

Decision

程序代写代做 CS 编程辅导

pickling up



time:

S_w

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

0 0 0	x x
* 0 0	5 5 x
x x x	0 0 0
x x x	0 0 0

For 3 ~~https://tutorcs.com~~ leaves any instance
that fall in them can be automatically
labeled without need to do $NN \rightarrow$ all 3
of them are pure

Advantage \rightarrow for those 3 leaves

you don't need to store the \Rightarrow

coordinate for any of the instances, since you skip doing NN for any instance that fails



that take further : for S_{NE} , $\Pr(\bar{x}) = \frac{5}{8}$ and $\Pr("0") = \frac{3}{8}$

So, in this case we may say that a new instance in that leaf is labeled " x " with $\Pr = \frac{5}{8}$

Assignment Project Exam Help

~~OR~~

We could continue to split S_{NE} until we reach pure leaves

Email: tutorcs@163.com

QQ: 749389476

Question: Is that always possible?

Answer: in the general case of noise (two instances with some \vec{x} have differently)

No.

Adding noise to resolve label conflict

Counterintuitively, we add a negligible amount of

zero-mean Gaussian noise ($\sigma^2 = 1 \rightarrow$ to all \vec{x}).
→ as illustrated as conflict resolves we can
split until



Without this you can split data until you reach pure leaves.

WeChat: cstutorcs

→ you have to split until you have one instance per leaf → BAD because one instance

per leaf means you are overfitting the training data (high variance)
QQ: 749389476

<https://tutorcs.com>

1 if tree is too deep (many layers and leaves)
then we have high variance (overfitting, high test error)

2 if tree is too shallow (low depth), then we have high bias (underfitting, high error on both training set and test set)

New Goal Build a tree that is maximally compact
and only has ~~more~~ leaves.



This is a hard problem.

(No solution can be found in polynomial-time algorithms)
WeChat: cstutorcs

Assignment Project Exam Help
Fortunately, we can find suitable solution using a
greedy algorithm

Email: tutorcs@163.com

QQ: 749389476
greedy algorithm = any algorithm that chooses to
make local optimum choice at each stage, without
considering the global optimum.

We keep splitting the data to minimize an impurity
function that measures purity among subsets

How do we decide borders for splits?
程序代写代做 CS 编程辅导
we test which borders lead to greater purity using

impurity functions



- Entropy

- Gini impurity index

Impurity Functions WeChat: cstutorcs

We want to measure ambiguity of labels

within a leaf. Email: tutorcs@163.com

Ex A leaf has 100 "x" and 1 "o" \rightarrow low

ambiguity, so label a new instance "x"
<https://tutorcs.com>

Ex A leaf has 500 "x" and 500 "o" \rightarrow high ambiguity, so unsure how to label a new instance

Given a data set : $S = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$

$y_i \in \{1, 2, \dots, c\}$, $c = \text{number of classes} (\text{distinct labels})$

Let $S_k \subseteq S$ where $S_k = \{(x, y) \in S \text{ such that } y=k\}$

↓
subset



$\Rightarrow S$ is partitioned $\Rightarrow S = S_1 \cup S_2 \cup \dots \cup S_c$

↓
union

WeChat: cstutorcs $S_k \cap S_j = \emptyset$ if $k \neq j$

Assignment Project Exam Help

↓
intersection empty set

Def Gini impurity function of set S

$= G(S)$
QQ: 749389476

$$= \sum_{k=1}^C p_k (1 - p_k)$$

$$p_k = \frac{|S_k|}{|S|} = \frac{\text{number of instances with label } k}{\text{number of instances in } S}$$

