# LOGISTIC REGRESSION : Weight Regularization



$2x_1 - x_2 = 0$, $W_1 = 2$, $W_2 = -1$

to the decision boundary

What happens when you Scale $w$ with a positive Scaling factor?

Whether it is $W_1 = 2$, $W_2 = -1$
$W_1 = 20$, $W_2 = -10$
$W_1 = 2 \times 10^5$, $W_2 = -1 \times 10^5$

SAME DECISION BOUNDARY

Which of these solutions is a log regression model likely to converge to?

Recall in a log. regression model:

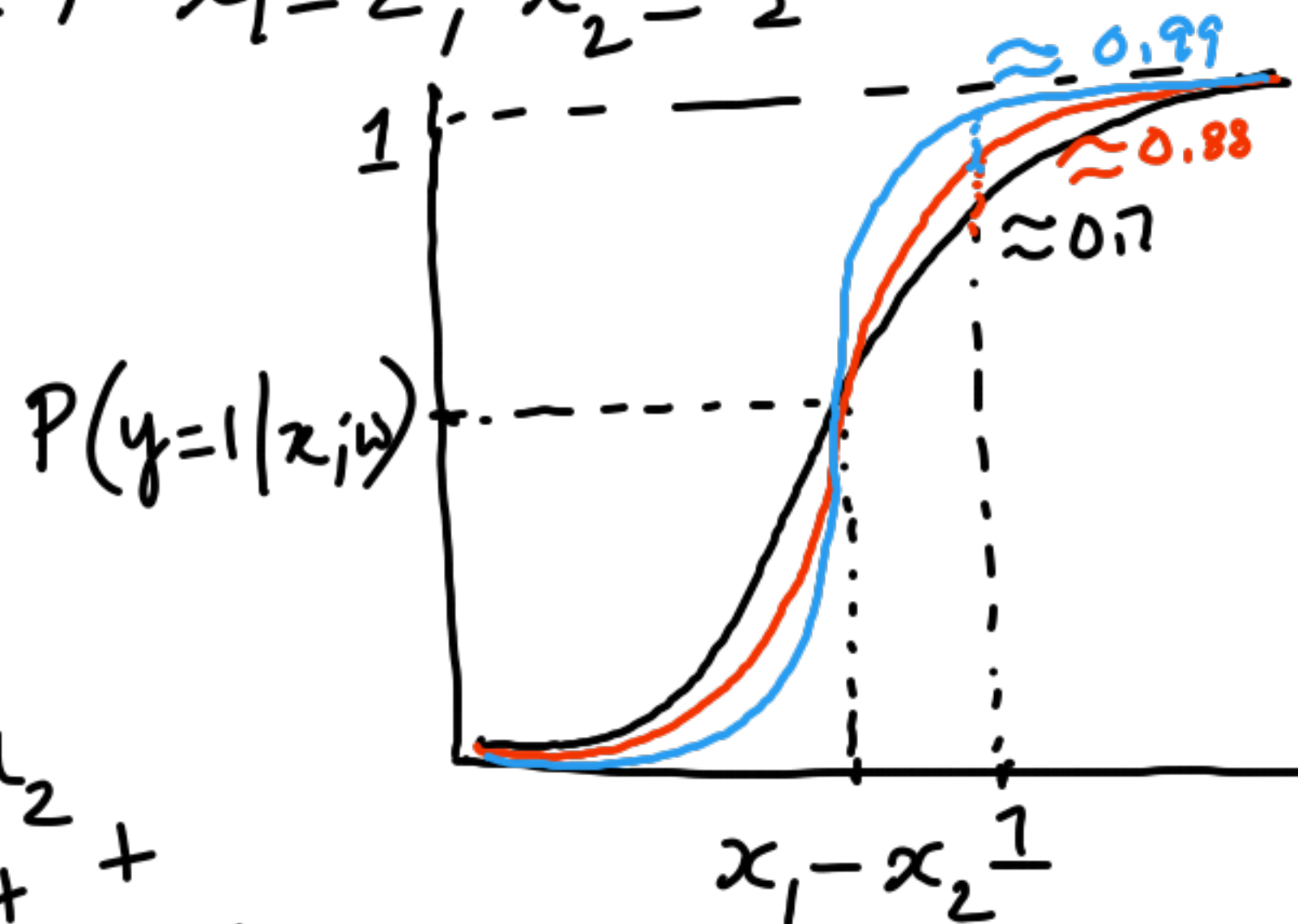$$P(y_i = 1 \mid x_i; w) = \sigma(w^T x_i) = \frac{1}{1 + e^{-w^T x_i}}$$

$$\longrightarrow Ⓐ$$

Objective function is to minimize CE loss i.e., maximize the conditional data likelihood

