# Lec 10: Decision Trees

Toy example:
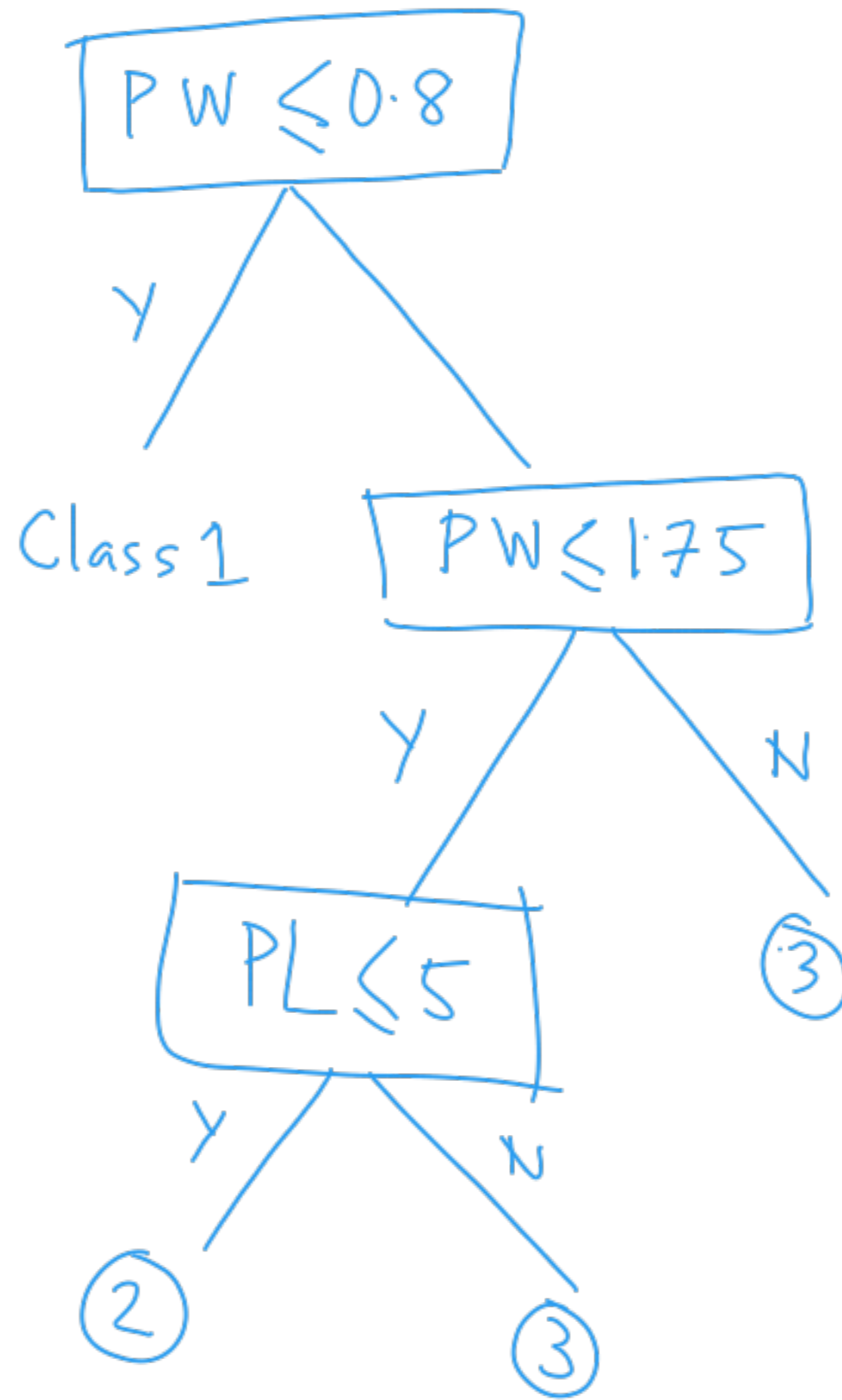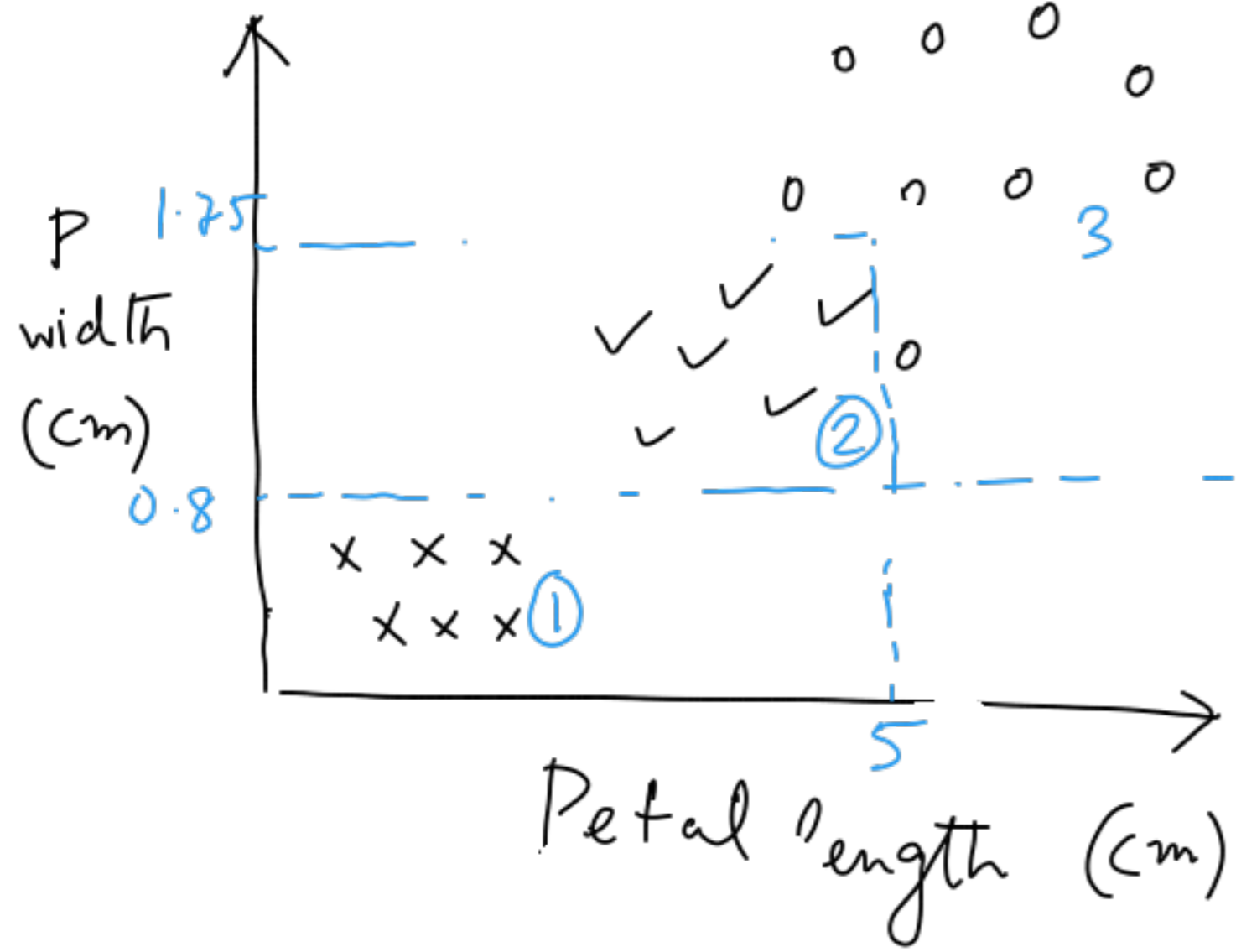
$$x$$

