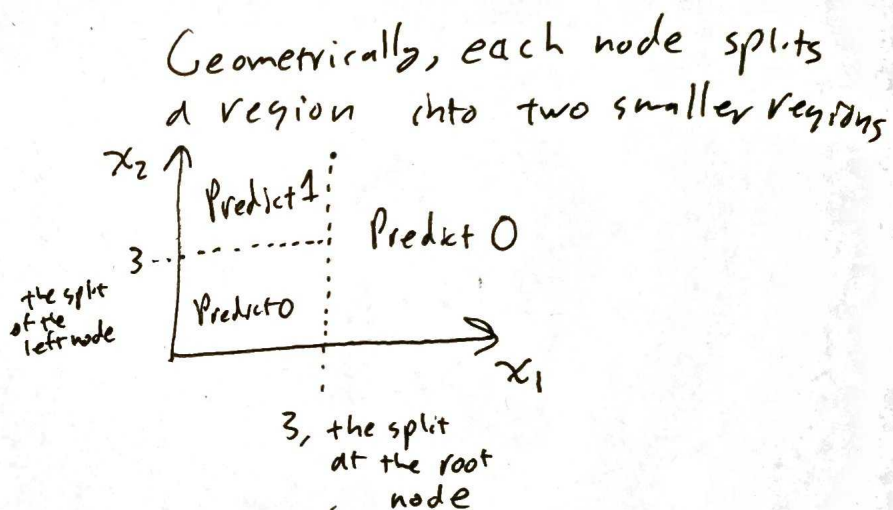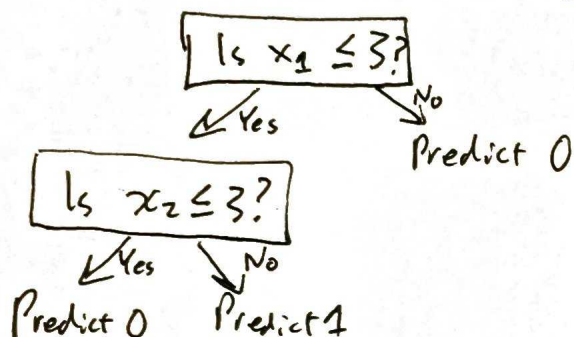
Prediction Rule: $f(x) =$ if ...
      then if ...
        :
        else if ...
          :
          then predict 0
          else predict 1

**Example**

If $(x_1 \leq 3)$
   then if $(x_2 \leq 3)$
      then predict 0
      else predict 1
   else predict 0

We can represent rules of this form as a tree.
Let $x_1, \ldots, x_d$ be the features of the vector $x$.

Is $x_1 \leq 3$?
   Yes / No
     No → Predict 0

Is $x_2 \leq 3$?
   Yes / No
Predict 0    Predict 1

Geometrically, each node splits a region into two smaller regions

$x_2$ ↑
   Predict 1    Predict 0
3 -- - - - - - - -
   Predict 0
            → $x_1$

the split of the left node

3, the split at the root node

We need to choose what kind of decisions can appear in a node.
We will follow these conventions:

- Every decision is of the form "Is $x_f \leq t$?" where $x_f$ is one feature of our feature vectors, and $t$ is a threshold value.

- The threshold $t$ is chosen only among midpoint values of neighboring values for that feature at that node.
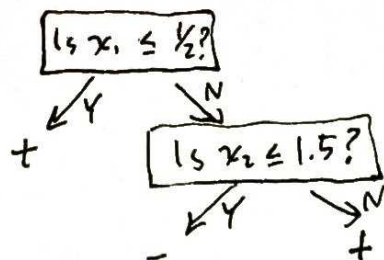
Note: with these conventions, when finding a decision rule for a node, there are only finitely many possibilities.

- Each leaf corresponds to label, which is the prediction made when $x$ lies in the region that reaches that leaf in the tree.
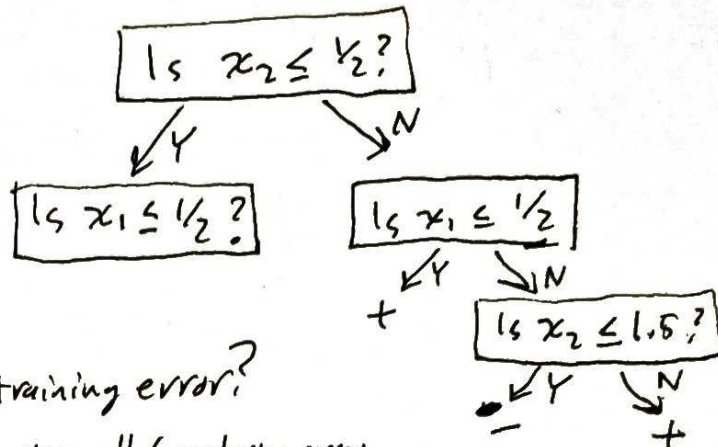
Note: multiple leaves might have the same predicted label.

Example: training data $((0,0),+)$ $((0,1),+)$ $((0,2),+)$
$((1,0),-)$ $((1,1),-)$ $((1,2),+)$

Possible Tree 1:

Is $x_1 \leq \frac{1}{2}$?
  Y → $+$
  N → Is $x_2 \leq 1.5$?
    Y → $-$
    N → $+$

Possible Tree 2:

Is $x_2 \leq \frac{1}{2}$?
  Y → Is $x_1 \leq \frac{1}{2}$?
  N → Is $x_1 \leq \frac{1}{2}$
    Y → $+$
    N → Is $x_2 \leq 1.5$?
      Y → $-$
      N → $+$

Which Tree is better? What is their training error?