Ex Case of $C=2$ labels

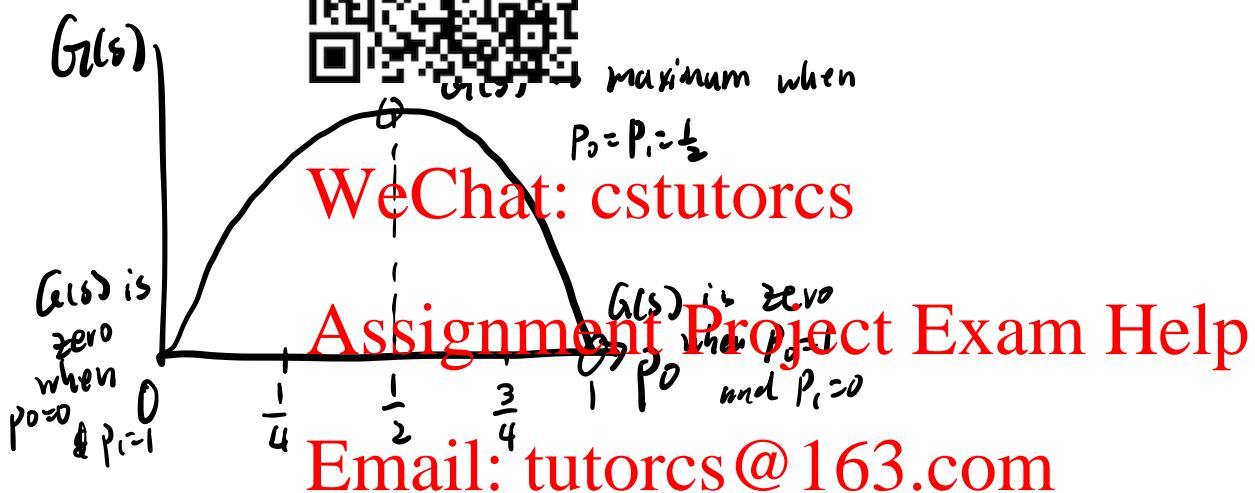
$$\Pr("1") = p_1$$

$$\Pr("0") = p_0 = 1 - p_1$$

$$\Rightarrow G(s) = P_0(1-P_0) + P_1(1-P_1)$$

程序代写代做CS编程辅导

$$= P_0(-P_0) P_0 \\ = 2$$



We want $G(\text{leads})$ to be low (at best, zero)

At worst $G(\text{leads}) = 0.5$ (large) when $P_0 = P_1 = \frac{1}{2}$

\Rightarrow this tells us nothing about a new instance falling in that leaf #

In general, the worst case scenario is when the labels are uniformly distributed

$$\Pr(\text{label } k) = \frac{1}{C} \quad \forall k = 1, 2, \dots, C$$

$$= g_k$$

$\{q_k\}$ 程序代写代做 CS 编程辅导

fancy



Ideally, we want $\sum_{k=1}^C p_k q_k = 0$

our goal is that to decide about a split,

we select $\min_{\{S, P\}} G(S)$
Assignment Project Exam Help

Email: ~~tutorcs~~@163.com

at least new leaf we want

QQ: ~~749389476~~ have a dominant label

$$p_k \gg p_j, \forall j \neq k$$

<https://tutorcs.com>

To measure our success, we use a tool called the Kullback-Leibler (KL) divergence

Def The Kullback-Leibler (KL) divergence between two distributions $\{p_k\}$ and $\{q_k\}$ is $KL(p||q) = \sum_{k=1}^C p_k \log \left(\frac{p_k}{q_k} \right)$

this can be
 $\begin{cases} \log(\cdot) & \text{for } a \rightarrow 0^+ \\ \log_a(\cdot) & \text{for } a > 0 \\ \ln(\cdot) & \text{for } a = e \end{cases}$