To maximize the prob in Ⓐ, we want $w^T x_i$ to be large (so that $e^{-w^T x_i}$ is small)

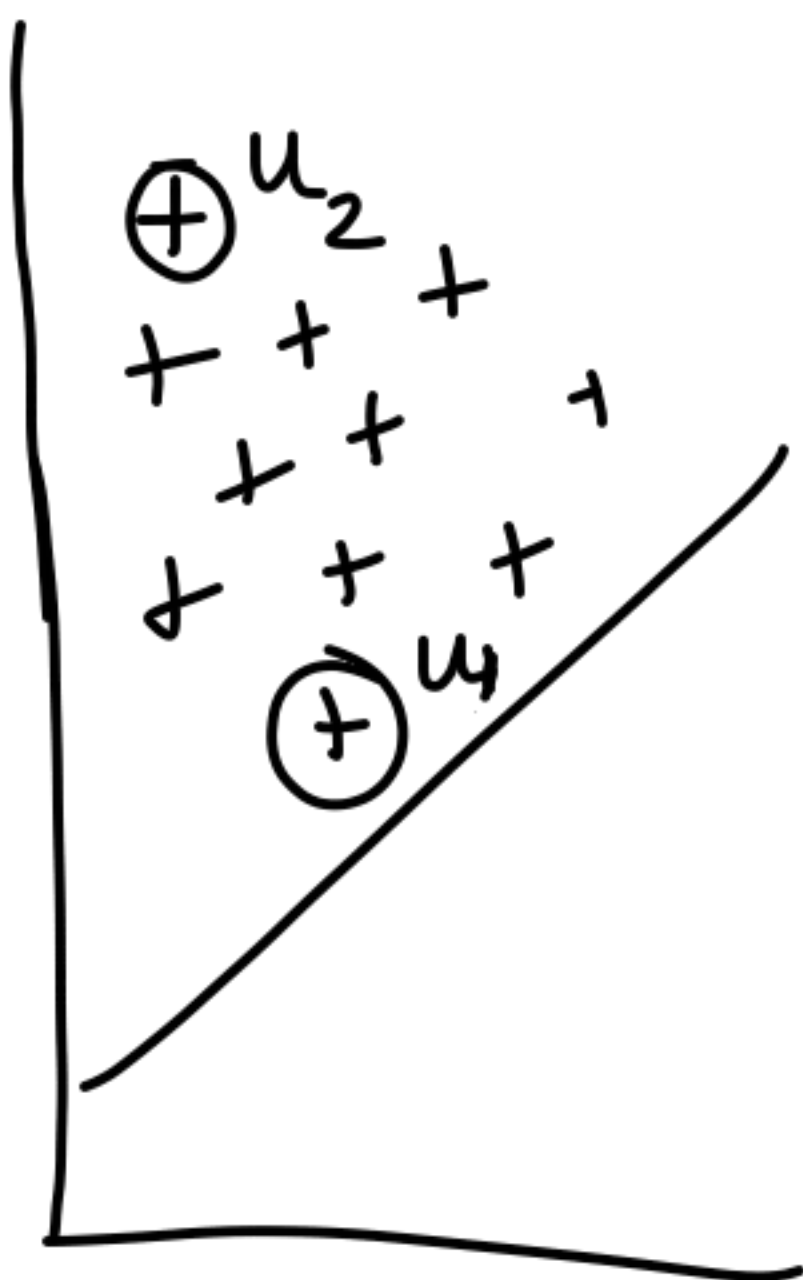$\Rightarrow$ we want $w$ values to be large

Is this a good idea?

$x$ ; $x_1 = 2$, $x_2 = 1$

≈ 0.99

≈ 0.88

1

≈ 0.7

$P(y=1|x;w)$

$x_1 - x_2$   1

$w_1 = 1$, $w_2 = 1$

$w_1 = 2$, $w_2 = 2$

$w_1 = 5$, $w_2 = 5$

⊕ $u_2$
+ +   +
+   +     +
+   +
+   +   +
⊕ $u_1$

Consider points $u_1$ and $u_2$ with the same label.

We want $w^T u_2$ and consequently the $P(y=1|u_2;w) \approx 1$

[ compared to point $u_1$ which is closer to the decision boundary ]

But with large $w$, the distinction between such points (e.g. $u_1, u_2$) becomes smaller.

