

Bagging (bootstrap aggregating)

THs: Louis: Series 8, 9, 10

Claude: 1, 5, 6, 7

Peter: 2, 3, 4

• Data:  $(X_1, Y_1), \dots, (X_n, Y_n) \sim P$  $X_i \in \mathbb{R}^p, Y_i \in \mathbb{R}$ • Base procedure:  $\hat{g}: \mathbb{R}^p \rightarrow \mathbb{R}$ 

eg. regression tree

• Take  $B$  bootstrap samples  $(X_1^{*b}, Y_1^{*b}), \dots, (X_n^{*b}, Y_n^{*b}), b=1, \dots, B$ • Then,  $\hat{g}^{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{g}^{*b}(x)$ , where  $\hat{g}^{*b}(\cdot)$  is the estimate based on the  $b^{\text{th}}$  bootstrap sample.• Note that  $\hat{g}^{\text{bag}}(\cdot) \rightarrow E^*[g^*(\cdot)]$  as  $B \rightarrow \infty \approx E_P[g^*(\cdot)]$ (Ref Slides Fig 8.11. (bagging)) If  $\exists$  learners (even mediocre ones) then agg. performance improves substantially)• Idea: Mimic sampling from the distribution  $P$ , and then aggregate estimates  
! \* Bagging does not do anything for linear predictors;  
linear predictors stay the same under bagging.• Simplest example:  $\hat{g}(x) = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$  for  $x$ 

$$\begin{aligned}
 \text{Then } E^*(g^*(x)) &= E^*[\bar{Y}^*] = E\left[\frac{1}{n} \sum_{i=1}^n Y_i^*\right] = \frac{1}{n} \sum_{i=1}^n E^*[Y_i^*] \stackrel{\text{i.i.d.}}{=} \\
 &\stackrel{\text{i.i.d.}}{=} \frac{1}{n} \sum_{i=1}^n E^*[Y_i^*] = E^*[Y_1^*] \quad \text{and given } Y_1^* = \begin{cases} Y_1 & \text{w.p. } 1/n \\ Y_n & \text{w.p. } 1/n \\ \vdots & \vdots \end{cases} \text{ so} \\
 &= \sum_{i=1}^n Y_i \cdot \frac{1}{n} = \bar{Y}
 \end{aligned}$$

• Bagging should be used for  
non-linear estimators!  
ex. trees.

(Key point: Bagging - decrease variance without increasing bias by aggregating the estimators, predicted outcome probabilities, ...)

• Bagging for Classification Given data  $(X_1, Y_1), \dots, (X_n, Y_n)$ ,  $X_i \in \mathbb{R}^p$ ,  $Y_i \in \{1, \dots, K\}$  <sup>categories</sup>• Base procedure:  $\hat{g}(\cdot): \mathbb{R}^p \rightarrow \{1, \dots, K\}$  (doesn't make sense to average the category #'s)Methods. (1) Majority Vote:  $\hat{g}^{\text{bag}}(x) = \underset{k=1, \dots, K}{\operatorname{argmax}} \sum_{b=1}^B \mathbb{1}_{\{\hat{g}^{*b}(x)=k\}}$ • ex. a thought experiment: ~~Given input  $x$  and the true class is 1~~  $\{1, 2, 3\}$ • For some given input  $x$  we have  $B$  independent classifiers  $\hat{g}^{*b}(x)$ ,  $b=1, \dots, B$  and each classifier has a misclassification rate of 0.4.• Note that this is a thought experiment; the classifiers are not independent for bagging because we take bootstrap samples from the same training data.• Assume wlog that the true class is 1. This implies that  $P(\hat{g}^b(x)=1) = 0.6$ ,  $P(\hat{g}^b(x)=2) = 0.2$ • Now we form a bagged classifier  $\hat{g}^{\text{bag}}(x) = \underset{k=1, 2}{\operatorname{argmax}} \sum_{b=1}^B \mathbb{1}_{\{\hat{g}^{*b}(x)=k\}}$ • What is the misclassification rate of  $\hat{g}^{\text{bag}}(x)$ ? ie. What is  $P(\hat{g}^{\text{bag}}(x) = 2)$ ?



