

CS 45780

ML Setup, Recap:

1. Data: $D = (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

x 's: Input instances $x \in \mathbb{R}^d$

y 's: Corresponding output $y \in \mathcal{C}$

Binary classification or $\mathcal{C} = \{0, 1\}$ or $\{-1, 1\}$, multi class $\mathcal{C} = \{1, 2, \dots, K\}$, regression $\mathcal{C} = \mathbb{R}$

2. \mathcal{H} : Set of models or hypotheses

each model $h \in \mathcal{H}$, given an input instance x , outputs $h(x)$

3. l : loss function measures performance of a model

eg: 0-1 loss: $l(h(x), y) = \mathbb{1}\{h(x) \neq y\}$

Squared loss: $l(h(x), y) = (h(x) - y)^2$

absolute loss: $l(h(x), y) = |h(x) - y|$

Goal of supervised learning:

Given data D find model h such that loss $l(h(x), y)$ on future instances (x, y) is small.

Asking for model that minimizes loss $l(h(x), y)$ for all possible future (x, y) is too much.

Why?

Future instances generated from some mechanism (often represented by some) distribution P written as

$$(x, y) \sim P$$

Formal Goal of supervised learning:

Given Data D find model h that minimizes

Risk or population loss: $E_{(x,y) \sim P} l(h(x), y)$
generalization loss

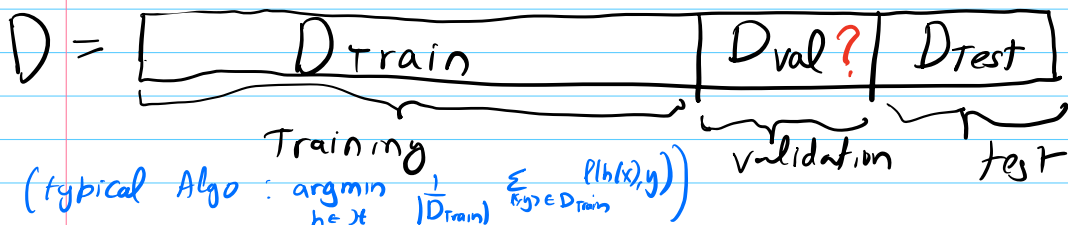
Learning Algo. is the procedure that tries to attain above goal.

What is a good proxy for

$$E_{(x,y) \sim P} [l(h(x), y)] ? \text{ (and why?)}$$

Test loss: For evaluation $\frac{1}{|D_{\text{test}}|} \sum_{(x,y) \in D_{\text{test}}} l(h(x), y)$

Split Data:



No free Lunch theorem:

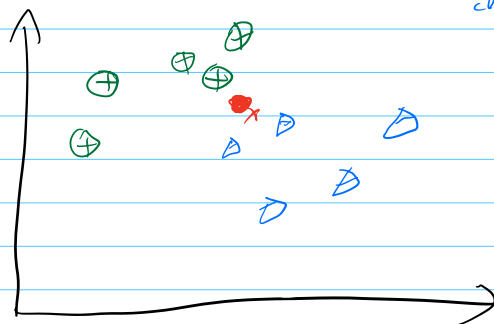
“Every ML Algorithm makes assumptions.”

CS 4/5780

K-Nearest Neighbors Classifier

Assumption: Similar points are likely to share same label

Classification Rule: For a test point x , assign the most common label amongst the k most similar (closest) training instances



Formally: Given $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ and test point x
Find $S_x \subseteq D$ st. $|S_x| = k$ and $\forall (x', y') \in D \setminus S_x,$

$$\text{dist}(x', x) \geq \max_{(x'', y'') \in S_x} \text{dist}(x'', x) \quad h(x) = \text{MODE}(y' : (x', y') \in S_x)$$

Pro-tip: In case of a tie, reduce k by 1 and repeat

Common distances: Minkowski distance or l_p distances

$$\text{dist}(x, x') = \left(\sum_{i=1}^d |x[i] - x'[i]|^p \right)^{1/p}$$

$p=2$ is Euclidean distance

QUIZ: What distances do $p=1$, $p \rightarrow \infty$ and $p \rightarrow 0$ correspond to.

Bayes Optimal Classifier

If we knew P how well could we do?

$$\text{Bayes Error} : \min_{\text{All possible } h} E_{(x,y) \sim D} \ell(h(x), y)$$

If you knew $P(y|x)$, Given point $x \in \mathcal{R}$, optimal classifier:

$$h_{opt}(x) = \operatorname{argmax}_{y \in \mathcal{C}} P(y=y|x=x)$$

$$\text{Bayes Error}(x) = 1 - \max_{y \in \mathcal{C}} P(y=y|x=x)$$

This is the Best we can do!

Quiz: Coin has probability p of heads.

1. If tossed twice, what is the probability q , of two different outcomes?
2. Conclude that $q \leq 2(1-p)$

1-NN classifier: Simplest case $k=1$

$$\text{Risk of 1-NN} \leq 2 \text{ Bayes ERROR}$$

Formal proof is involved, see Cover & Hart '67

Intuition:

1. Say P was a discrete distribution on a finite set of points. Then, as $n \rightarrow \infty$ every test point has already occurred in D_{train} . (Say we pick any one of previous occurrences as the nearest neighbor)
2. Risk of 1-NN Classifier is now given by the quiz question. Why?

We are asking the question, what is the probability that, label y of a new test instance x matches that of a randomly chosen training point $x_i \in D_{\text{train}}$ such that, $x_i = x$. Its label y_i is drawn independently from $P(y|x=x)$

$$\text{Hence, Risk of 1-NN} \leq 2(1 - \max_{y \in \mathcal{C}} P(y=y|x=x))$$

Claim: Given x , let $\hat{x} \in D_{\text{TRAIN}}$ be the 1-NN of x in D_T
 as $n \rightarrow \infty$, $\text{dist}(x, \hat{x}) \rightarrow 0$, $\hat{x} \rightarrow x$

Risk of 1-NN ≤ 2 Bayes ERROR

RAIN

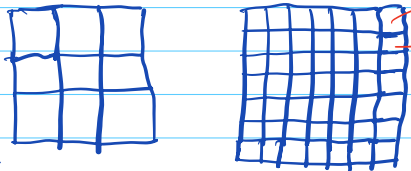
K-NN For general $K > 1$

a. Larger K , as $n \rightarrow \infty$,

Intuition: For point x , we predict as majority of K draws from $P(Y|X=x)$

b. But, if K grows too fast, the more we rely on farther points to predict label for x

The Curse of Dimensionality:



1. Each sub-cube has same area

2. if we randomly pick a few we will pick more cubes near surface!

1. Assume points are drawn uniformly at random from a unit hypercube $[0,1]^d$

2. Hypercube of volume k/n within unit hypercube is expected to contain k out of n points

3. Length l of such cube is given by $l^d = k/n$

What does this mean?

$k=10$
 $n=1000$ \rightarrow

d	2	10	100	1000
l	0.1	0.63	0.955	0.9954

$l = \left(\frac{k}{n}\right)^{1/d}$