→ They both have 0 training error, getting all 6 predictions correct.

We prefer the simpler one (Tree 1).

In general finding the smallest decision tree that fits a training dataset is computationally difficult.

A commonly used algorithm that proceeds greedily is the ID3 Decision Tree algorithm

### ID3 Algorithm

1. Initialize with all the data at a single root node.

2. While there is an impure node in the tree:

   2.1. Pick any impure node $v$ ✓

   2.2. Pick the best decision rule to split the data at $v$. ("best" to be defined later)

   2.3. Modify the tree by replacing $v$ with this decision plus two child nodes

   Is $x_i \leq t$?
     Y → $v_L$
     N → $v_R$

   2.4. If either new child $v_L$ or $v_R$ are pure, make them a leaf that predicts their single label.

Impure node: a node with data points of different labels

Each current leaf node corresponds to a subset of the training data that would reach that node.

We will discuss how we select the best decision rule (step 2.2) later.

**Example:** Training data $((0,0), 1)$ $((0,1), 0)$ $((1,0), 0)$ $((2,0), 0)$

            ↑ a        ↑ b        ↑ c      ↑ d
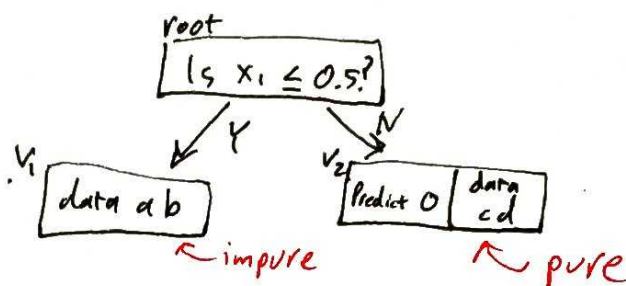
**Initially,** just a root node

root
| data abcd |

It is impure (label 0 and label 1 are present)

We select it to split. We have 3 candidate decision rules.

( Is $x_1 \leq 0.5$? ) ( Is $x_1 \leq 1.5$? ) ( Is $x_2 \leq 0.5$? )

Suppose we pick ( $x_1 \leq 0.5$? ) as our decision.

**Update the tree:**

root
| Is $x_1 \leq 0.5$? |

     Y ↙        ↘ N

$V_1$ | data a b |       $V_2$ | Predict 0 | data c d |
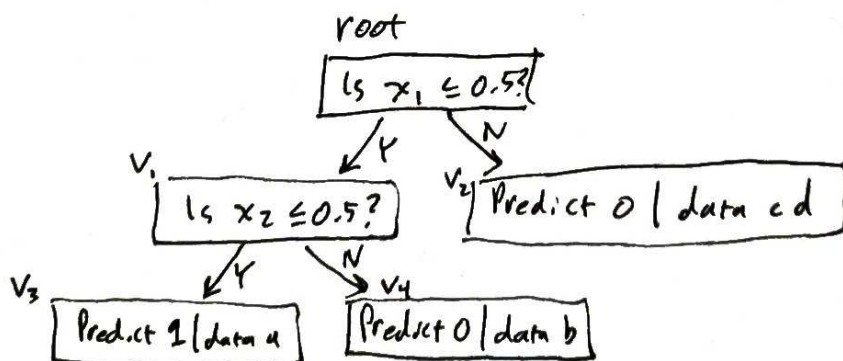
       ← impure           ↖ pure