## Bagging for Classification, cont'd.

(ex. continued) (What is the misclassification rate of  $\hat{g}^{\text{bag}}(x)$ ? i.e.:

$$P(\hat{g}^{\text{bag}}(x) = 2) = P\left(\underbrace{\sum_{b=1}^B \mathbb{1}_{\{\hat{g}^b(x) = 2\}}}_W > B/2\right)$$

$$\text{Let } W = \sum_{b=1}^B \mathbb{1}_{\{\hat{g}^b(x) = 2\}} \quad = P(W > B/2)$$

$$W \sim \text{Binom}(B, 0.6)$$

$$E[W] = 0.6B$$

$$\text{Var}[W] = B \cdot 0.6 \cdot 0.4 = 0.24B$$

$$\text{so then } W \approx \mathcal{N}(0.6B, 0.24B)$$

(Gaussian approximation)

$$\rightarrow = P\left(\frac{W - 0.6B}{\sqrt{0.24B}} > \frac{B/2 - 0.6B}{\sqrt{0.24B}}\right)$$

$$\approx P\left(Z > \underbrace{-\frac{0.1B}{\sqrt{0.24B}}}_{\rightarrow -\infty}\right) \rightarrow 1 \text{ as } B \rightarrow \infty$$

So, we get a perfectly bad classifier.

- Bagging a good classifier can improve performance, but bagging a bad classifier can degrade performance

(Suppose we want to know the class probabilities)

↳ Voting probabilities  $\frac{1}{B} \sum_{b=1}^B \mathbb{1}_{\{\hat{g}^b(x) = k\}}$ ,  $k = 1, \dots, K$  are not good estimates of the estimates of the class probabilities  $P(Y=k|X=x)$

- ex. Suppose that  $P(Y=1|X=x) = 0.75$  and that all  $\hat{g}^{*b}(x) = 1 \quad \forall b=1, \dots, B$   
Then voting probability  $\frac{1}{B} \sum_{b=1}^B \mathbb{1}_{\{\hat{g}^{*b}(x) = 1\}} = 1$

- Each tree gives predicted class probabilities  $\hat{p}_k^{*b}(x)$ , where

$\hat{p}_k^{*b}(x) \equiv$  proportion of observations with class  $k$  among all observations in the leaf node of tree  $g^{*b}$  that contains  $x$ ,  $\forall$  classes ( $k=1, \dots, K$ )

- We can bag the predicted class probabilities  $\hat{p}_k^{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^B \hat{p}_k^{*b}(x)$
- The final bagged classifier chooses the class with the highest  $\hat{p}_k^{\text{bag}}(x)$   
ie.  $\hat{g}^{\text{bag}}(x) = \underset{k=1, \dots, K}{\operatorname{argmax}} \hat{p}_k^{\text{bag}}(x)$

- Advantages of bagging: We get estimated class probabilities, sometimes better prediction performance

(Slides: Figure 8.10) - consensus majority vote vs. bagged (from this figure)  
↳ get an idea about the uncertainties

(Ref Rcode13)

- line 130: can use randomForest for regular bagging by setting mtry (# of variables to consider for each split)

- line 142: default uses 1/3 of variables @ each split





## Random Forests

- Actual random forests: trees will be more independent (than from bagging) so you get better performance.
- Pros/Cons of random forests: better classification and regression performance, but harder to interpret the model.
- "Out of bag" (OOB) (~~from~~ samples not bagged in training) - provides error estimates so we don't need to do X Valid'n. train  
 $\#$  OOB values are comparable to test OOB values.  
~~OOB is not~~ OOB  $\neq$  training MSE!  
 error estimates (MSE)
- Comparing Bagged trees and random forests:

	Trees	Bagged trees	Random Forests
Interpretation	+	-	- (we however can see/get variable importance.)
Computation	+	-	~
Performance	-	+	++
OOB error estimates	need xv)	+	+
Overfitting	Yes (need xv)	OK	OK

Variable Importance: Metric for # times used to split, ~~area~~ a variable was used in distinguishing data //