Solution: Regularized logistic regression

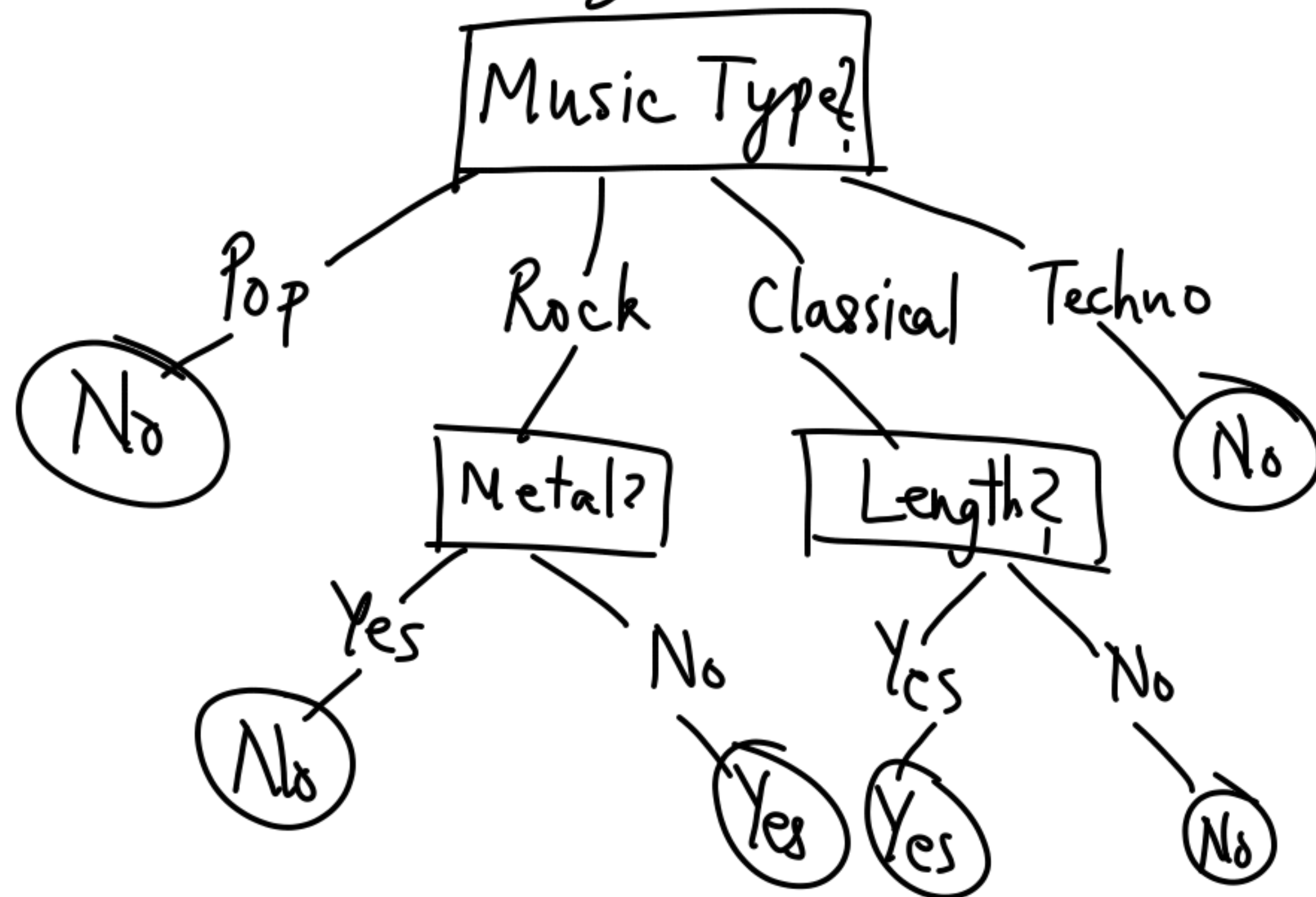$$w^*_{\text{Reg-LR}} = \underset{w}{\text{argmin}} \sum_i -\log P(y_i | x_i ; w) + \lambda \|w\|_2^2$$