Now $V_1$, containing a,b, is impure.
Suppose we select it to split. We have only one candidate decision rule.
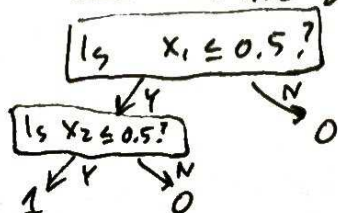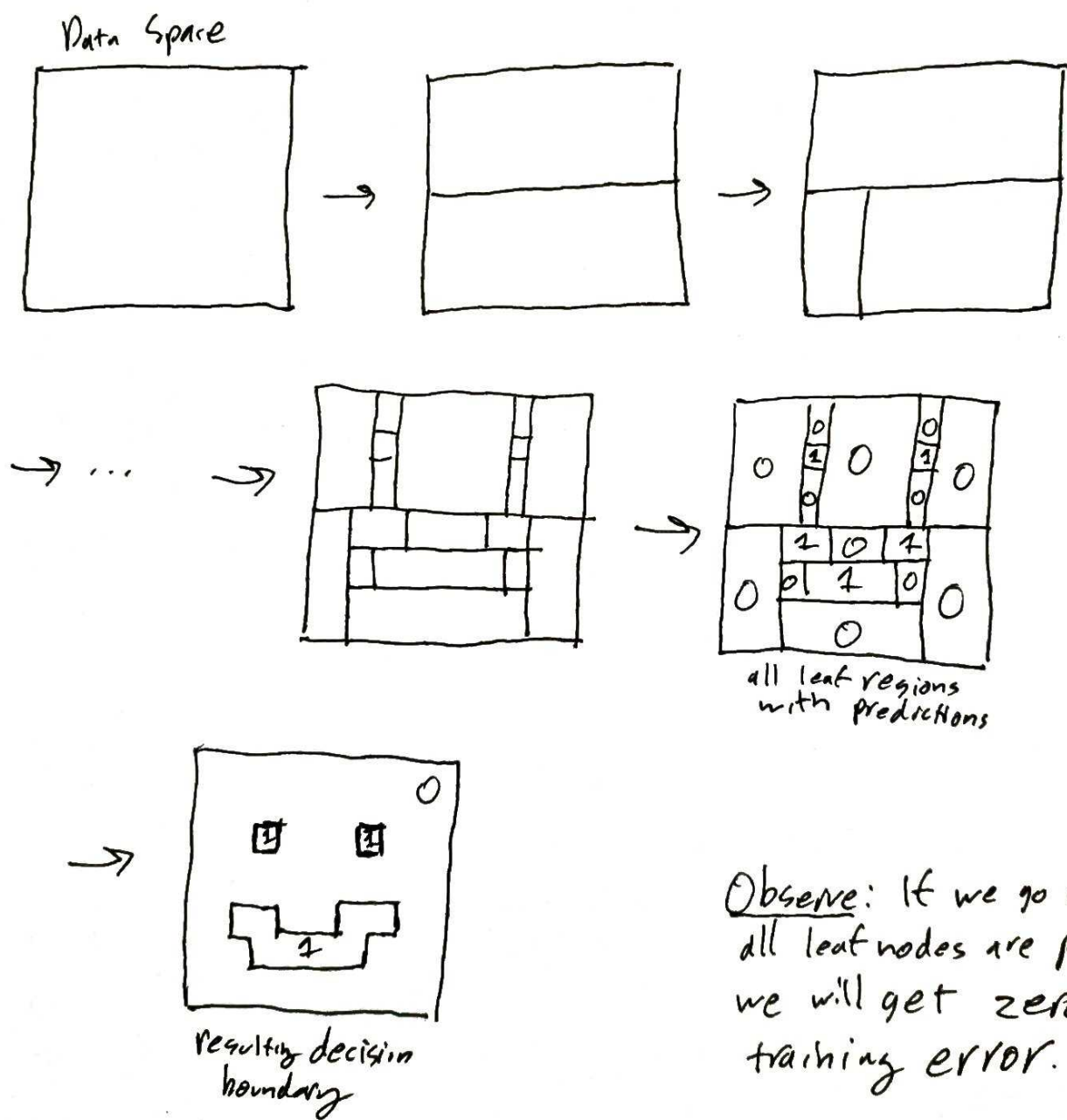We pick ( $x_2 \leq 0.5$? ) as the decision rule for $V_1$

**Update the tree:**

root
| Is $x_1 \leq 0.5$? |

   Y ↙     ↘ N

$V_1$ | Is $x_2 \leq 0.5$? |      $V_2$ | Predict 0 | data c d |

   Y ↙    ↘ N   $V_4$

$V_3$ | Predict 1 | data a |     | Predict 0 | data b |

Now all our leaves are pure. We are done.

**Final tree:**

| Is $x_1 \leq 0.5$? |

   Y ↙    ↘ N
          0

| Is $x_2 \leq 0.5$? |

  Y ↙   ↘ N
  1     0

# Geometry of Decision Trees

With our convention of "is $x_f \leq t$?" decision rules, each node splits a region into two regions, where the boundary is axis-aligned. With no limit on the number of nodes, this allows the decision tree to carve out arbitrarily complicated decision boundaries

Data Space



all leaf regions with predictions



resulting decision boundary

Observe: If we go until all leaf nodes are pure, we will get zero training error.

# Selecting a Decision Rule

In the ID3 algorithm we select the decision rule that reduces uncertainty the most.

We measure uncertainty through a quantity from information theory called entropy

Entropy: Let $X$ be a random variable that takes values $v_1, \ldots, v_K$ with probabilities $p_1, \ldots, p_K$. Then, the entropy of $X$, denoted by $H(x)$, is defined as

$$H(x) = -\sum_{i=1}^{K} p_i \log p_i = -p_1 \log p_1 - p_2 \log p_2 - \cdots - p_K \log p_K$$

note: by convention, we will declare $0 \log 0 = 0$ even though $\log 0$ is undefined.

note: by convention, we will only use the natural log, $(\ln$ or $\log_e)$, with the base $e$ in entropy calculations

Example $X$: a $0/1$ random variable (r.v.) with $Pr(x=1) = \frac{1}{2}$

$$v_1 = 0 \quad v_2 = 1 \quad P_1 = \frac{1}{2} \quad P_2 = \frac{1}{2}$$

$$H(x) = -P_0 \log P_0 - P_1 \log P_1 = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = -\log \frac{1}{2} = \log 2$$

$$H(x) = \log 2 \approx 0.69$$

Example $X$: a $0/1$ r.v. with $Pr(x=1) = 0$

$$v_1 = 0$$
$$v_2 = 1 \quad P_1 = 1 \quad P_2 = 0$$

$$H(x) = -1 \log 1 - 0 \log 0 = -0 - 0 = 0$$

Observe: for a 0/1 r.v., the closer $P_1(x=1)$ is to $1/2$, the higher the entropy.

