

Summary and discussion of: “Confidence Intervals for Random Forests: The Jackknife and the Infinitesimal Jackknife”

Journal club report

Zhipeng LIANG

1 Summary

1.1 Bagging and Random Forest

Suppose that we have training examples $Z_1 = (x_1, y_1), \dots, Z_n = (x_n, y_n)$, the dataset $\mathcal{D} = \{Z_i\}_{i=1}^n$, an input x to a prediction problem, and a base learner $\hat{\theta} = t(x; Z_1, \dots, Z_n)$.

For bagging method, we stabilize the base learner t by resampling the training data. The bagging can be viewed as the solution to

$$\hat{\theta}^B(x) = \frac{1}{B} \sum_{b=1}^B t_b^*(x), \text{ where } t_b^*(x) = t(x; Z_{b1}^*, \dots, Z_{bn}^*), \quad (1)$$

where the Z_{bi}^* are the i -th elements in the b -th bootstrap sample.

Moreover, the random forest constructs trees base on not only bootstrap samples but also random feature selection. Thus it can be viewed as the solution to the following problem

$$\hat{\theta}^{RF}(x) = \frac{1}{B} \sum_{b=1}^B \frac{1}{K} \sum_{k=1}^K t_b^*(x; \xi_{kb}, Z_{b1}^*, \dots, Z_{bn}^*) \text{ with } \xi_{kb} \stackrel{\text{iid}}{\sim} \Xi, \quad (2)$$

where ξ_{kb} can be viewed as the modeling for the feature selection noise.

To be specific, we introduce the generalized bagged version of $\hat{\theta}^\infty(x)$ is defined as

$$\hat{\theta}^\infty(x) = \mathbb{E}_{Z_i \sim P^*, \xi \sim \Xi} [t(x; \xi, Z_1^*, \dots, Z_n^*)], \quad (3)$$

where P^* is the bootstrap distribution and Ξ is the auxiliary noise distribution. Then (1) can be viewed as a Monte Carlo approximation to the bootstrap distribution P^* with trivial distribution Ξ while random forest is a Monte Carlo approximation to both the bootstrap distribution P^* and nontrivial auxiliary noise distribution.

First, to understand the bootstrap distribution, note that the empirical distribution is

$$P_\delta(x, y) = \frac{1}{n} \sum_{i=1}^n \delta(x = x_i, y = y_i),$$

where $\delta(\cdot)$ is the Dirac mass centered at (x_i, y_i) . Then the bootstrap distribution is

$$P_*(x, y) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^n k \delta(x = x_i, y = y_i) \frac{1}{k!} \exp^{-1}.$$

The derivation is from the approximation of the poisson distribution with parameter $\lambda = 1$ Poisson(1) to the multinomial distribution $\text{Multi}(n, \frac{1}{n})$ when $n \rightarrow \infty$ since we use Poisson bootstrap here.

1.2 IJ and J estimator

The goal of this paper is to study the sampling variance of bagged learners

$$V(x) = \text{Var}(\hat{\theta}^\infty(x)).$$

This paper consider two basic estimates of $V(\cdot)$:

1. The Infinitesimal Jackknife estimate

$$\hat{V}_{IJ}^\infty = \sum_{i=1}^n \text{Cov}_* [N_i^*, t^*(x)]^2$$

where $\text{Cov}[N_i^*, t(x)]$ is the covariance between $t^*(x)$ and the number of times N_i^* the i -th training example appears in a bootstrap sample. This is the direct application of the Theorem 1 in [1].

2. The Jackknife-after-Bootstrap estimate

$$\hat{V}_J^\infty = \frac{n-1}{n} \sum_{i=1}^n \left(\bar{t}_{(-i)}^*(x) - \bar{t}^*(x) \right)^2$$

where $\bar{t}_{(-i)}^*(x)$ is the average of $t^*(x)$ over all the bootstrap samples not containing the i -th example and $\bar{t}^*(x)$ is the mean of all the $t^*(x)$.

In practice, we can only ever work with a finite number B of bootstrap replicates. The natural Monte Carlo approximations to the estimators introduced above are

$$\hat{V}_{IJ}^B = \sum_{i=1}^n \widehat{\text{Cov}}_i^2 \text{ with } \widehat{\text{Cov}}_i = \frac{\sum_b (N_{bi}^* - 1) (t_b^*(x) - \bar{t}^*(x))}{B}, \quad (4)$$

and

$$\hat{V}_J^B = \frac{n-1}{n} \sum_{i=1}^n \hat{\Delta}_i^2,$$

where

$$\hat{\Delta}_i = \hat{\theta}_{(-i)}^B(x) - \hat{\theta}^B(x)$$

and

$$\hat{\theta}_{(-i)}^B(x) = \frac{\sum_{\{b: N_{bi}^*=0\}} t_b^*(x)}{|\{N_{bi}^*=0\}|},$$

where N_{bi}^* indicates the number of times the i -th observation appears in the bootstrap sample b .