$\hookrightarrow L_2$-regularized Log. Regression model

Desiderata for classification models:
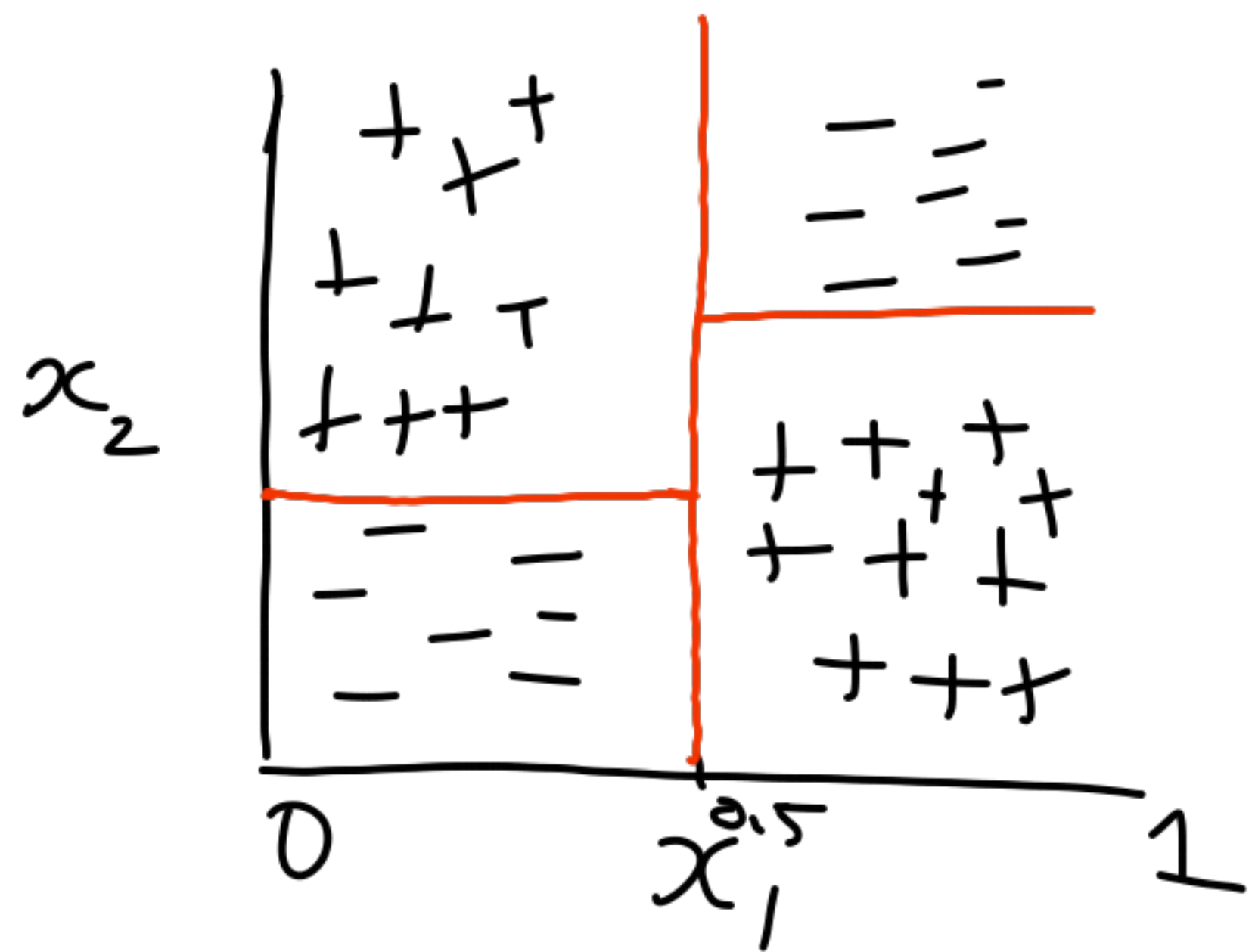① Ability to learn complex decision boundaries
② Interpretable

DECISION TREE CLASSIFIERs satisfy both criteria!