Does this match what one might think, with the entropy being a measure of uncertainty?

Example $X$ is a r.v. that takes values $1, 2, ..., k$ each with probability $1/k$

$$H(x) = -\frac{1}{k}\log\frac{1}{k} - ... - \frac{1}{k}\log\frac{1}{k} = -K\frac{1}{k}\log\frac{1}{k} = -\log\frac{1}{k}$$
$$= \log K$$

Properties of Entropy

- $H(x)$ does not depend on the values taken by $X$ $(v_1, ..., v_k)$ but only on the probabilities $(P_1, ..., P_k)$ of these distinct values.

- $H(x) \geq 0$ always. $\qquad H(x) = -\sum_i P_i \log P_i$

$0 \leq P_i \leq 1$ $\qquad$ Is $P_i$ positive? Is $\log P_i$ positive?

- If $X$ takes $K$ values, $H(x) \leq \log K$. This maximum is achieved only when each value happens w.p. $1/k$

Conditional Entropy

This measures the uncertainty when we can condition on other information.

Let $X, Z$ be two r.v.s. The conditional entropy of $X$ given $Z$ is defined as

$$H(X|Z) = \sum_v Pr(Z=v) H(X|Z=v)$$

<span style="color:red">uses $Pr(X=v'|Z=v)$ instead of $Pr(X=v')$ in the entropy calculation.</span>

# Example   $X, Z$   joint distribution

| Z\X | 0 | 1 |
|---|---|---|
| 0 | ⅓ | ⅓ |
| 1 | ⅓ | 0 |

Question: Are $X$ and $Z$ independent?
    (No.)

What is   $H(X|Z)$?

We need   $Pr(Z=0), Pr(Z=1)$, as well as ~~the~~ $H(X|Z=0), H(X|Z=1)$

$Pr(Z=0) = \frac{2}{3}$      $Pr(Z=1) = \frac{1}{2}$

$Pr(X=v | Z=0) = \dfrac{Pr(X=v, Z=0)}{2/3}$

$Pr(X=v | Z=1) = \dfrac{Pr(X=v, Z=1)}{1/3}$

$\longrightarrow$   $Pr(X=0 | Z=0) = \frac{1}{2}, \quad Pr(X=1 | Z=0) = \frac{1}{2}$
     $Pr(X=0 | Z=1) = 1, \quad Pr(X=1 | Z=1) = 0$

$H(X | Z=0) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = \log 2$

$H(X | Z=1) = -1 \log 1 - 0 \log 0 = 0$

$H(X|Z) = \frac{2}{3} \log 2 + \frac{1}{3} \cdot 0 = \frac{2}{3} \log 2 \approx 0.462$

Exercise: How does the value of $H(X|Z)$ compare to
    the value of $H(X)$?
    Does knowing $Z$ reduce our uncertainty about $X$?

# Example

Joint Distribution:

| $\frac{}{}$ X $\frac{}{}$ Z | 0 | 1 | 2 |
|---|---|---|---|
| 0 | 0 | 1/2 | 1/4 |
| 1 | 1/12 | 1/12 | 1/12 |

$Pr(z=0) = 3/4$

$Pr(z=1) = 1/4$

Conditional Distributions

$Pr(x=0 \mid z=0) = 0$

$Pr(x=1 \mid z=0) = 2/3$

$Pr(x=2 \mid z=0) = 1/3$

$Pr(x=0 \mid z=1) = 1/3$

$Pr(x=1 \mid z=1) = 1/3$

$Pr(x=2 \mid z=1) = 1/3$

$$H(x \mid z=0) = -0 \log 0 - \frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3}$$

$$= 0 - \frac{2}{3} \log 2 + \frac{2}{3} \log 3 - \frac{1}{3} \log \frac{1}{3}$$

$$= 0 - \frac{2}{3} \log 2 + \log 3 = -\frac{2}{3} \log 2 + \log 3$$

$$H(x \mid z=1) = -\frac{1}{3} \log \frac{1}{3} - \frac{1}{3} \log \frac{1}{3} - \frac{1}{3} \log \frac{1}{3} = \log 3$$

$$H(x \mid z) = Pr(z=0) H(x \mid z=0) + Pr(z=1) H(x \mid z=1)$$

$$= \frac{3}{4} \left( \log 3 - \frac{2}{3} \log 2 \right) + \frac{1}{4} \log 3$$

$$= \log 3 - \frac{1}{2} \log 2$$

# Properties of Conditional Entropy

- Suppose $X = Z$, what is $H(X|Z)$?

$$H(X|Z) = \sum_v Pr(Z=v) H(X|Z=v)$$

Given $Z=v$, we know $X=v$ w.p. $1$

So $H(X|Z=v) = 0$

Thus $H(X|Z) = \sum_v Pr(Z=v) \cdot 0 = 0$

- Suppose $X$ and $Z$ are independent, what is $H(X|Z)$?

Given independence, $Pr(X=v'|Z=v) = Pr(X=v')$
for any values $v, v'$

This means $X|Z=v$ has the same distribution as $X$

So $H(X|Z=v) = H(X)$

$$\rightarrow H(X|Z) = \sum_v Pr(Z=v) H(X|Z=v) = \sum_v Pr(Z=v) H(X)$$

$$= H(X) \sum_v Pr(Z=v) = H(X) \cdot 1 = H(X)$$

# Information Gain

The Information Gain of a variable $Z$ on another variable $X$
is defined as the difference $H(X) - H(X|Z)$

i.e. how much of the entropy (uncertainty) of $X$ is reduced
by knowing $Z$.