| Exam result $y$ | Online Courses? | Background | Mock Tests |
|---|---|---|---|
| P | Y | Math | N |
| F | N | M | Y |
| F | Y | M | Y |
| P | Y | CS | N |
| ⋮ | ⋮ | ⋮ | ⋮ |

Goal: create a classifier on whether a given student will pass/not depending on the data

## Mock Tests

```
        Mock Tests
        /        \
      Y/          \N
      /            \
  [5P, 4F]      [5P, 1F]
```

**1st stage**

Background
———————
is better to
split on

```
              Background
            /      |      \
        Math/    CS|       \Others
          /        |        \
     [4P, 3F]   [4P, 0F]   [0P, 4F]
     Mock Test
      /      \
     Y/       \N
     /         \
 [2P, 2F]    [3P, 0F]
 Online courses
   /      \
  Y/       \N
  /         \
[1P, 1F]  [1P, 1F]
```

**Q:** how can we make such a splitting scheme more systematic?

*not much useful info at this split.*

# Ex. 2    Iris dataset

P width (cm)

1.75

0.8

Petal length (cm)

5

×  ×  ×
×  ×  ×①

✓  ✓  ✓
✓  ✓  ✓②

o  o  o
o  o  n  o
o  o  ③  o

$PW \leq 0.8$

Y

Class 1

$PW \leq 1.75$

Y          N

$PL \leq 5$          ③

Y       N

②          ③

Q1: How to build the tree?