Ex] $kL(p \parallel q)$ 程序代写代做 CS 编程辅导
etc



$$\log(1) \\ := 0$$

WeChat: cstutorcs

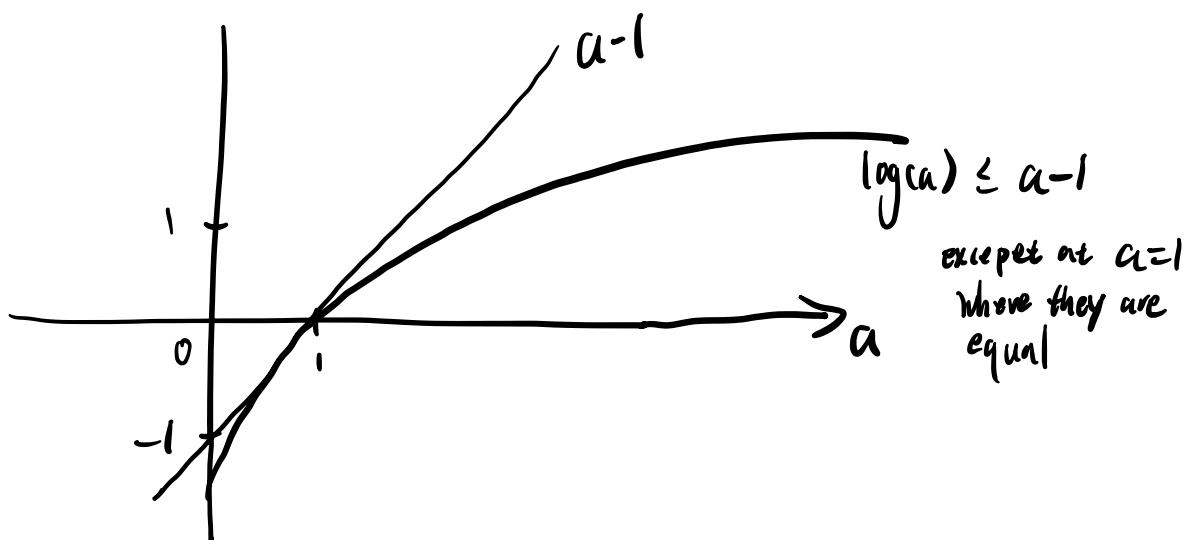
property of kL :

- kL is ~~NOT~~ a distance measure, since

$$kL(p \parallel q) \neq kL(q \parallel p)$$

$$kL(p \parallel q) \geq 0 \quad \text{with equality} \Leftrightarrow p = q$$

proof <https://tutorcscs.com>



$$-KL(p||q) = \sum_{k=1}^C p_k \log\left(\frac{p_k}{q_k}\right) = \sum_{k=1}^C p_k \log\left(\frac{q_k}{p_k}\right)$$



$$\left(\frac{q_k}{p_k} - 1 \right) = \sum_{k=1}^C (q_k - p_k)$$

$$= \sum_{k=1}^C q_k - \sum_{k=1}^C p_k = 1 - 1 = 0$$

WeChat: cstutorcs

$$= 1 = 1$$

Assignment Project Exam Help

$$\Rightarrow -KL(p||q) \leq 0$$

Email: tutorcs@163.com

with equality if $\frac{p_k}{q_k} = 1 \forall k$

QQ: 749389476

$$\Rightarrow KL(p||q) \geq 0$$

<https://tutorcs.com>

equality if both distributions are identical

$p_k \neq q_k \Rightarrow KL(p||q) > 0$

$p_k = q_k \Rightarrow KL(p||q) = 0$

For a given leaf, we measure the quality of purity as the KL divergence from the worst case,

q_{ik} (uniform) 程序代写代做 CS 编程辅导
Want to join?

$$H(p) = -\sum_{k=1}^c p_k \log\left(\frac{p_k}{q_k}\right) = \sum_{k=1}^c p_k \log\left(\frac{p_k}{\frac{1}{c}}\right)$$


↓
 $q_k = \frac{1}{c} \forall k$

$$= \sum_{j=1}^c p_k \log(c_j) = \sum_{k=1}^c p_k \log(c_k) = \log(c) \underbrace{\sum_{j=1}^c p_j}_{=1}$$

$$= \sum_{k=1}^c p_k \underbrace{\text{Assignment Project Exam Help}}$$

Email: tutorcs@163.com
Independent of our

QQ: [749389476](https://tutorcs.com)
Choice of split

In deciding about a split in a leaf, our goal is to
maximize $H(p)$ which is the same as

↓
uniform

$$\text{minimizing } - \sum_{k=1}^c p_k \log(p_k)$$

Def the ENTROPY of a set with probability
distribution $p = \{p_1, p_2, \dots, p_c\}$

is designed as 程序代写代做 CS 编程辅导

$$H(S) = - \sum_{k=1}^C p_k \log(p_k)$$



- * Our goal for every split the leaves have $\min H(S) \{S_1, S_2\}$
WeChat: cstutorcs