It can be shown the information gain is always $\geq 0$.

Back to our decision trees...

Selection Rule for Decisions: The "_best_" decision rule
is the one with the largest information gain,
i.e. the decision whose result will most reduce the uncertainty
we have in predicting a label.

Let $Y$ be a r.v. corresponding to the distribution of labels
among the data at a node. (Ex. $\frac{2}{5}$ data points have label 0
$3/5$ datapoints have label 1)

Let $Z$ be a r.v. corresponding to the distribution of
answers to a decision rule among the data at a node.
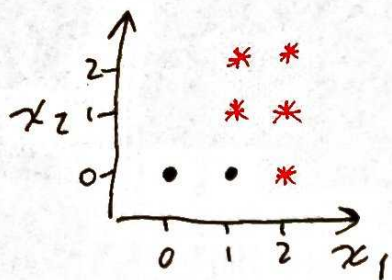ex. $Z=0$ corresponds to "no", $Z=1$ corresponds to "yes"

$Y | Z=0$: distribution of labels at this node
only among those data points with $x_\varepsilon > t$.

$Y | Z=1$: distribution of labels at this node
only among those data points with $x_\varepsilon \leq t$.

We pick the decision rule such that its r.v. $Z$ maximizes
the information gain on $Y$.

maximize $\quad H(Y) - H(Y|Z)$

# Example

Consider these seven training points



what is $H(Y)$? 5 points have label $*$
2 points have label $\bullet$

$$H(Y) = -\frac{5}{7}\log\frac{5}{7} - \frac{2}{7}\log\frac{2}{7} \approx 0.598$$

How many candidate decision rules do we have?

4: $x_1 \le 0.5?$    $x_2 \le 0.5?$
    $x_1 \le 1.5?$    $x_2 \le 1.5?$

Let's calculate the information gain for "$x_2 \le 1.5?$"

$$Z=0 \iff \text{"no."} \iff x_2 > 1.5$$
$$Z=1 \iff \text{"yes."} \iff x_2 \le 1.5$$



region where $Z=0$

$x_2 = 1.5$

region where $Z=1$

So $Y | Z=0$ has 2 $*$ labels and 0 $\bullet$ labels

$$Pr(Y=* \mid z=0) = 1$$
$$Pr(Y=\bullet \mid z=0) = 0$$
$$H(Y \mid z=0) = 0$$

and $Y | Z=1$ has 3 $*$ labels and 2 $\bullet$ labels
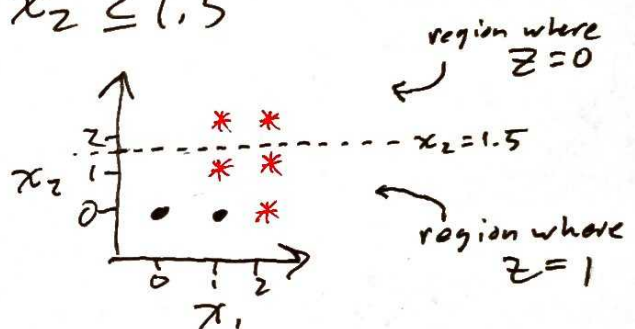
$$Pr(Y=* \mid z=1) = 3/5$$
$$Pr(Y=\bullet \mid z=1) = 2/5$$
$$H(Y \mid Z=1) = -\frac{3}{5}\log\frac{3}{5} - \frac{2}{5}\log\frac{2}{5} \approx 0.673$$

Also $Pr(z=0) = \frac{2}{7}$, $Pr(z=1) = \frac{5}{7}$

So $H(Y|Z) = Pr(z=0)H(Y|z=0) + Pr(z=1)H(Y|z=1) \approx 0.48$

Information Gain: $H(Y) - H(Y|Z) \approx 0.12$

Let's quickly consider the other decision rules.

$x_1 \le 0.5?$ :  $z=1: 0 *, 10 \bullet$
$\qquad\qquad z=0: 5 *, 1 \bullet$
$\qquad\qquad H(Y|z) \approx 0.38$
$\qquad\qquad$ info gain $\approx 0.22$

The distributions of counts of labels on the two sides of the decision are the important part in calculating information gain.

---

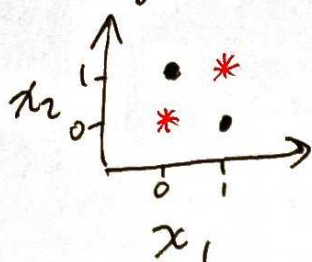$x_1 \le 1.5?$ :  $z=1: 2 *, 2 \bullet$
$\qquad\qquad z=0: 3 *, 0 \bullet$
$\qquad\qquad H(Y|z) \approx 0.39$
$\qquad\qquad$ info gain $\approx 0.21$

**Exercise:** try computing these values.

---

$x_2 \le 0.5?$ :  $z=1: 1 *, 2 \bullet$
$\qquad\qquad z=0: 4 *, 0 \bullet$
$\qquad\qquad H(Y|z) \approx 0.27$
$\qquad\qquad$ info gain $\approx 0.33$

It might help to plot the boundaries of these decision rules

---

$x_2 \le 1.5?$ :  $z=1: 3 *, 2 \bullet$
(we computed
this one earlier)  $z=0: 2 *, 0 \bullet$
$\qquad\qquad H(Y|z) \approx 0.48$
$\qquad\qquad$ info gain $\approx 0.12$

---

So ID3 will select "Is $x_2 \le 0.5?$"

to get

Is $x_2 \le 0.5$?