However, these finite- B estimates of variance are often badly biased upwards if the number of bootstrap samples B is too small. To correct the bias, we investigate the bias

$$\mathbb{E}_* [\hat{V}_{IJ}^B] - \hat{V}_{IJ}^\infty = \sum_{i=1}^n \text{Var}_* [C_i], \text{ where } C_i = \frac{\sum_b (N_{bi}^* - 1) (t_b^* - \bar{t}^*)}{B}.$$

They use the approximation

$$\text{Var}_* [(N_{bi}^* - 1) (t_b^* - \bar{t}^*)] \approx \text{Var}_* [N_{bi}^*] \text{Var}_* [t_b^*]$$

which holds when $t(x; Z_1^*, Z_2^*, \dots, Z_n^*) = \frac{1}{n} \sum_{i=1}^n Z_i^*$ as the mean estimator and Poisson bootstrap. To be specific, the approximation error is dominated by the approximated term.

Thus this bias can be approximated as

$$\mathbb{E}_* [\hat{V}_{IJ}^B] - \hat{V}_{IJ}^\infty \approx \frac{n}{B^2} \sum_{b=1}^B (t_b^*(x) - \bar{t}^*(x))^2.$$

Thus we have the bias-correction version

$$\hat{V}_{IJ-U}^B = \hat{V}_{IJ}^B - \frac{n}{B^2} \sum_{b=1}^B (t_b^*(x) - \bar{t}^*(x))^2, \quad (5)$$

and

$$\hat{V}_{J-U}^B = \hat{V}_J^B - (e - 1) \frac{n}{B^2} \sum_{b=1}^B (t_b^*(x) - \bar{t}^*(x))^2. \quad (6)$$

1.3 Upward Bias and Downward Bias

Suppose that we have data Z_1, \dots, Z_n drawn independently from a distribution F , and compute the estimate $\hat{\theta}^\infty$ from the data as in 3. Then using the ANOVA decomposition [4] we know

$$\text{Var}_F [\hat{\theta}^\infty] = V_1 + V_2 + \dots + V_n, \quad (7)$$

where V_k are all non-negative.

[4] also show that, under general condition, the Jackknife estimator of variance is biased upward. To be specific, the Jackknife estimator for the $n+1$ -th data point has the variance

$$\mathbb{E}_F [\hat{V}_J^\infty] = V_1 + 2V_2 + 3V_3 + \dots + nV_n. \quad (8)$$

However, as pointed out by [5, 1], \hat{V}_{IJ}^∞ is equivalent to the variance of a "bootstrap Hjek projection", i.e.,

$$\hat{V}_{IJ}^\infty \approx \sum_{i=1}^n h_F^2(Z_i) \quad (9)$$

where $h_F(Z) = \mathbb{E}_F[\hat{\theta}^\infty | Z_1 = Z] - \mathbb{E}_F[\hat{\theta}^\infty]$. (9) suggests that

$$\mathbb{E}_F[\hat{V}_{IJ}^\infty] \approx V_1. \quad (10)$$

Compare (7), (8) and (10) we can conclude that Jackknife estimator has a upward bias by double the second-order term, triple the third-term and so on while Infinite Jackknife estimator has a downward bias by dropping higher-order terms. Thus this paper suggests to use the arithmetic mean of \hat{V}_J^∞ and \hat{V}_{IJ}^∞ .

2 Result and Discussion

2.1 Empirical Bayes Calibration

I learnt Stefan Wager's randomForestCI and sklearn's randomForestCI and wrote my own code. All the figures can be reproduced using the code in this repository. Stefan Wager's code is for R while sklearn's is based on python. However, neither of them implemented Jackknife-after-Bootstrap estimator.