Impurity Assignment Project Exam Help

Gini Impurity of a tree	Email: tutorcs@163.com	Entropy of a tree
We seek the split that results in minimum G(tree)	QQ: 749389476	" $H(\text{tree})$
Step 1: choose a split. It partitions S . For this, you'll need to isolate feature (dimension) X_i and sort all instances along dimension i .	Step 1: same	" "
Step 2: For each of the two subsets, estimate Gini impurity index $P_{w,x} = \frac{ S_{w,x} }{ S_w }, P_{w,0} = \frac{ S_{w,0} }{ S_w } \Rightarrow G(S_w)$	Step 2: same, except we P's to get entropy " " " "	

and

$$P_{E,x} = \frac{|S_{E,x}|}{|S_E|} P_{E,0} = \frac{|S_{E,0}|}{|S_E|}$$



Step 3: $G(\text{tree}) =$
Weighted average optimism index

WeChat: cstutorcs

Step 3:

$$H(\text{tree}) = \frac{|S_w|}{|S|} H(S_w) + \frac{|S_E|}{|S|} H(S_E)$$

↑
number of instances
in original before split

Step 4: try all features

(dimensions) and all possible

$N-1$ splits

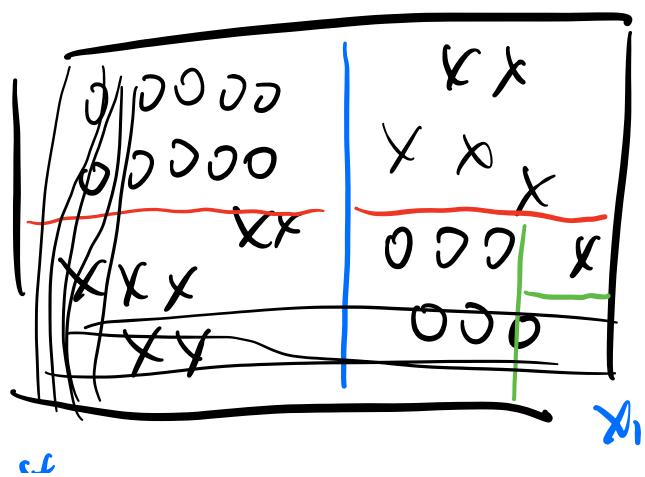
Email: tutorcs@163.com

Pick the split that
minimizes impurity.

QQ: 749389476

<https://tutorcs.com>

Ex 5



1st split along X
程序代写代做CS编程辅导

2nd split. Yields better entropy purity



in E

3rd split yields purity in all leaves

keep splitting until entropy (or Gini) is zero
WeChat: cstutorcs

try every single coordinate split, all $N-1$ of them in V_1 ,
then go to V_2 and try all $N-1$ coordinate splits

Email: tutorcs@163.com

For each coordinate,

QQ: 749389476

$$\text{calculate } \begin{cases} P_{N,0} = \frac{3}{4} \\ P_{N,X} = \frac{1}{4} \end{cases} \} H(S_N)$$

$$\text{calculate } \begin{cases} P_{S,0} = \frac{2}{3} \\ P_{S,X} = \frac{1}{3} \end{cases} \} H(S_S) \} H(\text{tree})$$

$$= \frac{4}{7} G(S_N) + \frac{3}{7} G(S_S) \quad (\text{or you could use } H \text{ instead of } G)$$

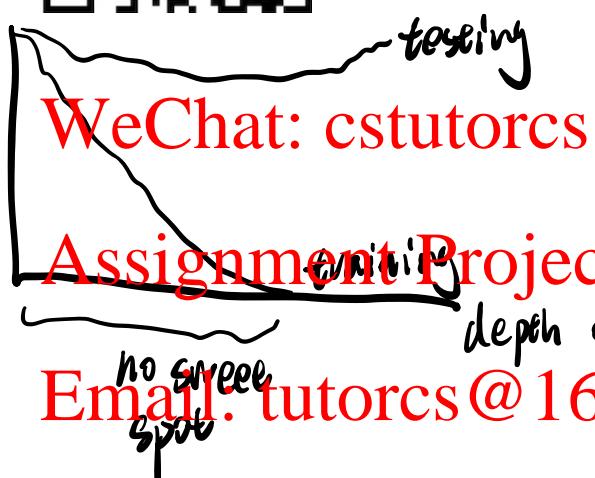
Danger of too few splits: high bias (underfitting)

"many" : high variance (overfitting)

When tree is too shallow \Rightarrow you'll see both training and testing error are high



you'll see training error low
and testing error high



Assignment Project Exam Help

Email: tutorcs@163.com

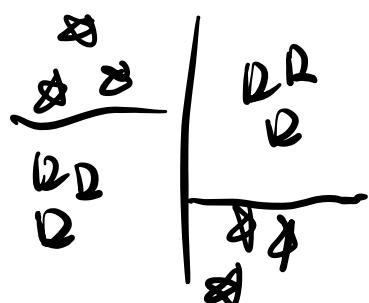
What if I go through splits along X;

but entropy not improve?

