

May 0

Bias - Variance Decomposition

(A) One fixed $\langle X, y \rangle$, one \vec{x}^* ^{r.v.} $y = f(x) + \delta$
 $\Rightarrow \text{MSE}(\vec{x}^*) = \sigma^2 + (\text{Bias}[g(\vec{x}^*)])^2$

where $\text{Bias}(g(\vec{x}^*)) = (g(\vec{x}^*) - f(\vec{x}^*))$

(B) One fixed X , but multiple y , ^{r.v. now} one \vec{x}^*
 $\Rightarrow \text{MSE}(\vec{x}^*) = \sigma^2 + \text{Bias}[g(\vec{x}^*)]^2 + \text{Var}[g(\vec{x}^*)]$

where $\text{Bias}(g(\vec{x}^*)) = \underbrace{E[g(\vec{x}^*)]}_{\text{mean model}} - f(\vec{x}^*)$

$\text{Var}[g(\vec{x}^*)] = E[g(\vec{x}^*)^2] - E[g(\vec{x}^*)]^2$

(C) Multiple X , multiple y , multiple \vec{x}^*

$\Rightarrow \text{MSE} = \sigma^2 + E_x[\text{Bias}[g(\vec{x}^*)]^2] + E_x[\text{Var}[g(\vec{x}^*)]]$

General
Error

Trees have no bias, but high variance as a model. (Average of trees converges to actual f as number of trees $\rightarrow \infty$)

May 1

Bagging intro

Bagging (Breiman, 1994)

Imagine many models g_1, \dots, g_n which are averaged

$$g_{\text{avg}} = \frac{g_1 + g_2 + \dots + g_n}{M}$$

What is MSE?

$$\text{MSE} = \sigma^2 + E_x [\text{Bias}[g_{\text{avg}}]^2] + E_x [\text{Var}[g_{\text{avg}}]]$$

↓ won't write E_x (but assume is there)

$$= \sigma^2 + E \left[\frac{g_1 + \dots + g_n}{M} - f \right]^2 + \text{Var} \left[\frac{g_1 + \dots + g_n}{M} \right]$$

$$= \sigma^2 + E \left[\frac{(g_1 - f) + (g_2 - f) + \dots + (g_n - f)}{M} \right]^2 + \frac{1}{M^2} \text{Var}[g_1 + \dots + g_n]$$

$$= \sigma^2 + \frac{1}{M^2} (\text{Bias}(g_1) + \dots + \text{Bias}(g_n))^2 + \frac{1}{M^2} (\text{Var}(g_1) + \dots + \text{Var}(g_n))$$

① Assume all biases are same, i.e.

$$= \sigma^2 + \frac{1}{M^2} (M \cdot \text{Bias}(g_n))^2 + \frac{1}{M^2} \sum \text{Var}[g_m]$$

$$= \sigma^2 + \text{Bias}(g_n) + \frac{1}{M^2} \sum \text{Var}[g_m]$$

② $\text{Bias}[g_n] \approx 0$ (for each n)

i.e. \mathcal{A} for g_n is overfit

$$\Rightarrow \text{MSE} = \sigma^2 + \frac{1}{M^2} \sum \text{Var}[g_m]$$

Assume

③ Models g_1, \dots, g_m are independent and have same variance

$$\text{MSE} = \sigma^2 + \frac{\text{Var}(g_m)}{M} \quad (\text{m's canceled})$$

This is a cheat: only have indep. models if all from different data

If X, Y are indep

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$$

$$MSE = \sigma^2 + \frac{\text{Var}[g_m]}{M}$$

if all models use different data, then as $M \rightarrow \infty$, requires we have infinite data

IV

Let $M \rightarrow \infty$

$MSE = \sigma^2$, the theoretical limit fine, but this only works if infinite data

Can you do this with one ID with n observations?

No, but you can do something close

Let $ID_{(1)}$ be a n -size sample with replacement of the observations in ID . This is a "non-parametric bootstrap sample"

This sample has $\frac{2}{3}$ of the rows of ID and $\frac{1}{3}$ duplicates.

Let $ID_{(2)}$ be another independently drawn bootstrap sample from ID .

Let $ID_{(3)}$ be another.

Let $ID_{(4)}$ be another.