2.2 Figure 2

In the first experiment I testing the performance of my version of the bias-corrected infinitesimal jackknife estimate of variance for bagged predictors, defined in 5, on the simulation data. The experiment setup is the same as described in the paper while here I only outline some points. The data is generated via $y = f(x) + \epsilon$ where ϵ is a Gaussian noise and

$$f(x) = \begin{cases} 14 & 0.45 \leq x \leq 0.55 \\ 7 & 0.35 \leq x \leq 0.65 \\ 0 & 0 \leq x \leq 1. \end{cases} \quad (11)$$

There is a jump of level 7 at four discontinuous points, 0.35, 0.45, 0.55 and 0.65. Thus the ideal variance estimator should capture them. As we can see, Figure 2.2 shows that the infinitesimal jackknife estimator of the bagging predictors does achieve it. However, in Figure 2.2 my version of Jackknife-after-Bootstrap estimator although react to the discontinuous points but fail to capture the true variance.

Before we start to present other experiments, note that both in the original code by Stefan Wager and sklearn's implementation, the bias-corrected infinitesimal jackknife estimate in (5) is not directly used. Instead, Wager conducted empirical Bayes calibration for the infinitesimal jackknife estimate using the g-modeling in [2]. Without such a empirical Bayes calibration, most of the elements in the infinitesimal jackknife estimate will be negative in my implementation.

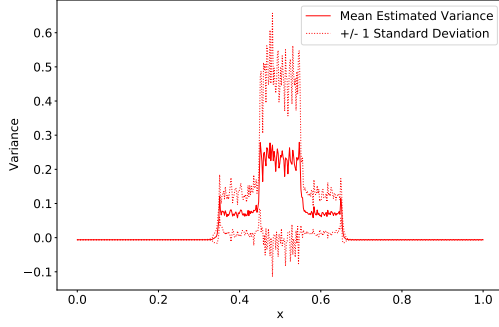


Figure 1: Jackknife-after-Bootstrap

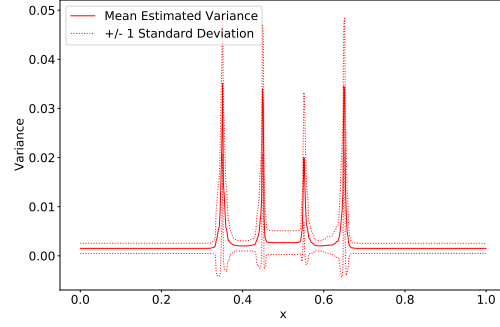


Figure 2: Infinitesimal Jackknife

2.3 Figure 3

In this experiment, I test the advantage of bias-correction version against. I test their performance on the cholesterol dataset [3]. I closely follow the experiment setup in [1]. For completeness I outline some key ingredients here. I use the bagging polyregression and the choice of degree follows C_p criterion [6], i.e., for polyregression regression with degree m , the C_p value is defined as

$$C_p(m) = \left\| \mathbf{y} - X_m \hat{\beta}_m \right\|^2 + 2\sigma^2 m,$$

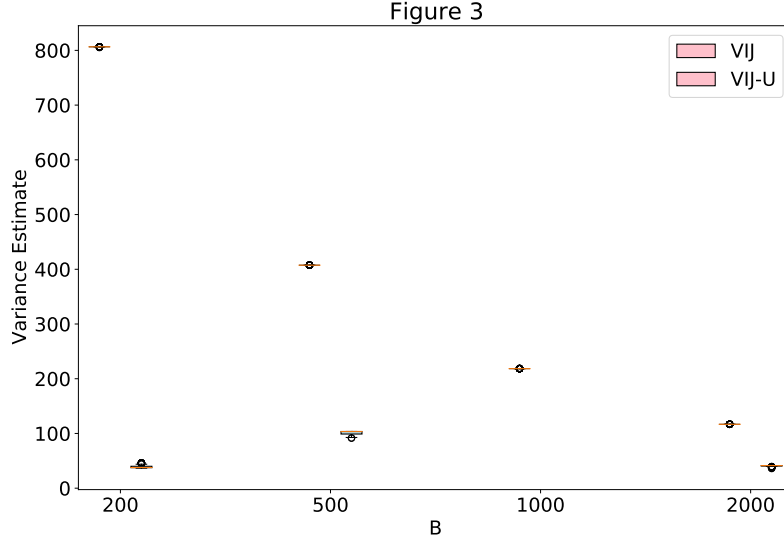
where X_m is the design matrix with polynomial factors whose degree at most m and

$$\hat{\beta}_m = \underset{\hat{\beta}}{\operatorname{argmin}} \left\| \mathbf{y} - X_m \hat{\beta} \right\|^2.$$

Follows the setup in [1], we use the value $\sigma = 22$ in this experiments. We use the best degree for each bootstrap sample for the final regression model and apply the infinite Jackknife 4 and bias-correction version 5. As we can see in Figure 2.3, the bias-correction version (on the right column) is significantly lower than the biased one. Moreover, as the bootstrap sample sizes increase, the Monte Carlo error decreases which leads to the closer performance between biased IJ and bias-correction one.