Q2: Where to stop?

# Ex 3

| $X_1$ | $X_2$ | $Y$ |
|-------|-------|-----|
| T | T | T |
| T | F | T |
| T | T | T |
| T | F | T |
| F | T | T |
| F | F | F |
| F | T | F |
| F | F | F |

$I(Y; X_1)$ vs $I(Y; X_2)$

$H(Y) - H(Y|X_1)$      $H(Y) - H(Y|X_2)$

$$\boxed{H(Y|X_1)} = \sum_{x_1 \in \{T,F\}} p(X_1 = x_1) H(Y|X_1 = x_1)$$

$= p(X_1=T) \, H(Y|X_1=T)^{\nearrow 0}$
$\quad + p(X_1=F) \, \underbrace{H(Y|X_1=F)}$

$X_1$

T / \ F

$Y=T : 4$     $Y=T : 1$
$Y=F : 0$     $Y=F : 3$

$X_2$

T / \ F

$Y=T : 3$     $X=T : 2$
$Y=F : 1$     $X=F : 2$

$p(y|X_1=T)$

$= P(Y=y | X_1 = T) = \begin{cases} 1 & y=T \\ 0 & y=F \end{cases}$

$= \frac{1}{2} \times \left( -\left( \frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right) \right)$

$= 0.4056$

$H(Y|X_2) = 0.9056$

If we divide w.r.t. $X_1$ or $X_2$
what can say about the classification
and with what "certainty"? $\rightarrow$ entropy
measurement of certainty

**Entropy:** measurement of randomness of a RV

$X$ be a categorical RV, $p(x) = P(X=x)$, $\forall x \in X$

$$H(X) = -\sum_{x \in X} p(x) \log_{|X|} p(x) \qquad X = \{0,1\} \rightarrow X \text{ is a Binary}$$

$$\underbrace{\qquad\qquad\qquad}_{\mathbb{E}\left[\log_{|X|} p(X)\right]}$$

① $H(X) \geqslant 0$

$\boxed{\mathbb{E}\left[\log_{|X|} p(X)\right]}$

② $H(X) \leqslant 1 \quad \rightarrow \quad$ Jensen's ineq.

$X = \{0,1\} \rightarrow X$ is a Binary

RV

and $H$ is measured in bits.

$f$ is convex $\qquad f(x)$

$\underline{\text{Convex}}$

$$\mathbb{E}f(x) \geqslant f(\mathbb{E}X)$$

concave $\rightarrow \quad \leqslant$

Conditional Entropy : Observe $Y$, a proxy of $X$

$$-\sum_y \sum_x p(x,y) \log p(x|y) \quad = \quad H(X|Y) \qquad \text{if } X \perp\!\!\!\perp Y$$

$$\hookrightarrow P(X=x \mid Y=y) \qquad H(X|Y)$$

$$= \sum_y P(y) \left( -\sum_x p(x|y) \log p(x|y) \right) \qquad = H(X)$$

$$\underbrace{\phantom{-\sum_x p(x|y) \log p(x|y)}}_{H(X|Y=y)} \qquad = \sum_y P(y) H(X|Y=y)$$