\vdots

Let $ID_{(M)}$ be another.

Let $g_1 = \mathcal{A}(ID_{(1)})$

$g_2 = \mathcal{A}(ID_{(2)})$

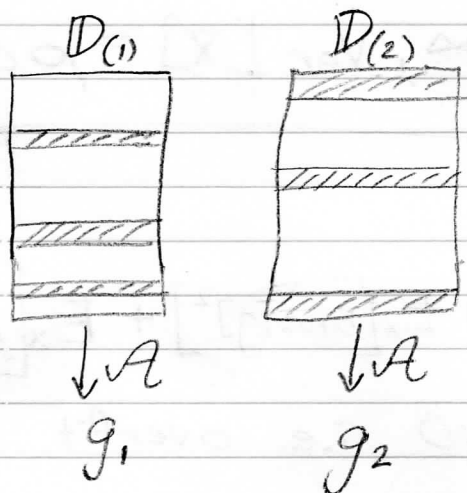
\vdots
 $g_n = \mathcal{A}(ID_{(M)})$

take average of the models (at every value of x)

and $g_{\text{avg}} = \frac{g_1 + \dots + g_M}{M}$

This is called a "bootstrap aggregator" or "bagging"

Doing this is "Meta-Algorithm"



Are g_1 and g_2 indep.?

No, but they are between fully indep. and fully dependent.

for arbitrary X_1, \dots, X_n r.v.'s

$$\text{Var}\left[\frac{X_1 + \dots + X_n}{n}\right] = \frac{1}{n^2} \left(\text{Var}[X_1] + \dots + \text{Var}[X_n] + \sum_{i \neq j} \text{Cov}[X_i, X_j] \right)$$

Assume/Use

① $\sigma^2 = \text{Var}[X_i] \forall i$

② $\sigma_{ij} = \text{Cov}[X_i, X_j] \forall i \neq j$

let $\rho = \frac{\sigma_{ij}}{\sigma^2}$ (ρ is unitless and scaled $\in [-1, 1]$)
 ↑
 correlation

so

$$\text{Var}\left[\frac{X_1 + \dots + X_n}{n}\right] = \frac{1}{n^2} \left(n\sigma^2 + n(n-1)\sigma_{ij} \right)$$

cancel n

$$= \frac{1}{n} \left(n\sigma^2 + (n-1)\sigma_{ij} \right) \Big|_{\sigma_{ij} = \sigma^2 \rho}$$

$$= \frac{1}{n} \left(\sigma^2 + (n-1)\sigma^2 \rho \right)$$

$$\text{Var}[\bar{X}] = \frac{1}{n} (\sigma^2 + n\sigma^2 \rho)$$

$$\text{Var}[\bar{X}] = \frac{\sigma^2}{n} + \sigma^2 \rho - \frac{\sigma^2 \rho}{n}$$

$$\text{Var}[\bar{X}] = \rho \sigma^2 + \frac{1-\rho}{n} \sigma^2$$

If $n \rightarrow \infty$, $\text{Var}[\bar{X}] = \rho \sigma^2$

Use g_{avg}

$$\text{MSE} = \sigma^2 + E_x[\text{Bias}[g]^2] + E_x\left[\rho \text{Var}[g] \cdot \frac{1-\rho}{M} \text{Var}[g]\right]$$

① $\text{Bias}[g] \approx 0$ i.e. overfit

② $\text{Var}[g]$ is same for all models

③ $M \rightarrow \infty$

get

$$\text{MSE} = \sigma^2 + E_x[\rho \text{Var}[g]]$$

$\rho < 1$, which is definitely the case during bagging

compare \uparrow without bagging, would be

$$\text{MSE} = \sigma^2 + E_x[\text{Var}[g]] \quad (\text{as } M \rightarrow \infty)$$

by
Breiman

Random Forests

At each iteration during regression tree algorithm, it tests every single $x_j < x_k$ rule at which there are $(n-1)p$

Breiman imagined a tree algorithm such that you only test

$$\{j_1, j_2, \dots, j_{n+n}\} \subset \{1, 2, \dots, p\}$$

covariates randomly chosen at each iteration

(trees less accurate)

so... bias increases, but still negligible
i.e. bias still ≈ 0

Thankfully, the trees are less correlated
(better savings)

$$MSE = \sigma^2 + E(\rho \text{Var}[g])$$

↑
decreased
this over
just doing Bagging

less
error
overall!