2.4 Figure 4

In this section we aim to gain some qualitative insight from the infinite jackknife estimator on the random forest. From Figure ?? we plot test-set predictions against IJ-U estimates of standard error for all three random forests with maximum number of features as 5, 19 and 57 (which reduce to the bagging version). The sample variance increases with the maximum number of select features. Thus the $m = 57$ forest tends to suffer overfitting under this limited sample size, which hinder it from achieving the best test error amongs them. In contrast, the $m = 5$ forest has small variance across all samples, which may



implies it suffices from bias due to the limited capacity of utilizing features. From Figure ?? we plot the prediction probability of $m = 5$ forest versus the difference between prediction probability of $m = 19$ against $m = 5$. To visualize the main trend from the scatters, we interpolate the scatters with a polynomial curve with degree 5 (red line). The red line remains negative when prediction is less than $1/2$ and remains positive when prediction is greater than $1/2$, which implies that the $m = 19$ forest has more confidence along the correct direction, and further verify the insight gained from the Infinite Jackknife estimator.

3 Conclusion

References

- [1] Bradley Efron. *Model Selection Estimation and Bootstrap Smoothing*. Division of Biostatistics, Stanford University, 2012.
- [2] Bradley Efron. Two modeling strategies for empirical bayes estimation. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 29(2):285, 2014.
- [3] Bradley Efron and David Feldman. Compliance as an explanatory variable in clinical trials. *Journal of the American Statistical Association*, 86(413):9–17, 1991.
- [4] Bradley Efron and Charles Stein. The jackknife estimate of variance. *The Annals of Statistics*, pages 586–596, 1981.
- [5] Louis A Jaeckel. The infinitesimal jackknife. 1972.
- [6] Colin L Mallows. Some comments on cp. *Technometrics*, 42(1):87–94, 2000.

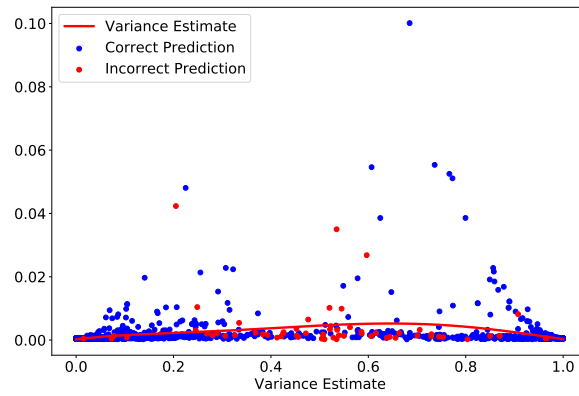


Figure 3: Jackknife-after-Bootstrap

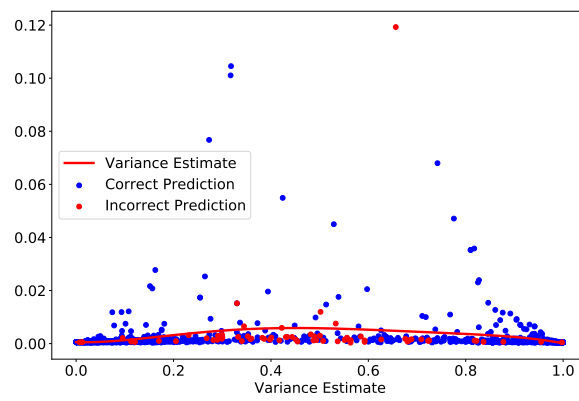


Figure 4: Infinitesimal Jackknife

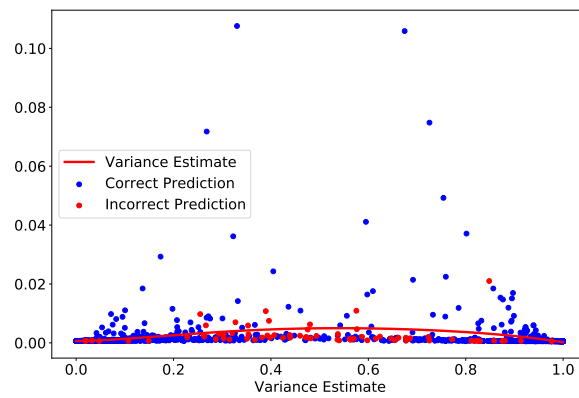


Figure 5: Infinitesimal Jackknife