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

mutual information

$\uparrow$ HW

# Algorithm for decision tree building

- Repeat until stopping criteria not met

   − find the feature that
     yields max information gain
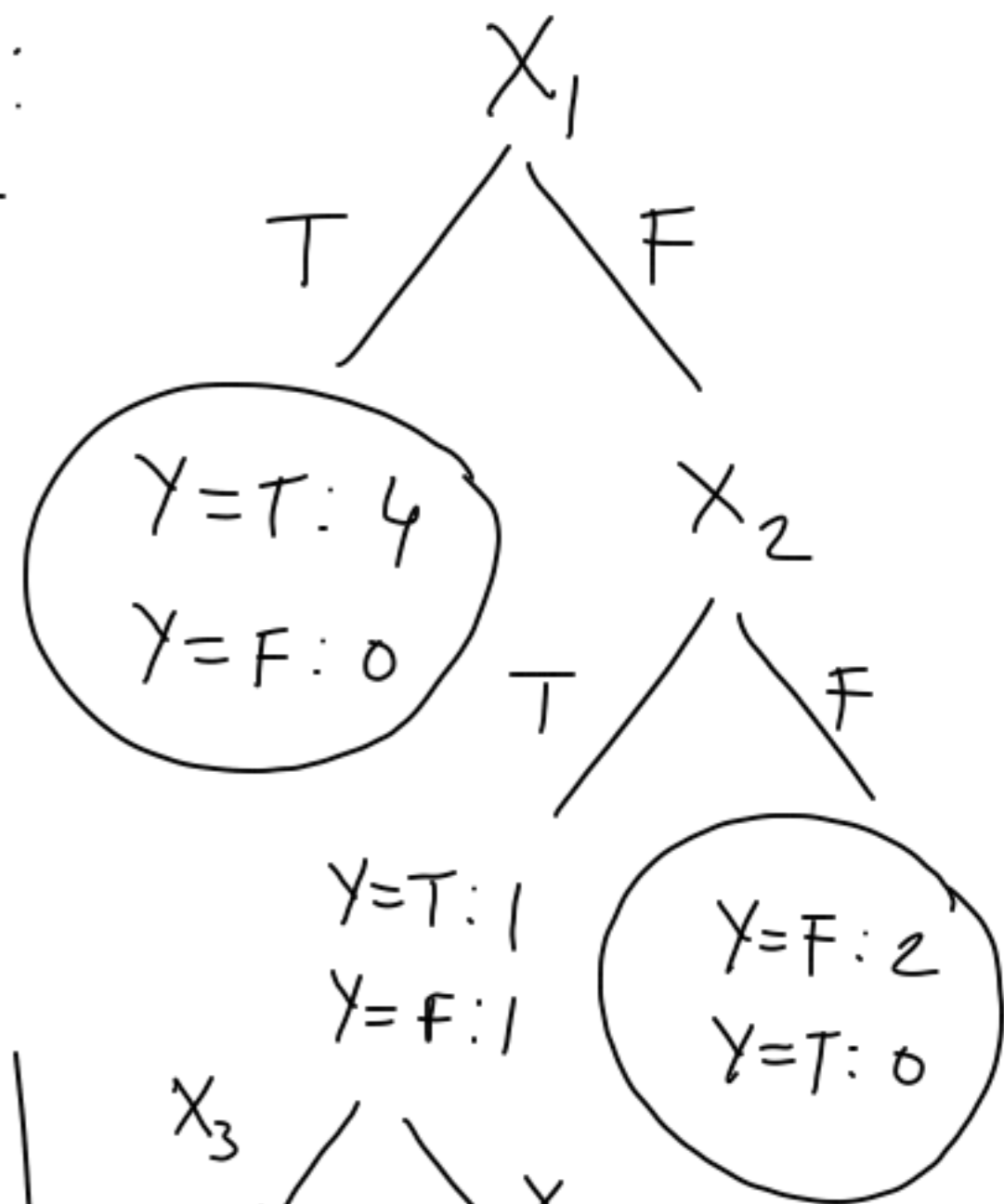     $\left(\text{min conditional entropy}\right)$

Other

**Remark:** Metric used : Gini index.

## Where to stop?

**Base case 1:**

node with
atomic distributions

$H(Y \mid node) = 0$

**Base case 2:**

When all remaining
features give
identical info
gain.



$X_1$

T / F

$Y=T: 4$
$Y=F: 0$

$X_2$

T / F

$Y=T: 1$
$Y=F: 1$

$Y=F: 2$
$Y=T: 0$

$X_3$ / $X_4$

info gain is same for all
remaining variables

Ex 4:

| $Z_1$ | $Z_2$ | $Y$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

$H(Y) = 1$

$H(Y|Z_1) = 1 = H(Y|Z_2)$

according to base case 2, this
shouldn't be split



Decision tree with root node $Z_1$. Branch 0 leads to node $Z_2$, branch 1 leads to node $Z_2$. Left $Z_2$: branch 0 → $Y=0$, branch 1 → $Y=1$. Right $Z_2$: branch 0 → $Y=1$, branch 1 → $Y=0$.

# Overfitting in decision trees

Shallow tree → not enough power to distinguish

deep tree → specific to training examples

## Three methods

○ Pre-pruning / Early stopping: hold a validation set

keep on creating the tree until test error goes up again.

' Post-pruning: allow the tree to fully grow and then reduce some of its branches.

' Ensemble method: using averages of various models

error    Test error

model complexity