Do I stop?

Answer is NO.

Ex XDR function, binary addition



1st split does not
improve the entropy

2nd split round 2 get purity

NOTE: Decision Trees are prone to high variance
 (why we use bagging)

Classification Trees (CART)

The algorithm is the same, but instead of using entropy or Gini impurity, we use average square loss

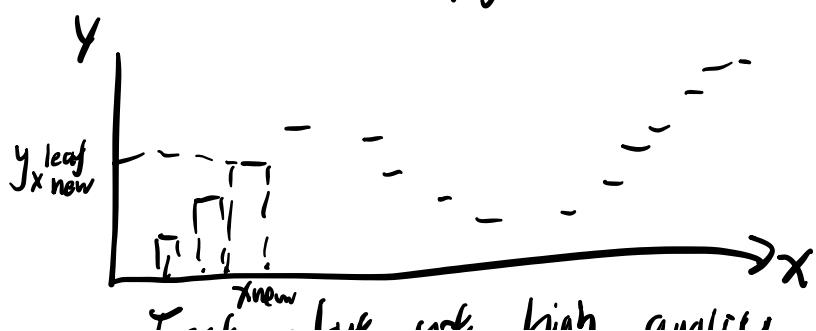
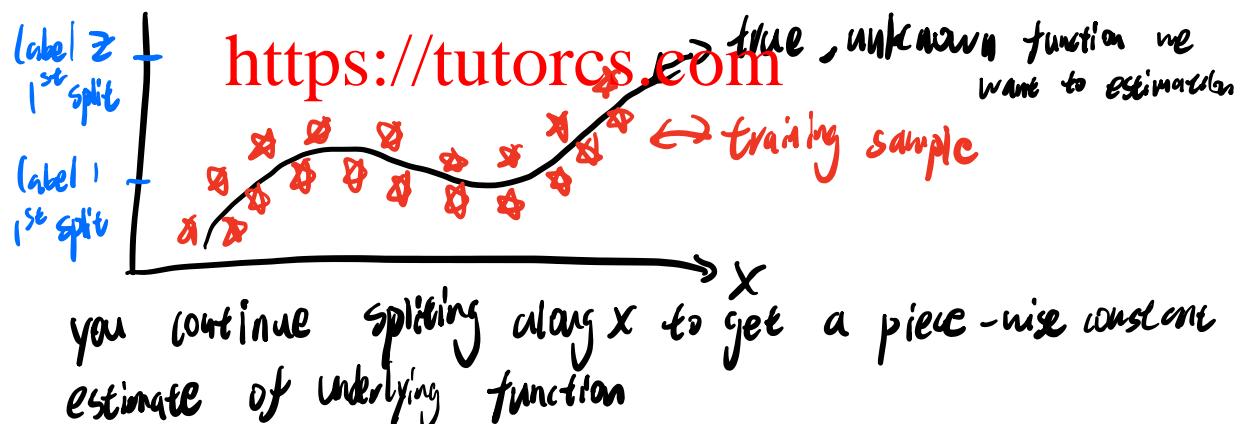
$$L(s) = \frac{1}{|S|} \sum_{(x_i, y_i) \in S} (y_i - \bar{y})^2 \quad \begin{matrix} \rightarrow \\ \text{average square distance} \\ \text{from label} \end{matrix}$$

WeChat: <https://tutorcs.com>

Assignment Project Exam Help

Email: tutorcs@163.com

Ex we took [QQ:749389476](https://tutorcs.com) to split the "x" feature space into 2



fast , our view "ij" quickly

程序代写代做 CS 编程辅导

Bagging (Bootstrap Aggregating)

