

Generalization = predicting power.

Built a decision tree: how well will it do on future data?

Lead us to PAC model.

What is the "true error" or generalization error of a classifier?

- Decision trees: Fix  $T$ ,  
probability distribution  $D$  on new examples.

$(x, y)$   
↑ challenge    ↑ label

$$\Pr[T(x) \neq y] \text{ of } T$$

$(x, y) \sim D$

True error /  
Generalization error.

Consider:

$x^1$	$y^1$
$\vdots$	$\vdots$
$x^m$	$y^m$

$S$  = training set.

Learner is given  $S$ .

Learner memorizes  $S$ , and hands it back.

On new example, it just outputs 0.

Now consider:

$x$	$y$
$x^1$	$y^1$
$\vdots$	$\vdots$
$x^m$	$y^m$

$S =$

training set

Each

$x^i \in \{0, 1\}^n$ , e.g. 011010 if  $n = 6$ .

$y^i \in \{0, 1\}$

Let's build a decision tree.

Assume  $x^i$  are distinct.

You can build a decision tree (size  $\geq |S|$ )  
that is consistent with all the points in  $S$ .

Question : How well does the tree generalize?  
what is the true error of this tree?

How can we estimate the true error of  
a classifier (in this case a decision tree)?

use a "hold-out" or a "validation set", call it "H".

$S$  = training set.

1) use  $S$  to build a decision tree.

$H$  = hold-out set.

2) Estimate tree's true error via  
its error on  $H$ .

Cross-validation : a way to reuse the data that  
you've held-out into a hold-out set.

---

Another approach :

trade-off training error with "model complexity."

Define another potential function  $\phi$

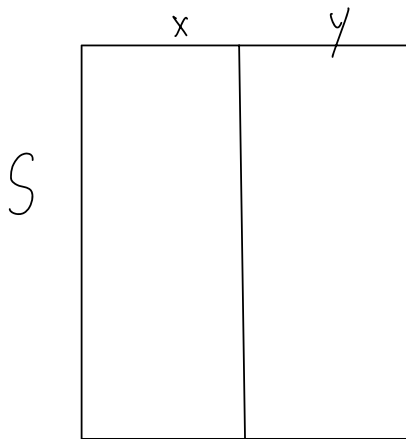
$\phi$  : trees  $\rightarrow \mathbb{R}$ . Given a training set  $S$

$$\phi(T) = \underbrace{\text{training error on } S}_{\text{hyperparameter,}} + \alpha \cdot \underbrace{\text{size}(T)}_{\substack{\text{size of} \\ \text{training} \\ \text{set.}}}$$

Goal : minimize  $\phi$  :

- min training error
- min size of tree
- increase  $\alpha$

Another approach: MDL / Minimum Description Length principle.



Some # of bits needed to encode  $S$ .

upper bound  $m \cdot (n + 1)$  label.

↑ training examples

↑ features

Build a tree  $T$ :

Let's say  $T$  is correct on 90% of  $S$  and incorrect on 10%.

We can encode  $S$  using  $\#$  bits ( $T$ ) and  $\#$  bits to encode that 10% we got wrong.

This captures the notion of compression.