A decision tree is an _interpretable_ model whose predictions can be written as a "disjunction of conjunctions" across the attribute values taken by the training instances.
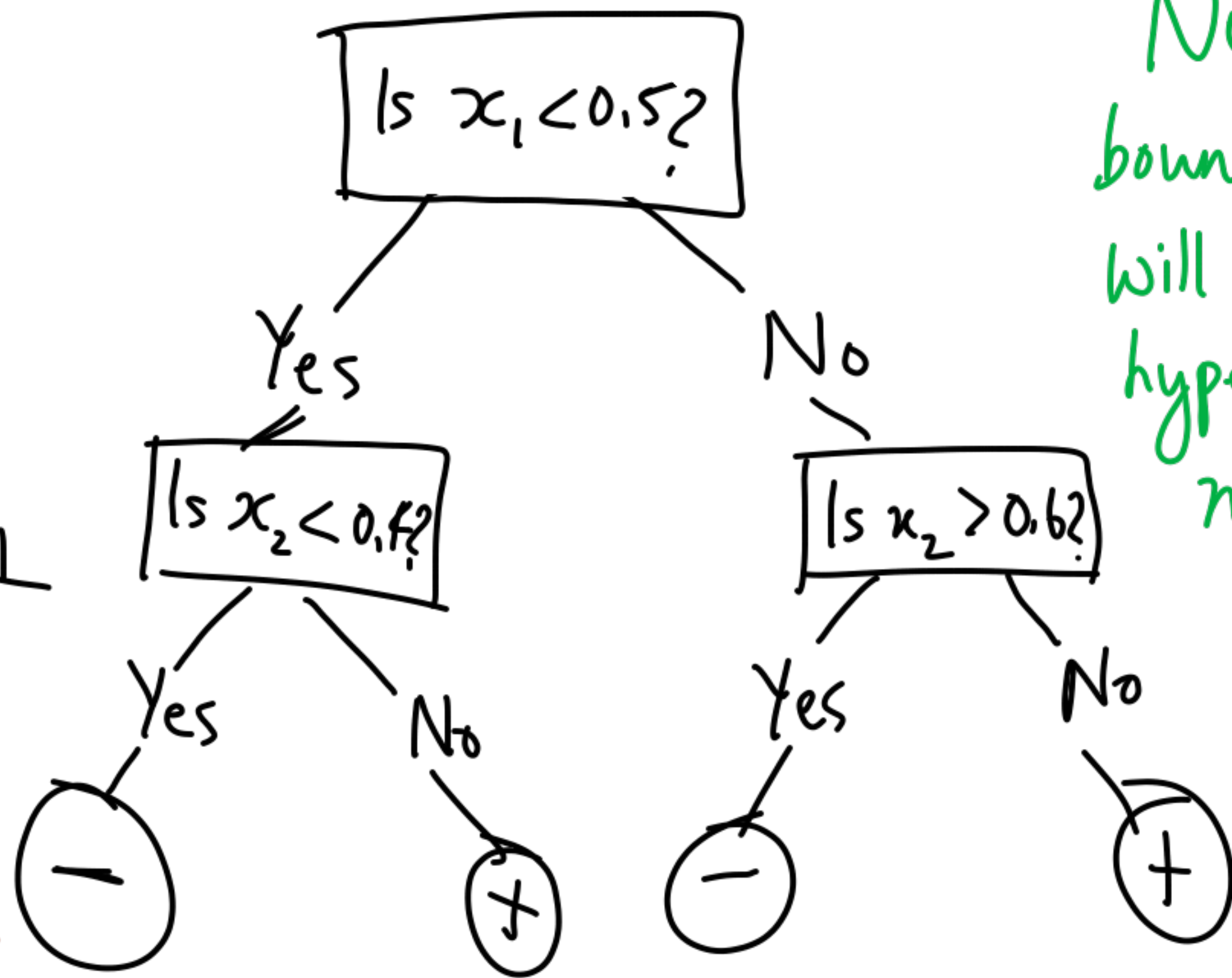
Music Type?

- Pop → No
- Rock → Metal?
  - Yes → No
  - No → Yes
- Classical → Length?
  - Yes → Yes
  - No → No
- Techno → No

Path in this DT is a conjunction of attributes;

$$\left( Type = Rock \wedge Metal = No \right)$$

∴ Tree is a _disjunction of conjunctions_

# What is the decision boundary of a DT?



Is $x_1 < 0.5$?

Yes — No

Is $x_2 < 0.1$?    Is $x_2 > 0.6$?

Yes    No    Yes    No

$-$    $+$    $-$    $+$

Decision boundary of DTs are axis-parallel hyperplanes [ if the nodes are functions of a single attribute)

Note: Decision boundary for DTs will be linear hyperplanes if the nodes are a linear combination of 2 or more attributes (e.g., $x_1 + x_2 > 0$)

## finding an optimal DT

Finding the optimal DT by optimizing an objective/loss function over the attributes is NP-hard!