Used to reduce variance in machine learning algorithms

We'll try to make $\bar{h}(x)$ closer to $h(x)$

by using the weak Law of Large Numbers

(the sample average of independent and identically distributed sample converges to the expectation or mean)

$$\frac{1}{N} \sum_{i=1}^N x_i \xrightarrow[N \rightarrow \infty]{\text{Email: tutorcs@163.com}} E(x) = \bar{x}$$

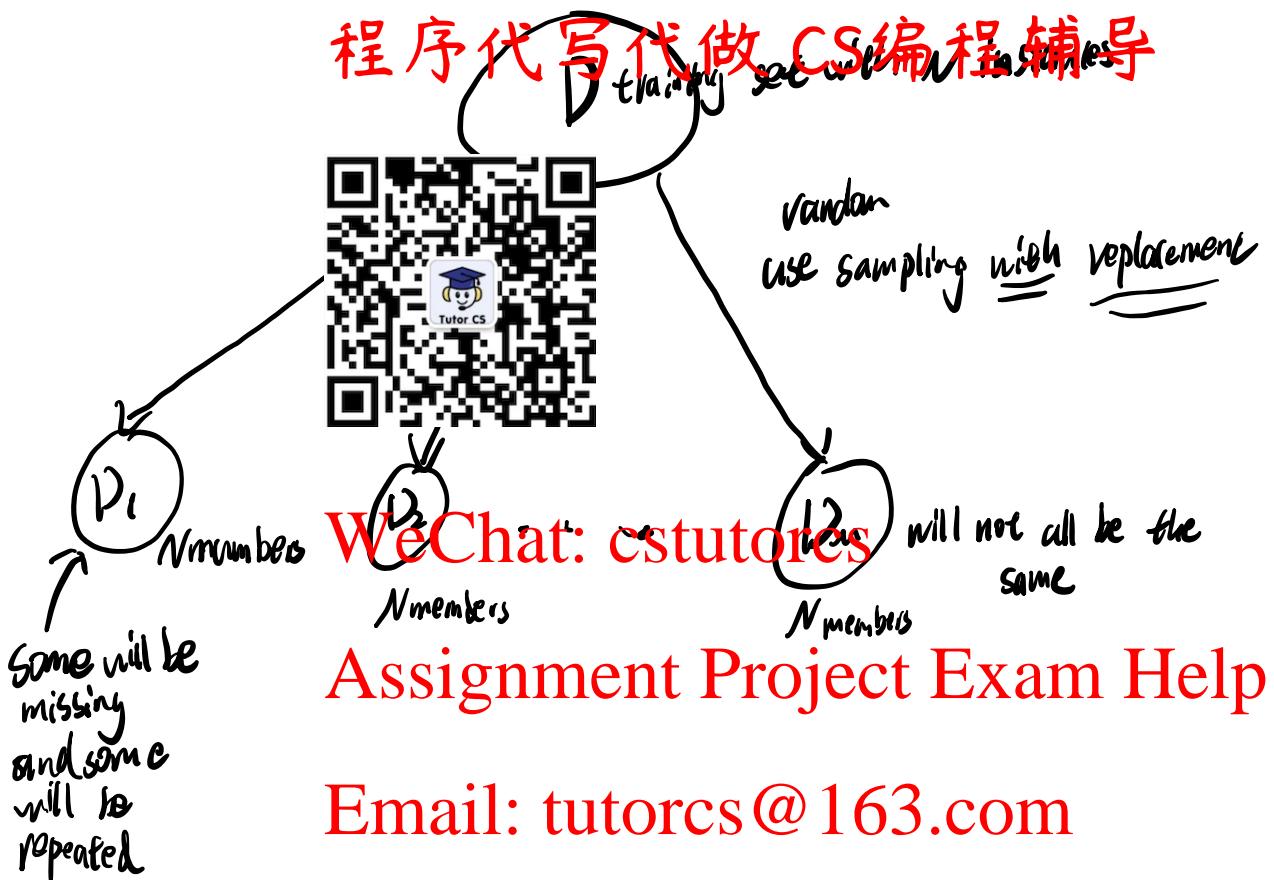
QQ: 749389476

If I have M data sets, then I can train M trees, each one based on sets $D_i, i=1 \dots M$