$V_L$ — impure — $(0,0,\bullet)$ $(0,1,\circ)$ $(0,2,\bullet)$

$V_R$ — Predict $*$ — $(1,1,*)$ $(1,0,*)$ $(2,1,*)$ $(2,0,*)$ — pure

... and the process will continue, since we still have an impure node, ...
**Exercise:** The node $V_L$ has only two candidate decision rules. What are they? Which will be chosen?

# Example:

Training Data



Initially, $2\bullet$ and $2*$
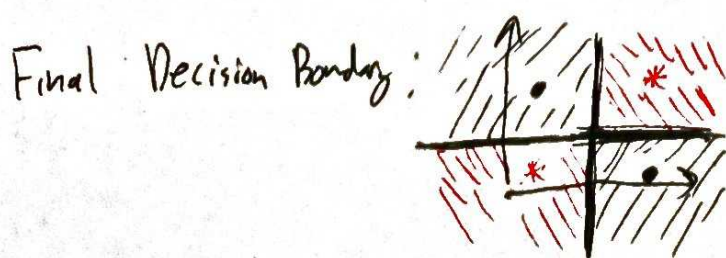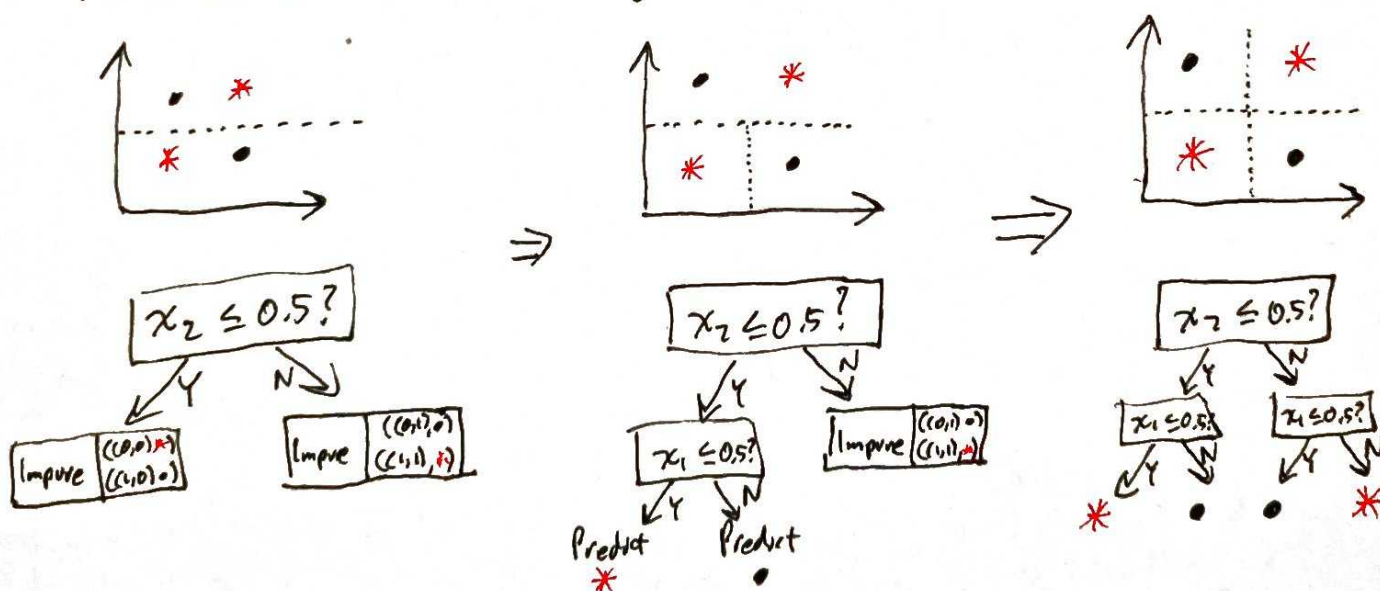
$$H(Y) = \log 2 \approx 0.69$$

Two candidate splits

$x_1 \le 0.5$? : $H(Y|Z) = \frac{1}{2}\log 2 + \frac{1}{2}\log 2 = \log 2 \approx 0.69$

info gain $= 0$

$x_2 \le 0.5$? : $H(Y|Z) = \frac{1}{2}\log 2 + \frac{1}{2}\log 2 = \log 2 \approx 0.69$

info gain $= 0$

The best info gain is 0. Do we stop? <u>No</u>. We still have impure nodes.

There is a tie? Randomly select one of the tied rules.



Final Decision Bondary:



Training Error: $0/4$

This is always possible if we assume every training point has a distinct feature vector.

When all are pure, we get zero training error.
But this might <u>overfit</u> the data.

Example      split until pure      decision boundary



That isolated ● point might be an outlier, or a misrecorded label. Having that isolated region of predicting ● might lead to many mistakes on unseen test data. A simpler tree might be better!



Recall the example about 1-NN vs. 3-NN vs. 5-NN.