DT construction is typically _greedy_. Here's the basic template.
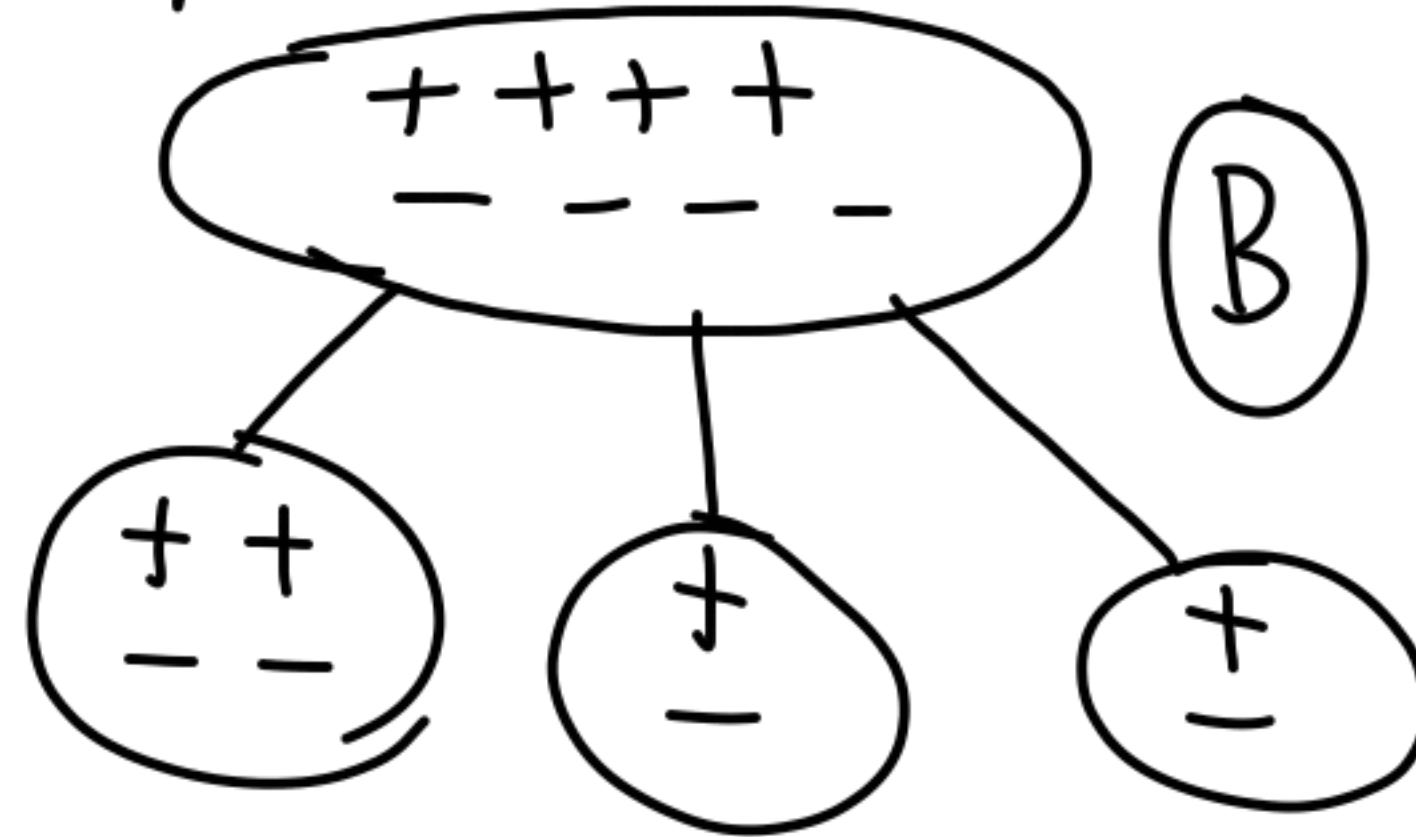