# k-nn and decision trees

k-nn and decision trees don't use GP.
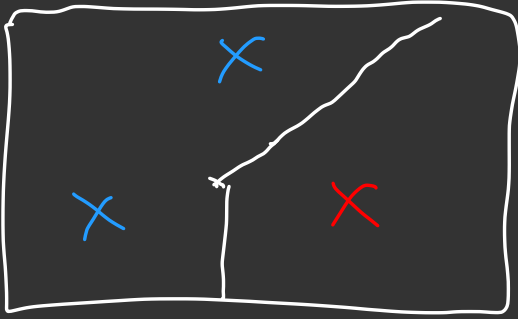
## 1-nn:

① pick a distance function

② memorize training set

③ output closet point

$$p(x, x_i) = \min_j p(x, x_j)$$

## k-nn

① pick a distance function and integer $k \geq 1$

② memorize training set

③ classification : output plurety labels among $k$ nois.

regression : average .....

1-nn
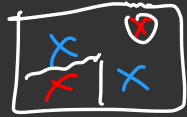
3nn: all red
└─┘
majority

Remarks:

① If $(x_i)_{i=1}^n$ distinct, 1-nn gets
→ False

    o    traing error,



    ② k-nn    may fail to get o training error.

    ③   why   k-nn.

       ↓

① higher  k  | smooth |  predictor. with less complex model.
                              (Those random blue labels
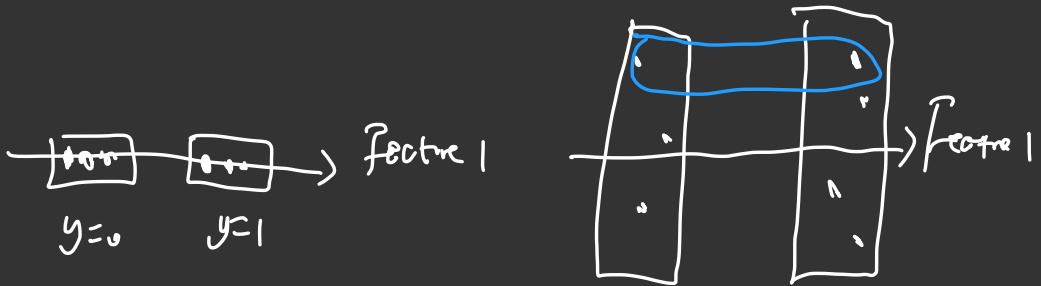                                    won't bother)



② carefully  chose  k   $o(\ln n)$

Test error of knn with L2 distance.

| K | 1 | 3 | 5 | 7 | 9 |
|---|---|---|---|---|---|
| err | 0.0309 | 0.0295 | 0.0312 | 0.0306 | 0.0341 |

Test error of 1-nn with diff distances

| Distance | L2 | L3 | Tangent | Shape |
|---|---|---|---|---|
| err | 3.09 | 2.83 | 1.10% | 0.63% |

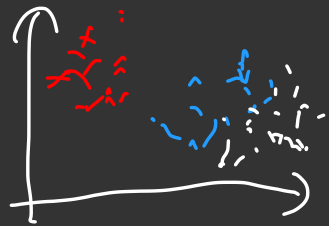k-nn  can be [broken] by bad features.



$y=0$    $y=1$    → feature 1

→ feature 1

Curse of dimension: Given poly $(d)$ random unit norm points in $\mathbb{R}^d$, with prob $> 99\%$, each is $2 \pm O(1/\sqrt{d})$ from all others.

# Decision trees

① binary tree which partitions input space.

② each tree node is associated a splitting rule.

③ leaf node associated with label $y$.

# Training decision trees



① pick uncertainty measure:

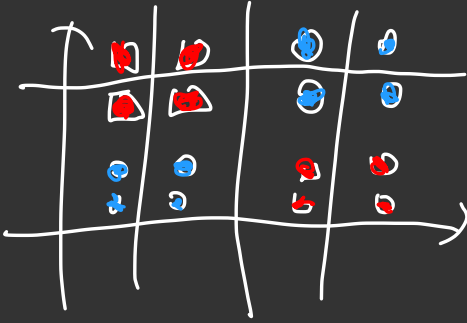$$u(T) = \frac{1}{n} \sum_{\text{leaf } S \in T} \frac{1}{|S|} u(S)$$

② put all pts at root

③ Loop (threshold)

    ① pick leaf $l$ and splitting rule $h$ that maximally reduces uncertainty

    ② split data in $l$ using $h$ and grow trees accordingly,
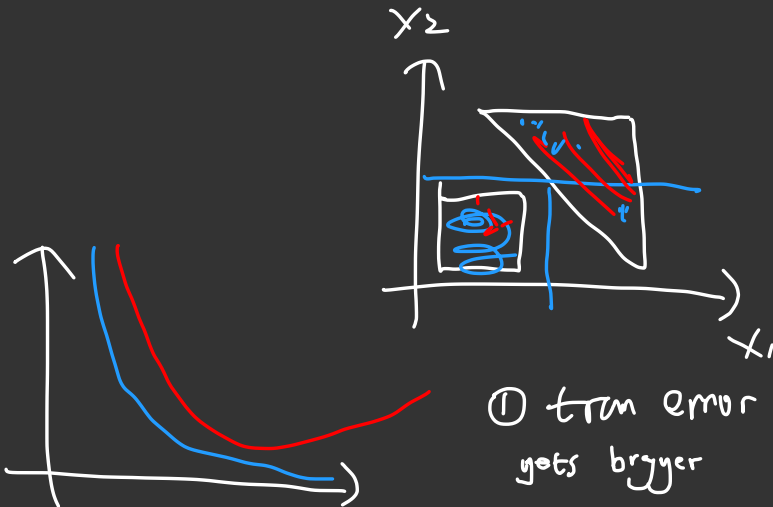
# The greedy algorithm may fail



Those vertical lines
fail half

## When to stop

① tree reaches size

② when every leaf is pure ( overfitting )



$X_2$

$X_1$

① train error → 0 when tree
gets bigger

② test error. decreases at first
but increases overfitting

Summary

|         | k-nn                              | Decision trees.                    |
|---------|-----------------------------------|------------------------------------|
| Train   | memorize data                     | greedy split                       |
| Test    | find k closest memorized pts      | traverse tree output knt label     |
| overfit | vary k                            | limit tree size.                   |