What is Overfitting?

Overfitting is when a model or prediction rule too closely matches the peculiarities of a specific set of data, at the cost of performing well on additional unseen data.

We can observe overfitting by looking at the <u>training error</u> versus the <u>true error</u> of a classifier, as we vary the complexity of the classifier.

Assume $(x, y) \sim D$, and we have a training set
$$(x_1, y_1) \cdots (x_n, y_n) \text{ sampled from } D$$

Training error, for a prediction rule $h(x)$
$$\hat{err}(h) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}(h(x_i) \neq y_i)$$

True error, for a prediction rule $h(x)$
$$err(h) = \Pr_{(x,y) \sim D} (h(x) \neq y)$$

$\mathbb{1}(\cdot)$ is an indicator function

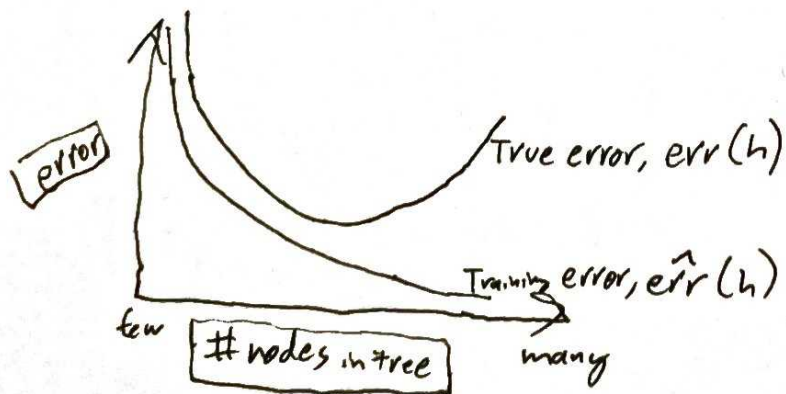$$= \begin{cases} 1 & \text{if condition is true} \\ 0 & \text{if condition is false} \end{cases}$$

(With an unseen test set, we can get a good estimate of $err(h)$ using the test set error rate.)

For a fixed learning algorithm, generally more data helps ensure
$$\hat{err}(h) \text{ approaches } err(h).$$ More data means outliers are less misleading.

If we plot $\hat{err}(h)$ vs $err(h)$ for a fixed training set, and vary the complexity of the model returned by the ~~test~~ learning algorithm (e.g. number of nodes allowed in a decision tree) we usually see this behavior:



True error, $err(h)$

Training error, $\hat{err}(h)$

few   #nodes in tree   many

As the complexity increases, the algorithm can learn more from the data, and the error decreases.

At some point, the true error stops improving and may get worse. Past this point, we say the algorithm is over fitting.
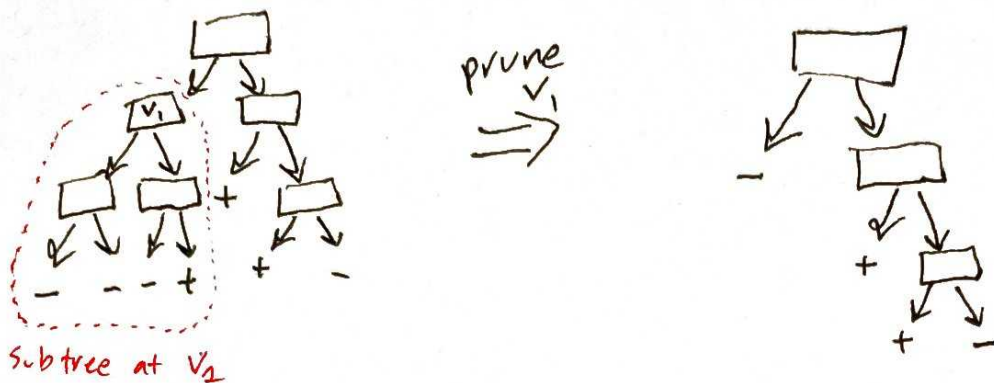
The distribution $D$ has some structure, and we want to learn it for making predictions. The training data shows some of this structure, but it might also have some spurious structure of its own by random chance, which we should not be modelling in our predictions.

# Reducing the Complexity of a Decision Tree: Pruning

After splitting until all nodes are pure, our tree may be very complex, many nodes.

Some of these nodes might be carving out ~~some~~ small precise regions that overfit.

So we can prune, or cut out entire subtrees at a node



Subtree at $V_1$

We replace a ~~leaf~~ decision node with a leaf that predicts the majority label of the data at that new leaf.

But how can we tell which nodes to prune?

With not-yet-seen validation data to estimate the true error.