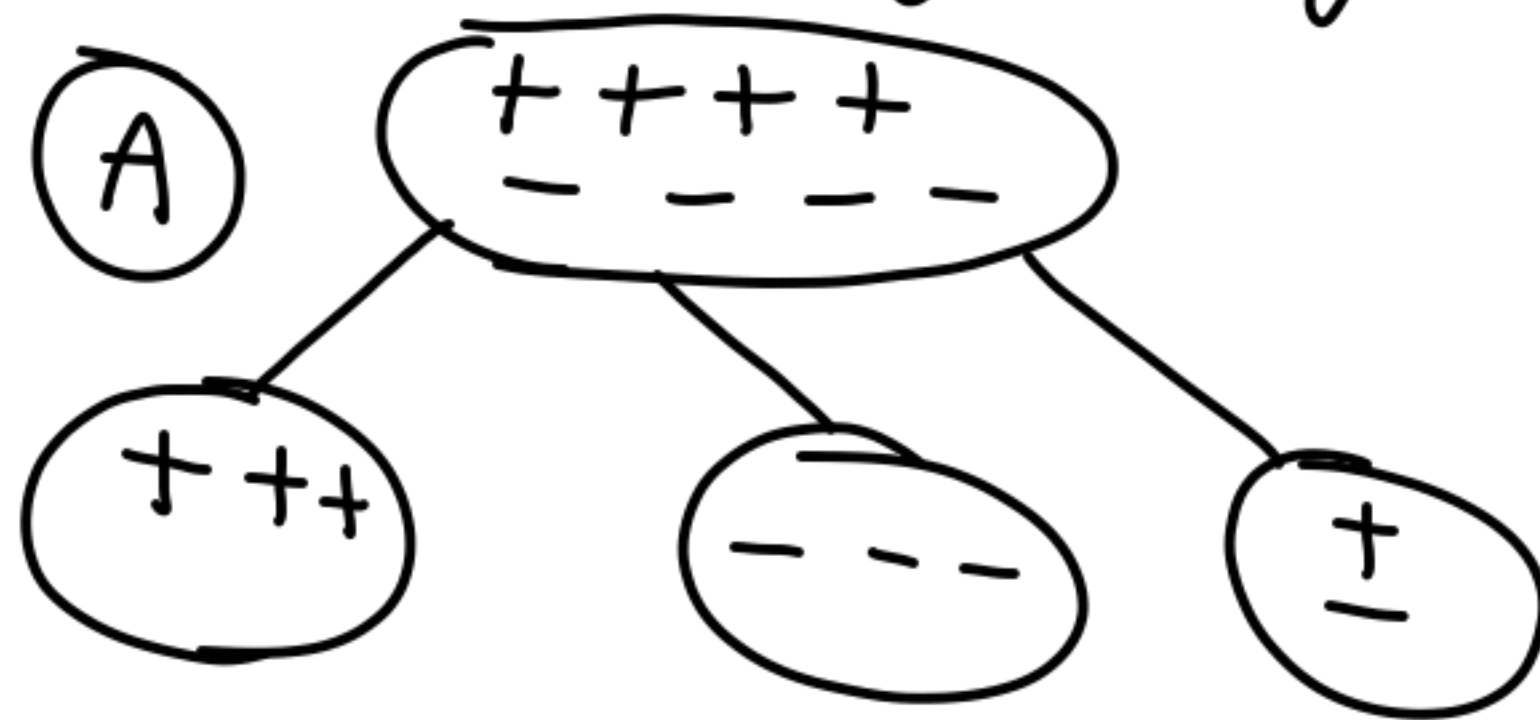
Step 1 : Start with an empty node