## Lecture

Fri. 01 June 2018

(Ref Rcode14): 3 training sets sampled from the data (~ line 20)

MSE	train1	train2	train3	
Bagging OOB	13.7	15.7	11.8	• Note that training MSE $\neq$ OOB <sup>MSE</sup> , and that OOB is comparable to test MSE.
Bagging testset	15.8	11.0	15.4	
Random forest OOB	12.4	14.7	12.8	• Also note that variability <del>is</del> a MSE (as a random var?) - ie. OOB is neither optimistic nor pessimistic (generally). - Train and Test MSEs are in similar ranges
Random forest test set	16.8	12.5	14.1	

(line 80) - randomForest(.) : maxnodes = # arg  $\equiv$  max size (~~# of data~~ depth of the tree. (ie. (# of data in each node))

(line 88) df: eg. 

0.044...	0.955...
----------	----------

  
 is proportion predicted class class 2  
 min node size (# at least, in each node)

- Deep trees: bias is less of an issue; the  
 - larger the node ( $\leadsto$  shallower tree) the more problematic  
 - deep tree results: more scatter/variability but less bias. (mean closer to true value (2/3))

(line 106) ranger (library). Using probability method.

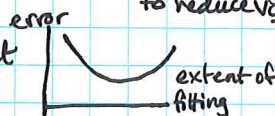
Shallow trees with prob. method (preserves) low bias while also lowering variability (lower var than in ~~high~~ deep trees (with prob. method))

★ Random Forests  
 with Probability Method  
 on Exam!





- Bagging - Start with a deep tree; flexible estimate  
 $\hookrightarrow$  low bias, high variance  
 - bootstrap aggregating reduces variance  
 - as  $B \rightarrow \infty$ , we don't have to worry about overfitting (since we are combining things to reduce variance)
- Boosting - Start with a simple estimator; ("tree stump") with 1 split  
 $\hookrightarrow$  high bias, low variance  
 - combine estimates to reduce bias  
 - as  $B \rightarrow \infty$ , overfitting tends to occur.



(Ref Algo 8.2 from ISLR - in slides)

- line (26):  $\lambda$  is usually a small value that is fixed (usually)  
 $\hookrightarrow$  This is a "slow learner"
- choose  $B$  by  $X$ Valid'n  $\hookrightarrow$  from big  $B$  you can do  $X$ Valid'n from small  $B$  subsets.
- reweight (increase weight of) misclassified outcomes in  $X$ Valid'n ~~to~~ so training improves performance more, i.e. "pay more attention" to them.

(Ref Rcode14) gbm library for boosting (line 423)

- `gbm()` arg:  $n.tree \equiv B$   
 $interaction.depth \equiv$  # of variables from splits that are allowed to interact
- line 130: partial dependence plot
- Adjust params ( $\lambda$ , # of trees, ...) to improve MSE. Can use  $X$ Valid'n to choose both  $\lambda$  and # of trees, but usually  $\lambda$  is fixed and  $X$ Valid'n is done for  $B$ .

## Section

Fri 07 June 18

(Series question is from ISLR). - see online section notes.

- check quality of model by  $X$ Valid'n.
- (Assume that we can randomly shuffle the data)
- "oof" := out-of-fold: vector of whole data set and we fit on ~~all data~~  $B$  all data \ {current fold}  $\rightarrow$  vector of prediction  
 $\xrightarrow{\text{get a}}$
- What is the probability threshold of decision? Use not 0.5 (fixed value) but the mean of oof.

• Boosting: fit tree based on ~~previous~~ <sup>the residues</sup> of previous trees

- Importance (plot): of variables - in splits, how many times was a variable used to split, distinguish the data?

• Bagging vs. RF.

- First question: Using trees, "How ~~do~~ we deal with overfitting?"  
 $\rightarrow$  combine / aggregate trees, i.e. bagging
- Can use random Forests library (with bagging as a special case.)  
 $\hookrightarrow$  so you don't have to use a specific bagging library.
- Bagging - instead of resampling features (each step, take all the features (for each step)).