Each yields a classifier $h_{D_i}(x)$

Next, average them for regression $h_{avg}(x) = \frac{1}{M} \sum_{i=1}^M h_{D_i}(x) \xrightarrow[M \rightarrow \infty]{\text{https://tutorcs.com}} h(x)$

Ok, if you are working with classifiers, then take the mode among all $h_{D_i}(x)$

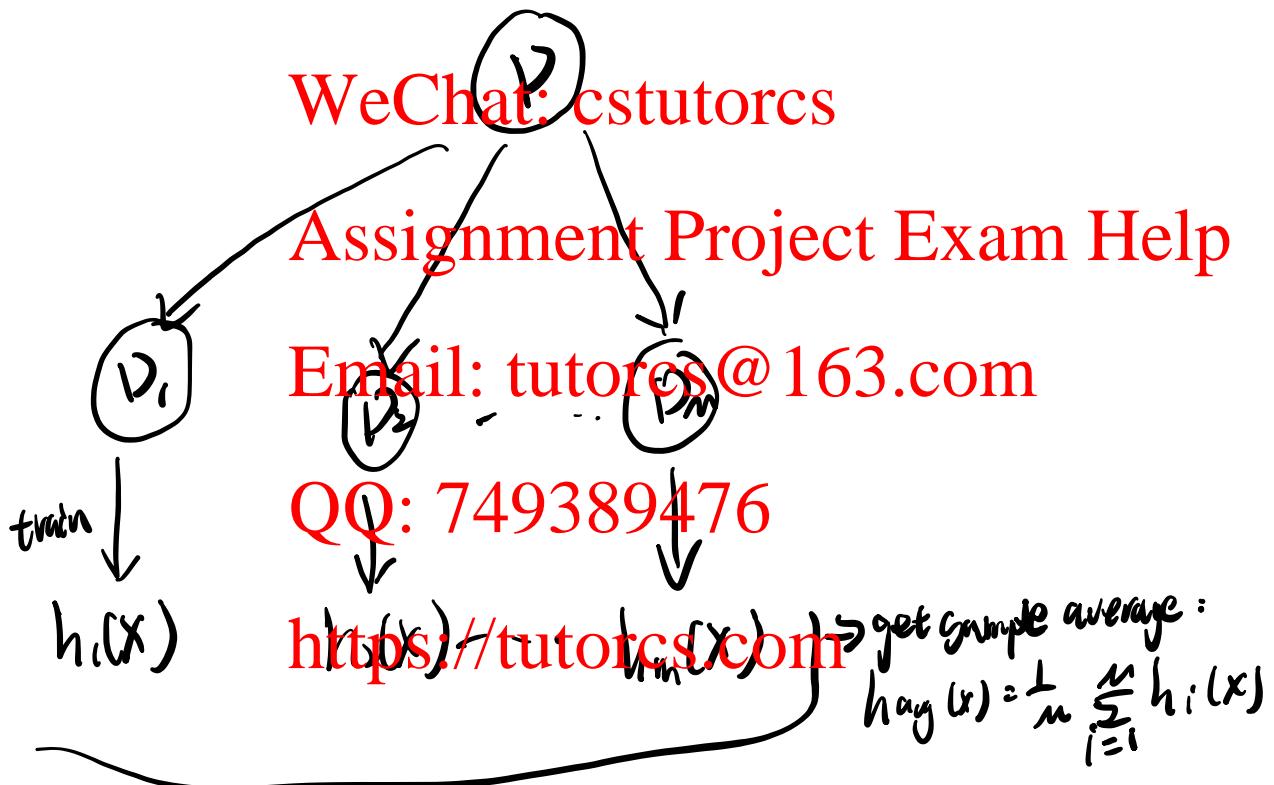


* Because we ~~QQ 749389476~~ sample with replacement
 randomizing which instances we pick from $S \dots$.
 some instances get repeated, and some get skipped
 This mitigates overfitting

Def "Bootstrapping" = resampling technique that involves repeatedly drawing samples from our source data with replacement.

we'll try to bring $h_{avg}(x)$ closer to $\bar{h}(x)$
by using the weak Law of Large Numbers

(the average of independent and identically
distributed random variables converges to the expectation or
mean)



Note, there gets one NOT independent \Rightarrow Weak Law
of Large Numbers does not apply
and sample average will not converge to \bar{h}

But bagging effectively improves Variance!