Step 2 : Pick the best attribute to split on

Step 3 : Repeat step 2 recursively on each node till a stopping criterion is met

Two Qs:
1. Choosing the "best" attribute to split on
2. Choosing a Stopping criterion

Q1  Picking the "best" attribute
Consider the following two splits i

A

++++
-- - - -

+++     ---     +

++++
-- - - -                    B

++      +     +
--      -     -

Which split is better?

A is better
because the two of
resulting data subsets
are homogenous in label

(A) is a better split because the tree depth is smaller compared to (B). We like smaller trees since they tend to generalize better [ lower of overfitting ].
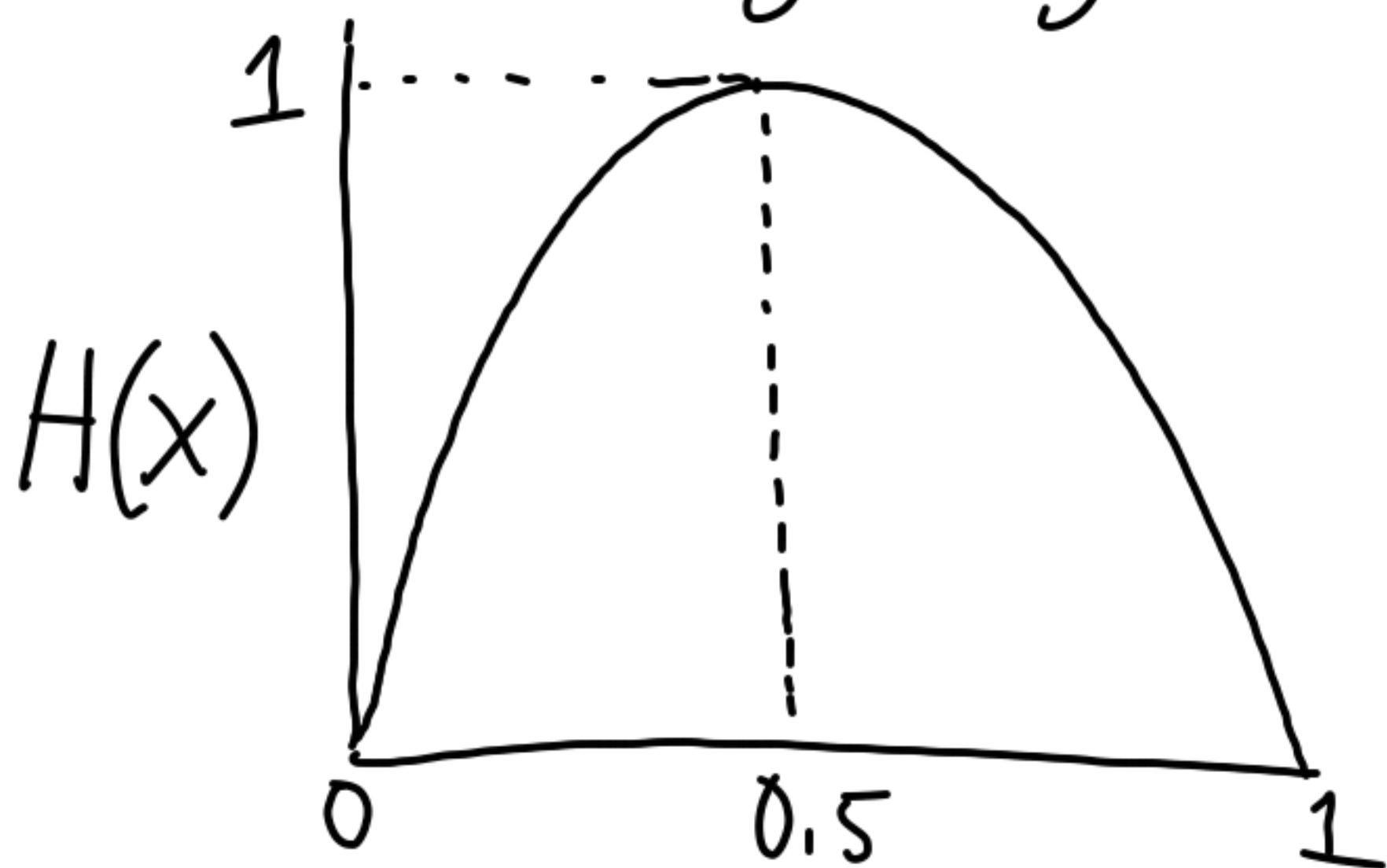
Intuition : We want splits to result in nodes that are homogenous in their labels

How do we quantify this intuition?

ENTROPY : Entropy of a random variable $X$ is a measure of uncertainty in $X$

for a random variable $X$ that takes values $x \in \mathcal{X}$, entropy $H(X)$ can be defined as:

$$H(X) = \sum_{x \in \mathcal{X}} - P(X=x) \log_{\textcircled{2}} P(X=x)$$

base 2 to measure entropy in bits

Illustrate entropy using a coin toss. $X$ is a binary r.v. taking values $\underset{(H)}{0}$ and $\underset{(T)}{1}$



High entropy indicates that the underlying distribution is nearly uniform

Low entropy indicates low uncertainty meaning there are well-defined modes in the underlying distribution

Entropy of a dataset $S$

$\Rightarrow$ Entropy of the underlying label distribution

$\Rightarrow H(S) = \sum_{i=1}^{K} -P_{S,i} \log P_{S,i} \quad \left[ K \text{ labels overall} \right]$

where $P_{S,i}$ is the probability that a random sample from $S$ will have label $i$

$\left( P_{S,i} \text{ is the relative count of } \# \text{ of instances with label } i \right)$

A good splitting criterion for DTs is "INFORMATION GAIN"

Consider an attribute "a" that can take values from $V(a)$, and a dataset $S$. Let $S_p$ be the subset of $S$ with all instances having its "a" attribute labeled as $\gamma$.

$$\text{Gain}(S, a) = H(s) - \sum_{\gamma \in V(a)} \frac{|S_\gamma|}{|S|} H(S_\gamma)$$