This is called

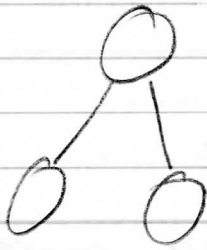
Averaging these kinds of trees
is called Random Forests

only
use some
of the
features
to make
each
tree

Regular Trees vs. Random Forest

→ check all possible splits for all features

→ check all possible split for random selection of features



$$x_1 < x_{11}$$

$$\vdots$$

$$x_1 < x_{1(n-1)}$$

$$x_2 < x_{21}$$

$$\vdots$$

$$x_2 < x_{2(n-1)}$$

$$\vdots$$

$$x_n < x_{n1}$$

$$\vdots$$

$$x < x_{p(n-1)}$$

$$p = 10$$

use p_m only

$$p_m = 5$$

example of possible set of 5 features selected

$$is \{1, 3, 4, 8, 10\}$$

$$\begin{matrix} \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ x_1 & x_3 & x_4 & x_8 & x_{10} \end{matrix}$$

used to make the trees

all $p(n-1)$ possible rules

only $p_m(n-1)$ possible rules

more bagging ... validation

How do we validate g_1 ?

$D = \underbrace{D_{(1)}}_{\substack{\text{also written} \\ D \setminus D_{(1)}}} \cup \underbrace{(D - D_{(1)})}_{\substack{= D_{\text{test}} \\ \text{about } \frac{1}{3} \text{ of } D}}$

to validate prediction $D - D_{(1)}$

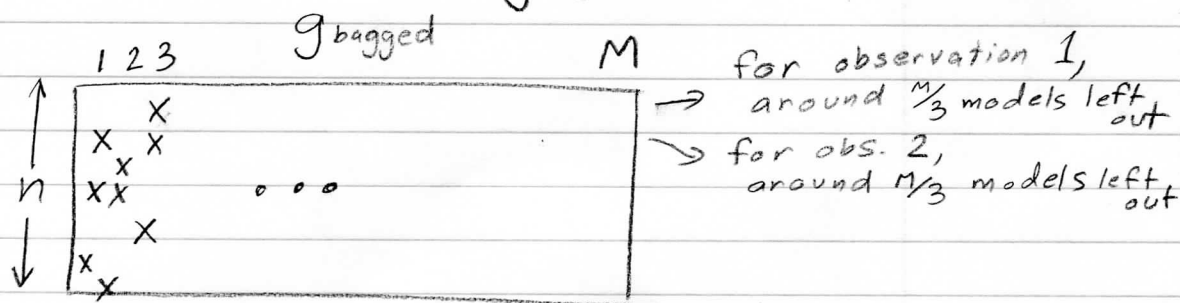
How do we validate g_2 ?

$D = \underbrace{D_{(2)}}_{\substack{\text{we use} \\ \frac{1}{3} \\ \text{is actually} \\ \text{more like} \\ \frac{1}{e}}} \cup \underbrace{(D - D_{(2)})}_{= D_{\text{test}} \text{ (about } \frac{1}{3} \text{ of } D)}$

use $D - D_{(2)}$ to validate prediction

How to validate g_{avg} ?

"hat problem"



If I predict for \vec{x}_1 , on the models left out and average, I get oos \hat{y}_1
 similarly for $\vec{x}_2, \vec{x}_3, \dots, \vec{x}_n$

for each obs.
Models left out^v are called "out of bag" models

If do
bagging,
get
validation
for
free

If I predict for \vec{x}_1 on these models
and average, I get oos \hat{y}_1 ("oob" out of bag estimate)

\vdots
If I predict for \vec{x}_2 on these models
and average, I get oos \hat{y}_2

\vdots

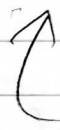
\hat{y}_3

\vdots
 \hat{y}_n

" " " for \vec{x}_n "

doing
this
is
like
 $k=2$
split

Average the residuals for all i ;
(each residual is $e_i = y_i - \hat{y}_i$)



Random forest does same thing,
but with modified trees.