## ID3 with Pruning

1. Split our labelled data into two sets: training set $S$ and validation set $V$

2. Build full tree $T$ with ID3 on the training set $S$.
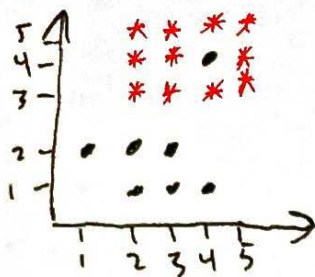
3. Prune the tree:
   For each node $v$ in $T$
   Make $T'$, which is $T$ modified to prune at $v$ $\left(\begin{array}{l}\text{replace subtree at } v \text{ with} \\ \text{a majority-predicting leaf}\end{array}\right)$

   If $error(T')$ on $V \leq error(T)$ on $V$, then update $T$ to keep this prune
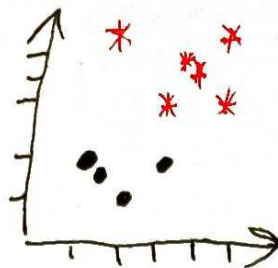   $\left(\text{update } T = T'\right)$

---

If $error(T') \leq error(T)$, then the validation data suggests ~~that~~ pruning this node is a good idea.
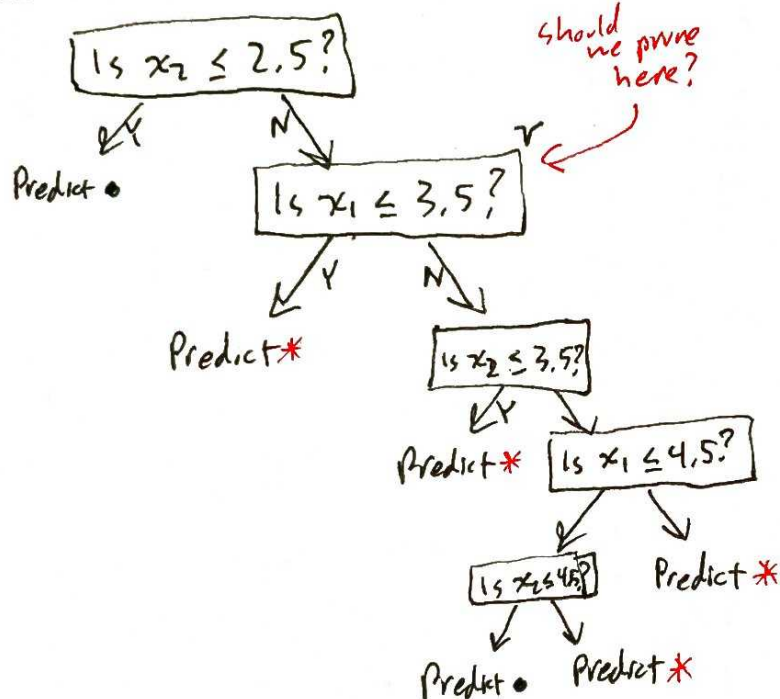
# Example

### Training Data



### Validation Data



## tree T:

Is $x_2 \leq 2.5$?
- Y → Predict •
- N → Is $x_1 \leq 3.5$?
  - Y → Predict *
  - N → Is $x_2 \leq 3.5$?
    - Y → Predict *
    - N → Is $x_1 \leq 4.5$?
      - Is $x_2 \leq 4.5$?
        - Predict •
        - Predict *
      - Predict *

should we prune here?

v

## Tree T':

Is $x_2 \leq 2.5$?
- Y → Predict •
- N → Predict *

majority vote,
11 *
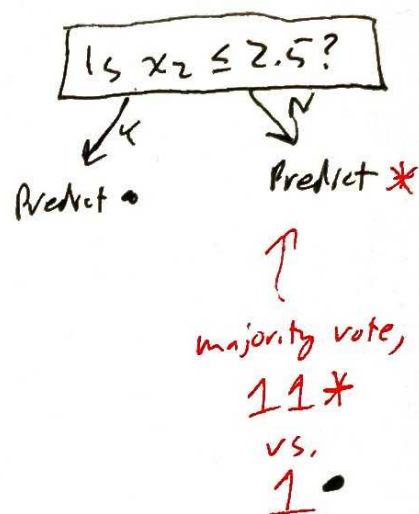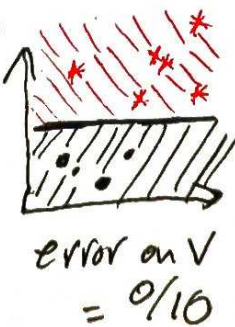vs.
1 •

Here our validation data suggests the isolated • point is misleading, and the error rates on V suggest we should prune at node v

**T on V**



error on V
$= 2/10$

**T' on V**



error on V
$= 0/10$

By improving our validation error, we have reason to believe that this pruning will also improve the true error, since the validation error is a measure of how well this tree performs on some unseen data drawn from $\mathcal{D}$.

Let's take a moment to compare our ID3 learning
algorithm with K-NN

| ID3 with pruning | K-NN |
|---|---|

**ID3 with pruning**

- longer training time (evaluate all candidate decisions at every node)

- shorter prediction run time (just answer some simple questions at each node as we descend the tree)

- unpruned tree gets zero training error

- we can prune to avoid overfitting
  Use validation data to select when and where to prune

- storing the model is just storing a tree of decisions

- somewhat interpretable, we can look at the tree to see what was found to be most informative.
  We can also look at regions of its decision boundary

- a little trickier to implement

**K-NN**

- short training procedure (memorize training data)

- long prediction run time (compute all distances...)
      scales poorly with #training points and #dimensions.

- 1-NN gets zero training error

- we can increase k to avoid overfitting
  Use validation data to select k

- storing the model involves storing all of the training data

- not very interpretable, but we could try to look at regions of its decision boundary.

- simple to implement

When might ID3 be better than k-NN? When might k-NN be better?

If we did not care about run time or storage costs, how might we use data
      to tell us which is a better choice for a certain task?