

The jackknife+: distribution-free prediction

Rina Foygel Barber

(with Emmanuel Candès, Aaditya Ramdas, Ryan Tibshirani)

<http://www.stat.uchicago.edu/~rina/>

Collaborators



Emmanuel Candès



Aaditya Ramdas



Ryan Tibshirani

- Thanks to American Institute of Math (AIM)

The prediction problem

Setting:

- Training data $(X_1, Y_1), \dots, (X_n, Y_n) \rightsquigarrow$ fit model $\hat{\mu}(X_i) \approx Y_i$
- Test point (X_{n+1}, Y_{n+1}) from same distribution
- If $\hat{\mu}$ overfits to training data,

$$|Y_{n+1} - \hat{\mu}(X_{n+1})| \gg \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{\mu}(X_i)|$$

The prediction problem

Goal: build prediction band C as a function of the training data, such that $C(X_{n+1})$ is likely to contain Y_{n+1}

- Want to be distribution-free — coverage holds w/o assumptions on distrib. of (X, Y)
- Want to be efficient — minimize width of interval $C(X_{n+1})$

Defining distribution-free coverage

Definition:

A method satisfies distribution-free coverage at level $1 - \alpha$ if

$$\mathbb{P} \{ Y_{n+1} \in C(X_{n+1}) \} \geq 1 - \alpha$$

w.r.t. $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1}) \stackrel{\text{iid}}{\sim} P$, for any P .

Defining distribution-free coverage

Definition:

A method satisfies distribution-free coverage at level $1 - \alpha$ if

$$\mathbb{P} \{ Y_{n+1} \in C(X_{n+1}) \} \geq 1 - \alpha$$

w.r.t. $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1}) \stackrel{\text{iid}}{\sim} P$, for any P .

The coverage rate is averaged over training data & over test points

1. Quantiles of a sample:

$Q_{1-\alpha}(x_1, \dots, x_n) = \lceil (1 - \alpha)(n + 1) \rceil$ -th smallest value of x_1, \dots, x_n

Abusing notation...

$Q_\alpha(x_1, \dots, x_n) = \lfloor \alpha(n + 1) \rfloor$ -th smallest value of x_1, \dots, x_n

1. Quantiles of a sample:

$Q_{1-\alpha}(x_1, \dots, x_n) = \lceil (1 - \alpha)(n + 1) \rceil$ -th smallest value of x_1, \dots, x_n

Abusing notation...

$Q_\alpha(x_1, \dots, x_n) = \lfloor \alpha(n + 1) \rfloor$ -th smallest value of x_1, \dots, x_n

2. Fitted regression function

$$\hat{\mu} = \mathcal{A}\left(\underbrace{(X_1, Y_1), \dots, (X_n, Y_n)}_{\text{invariant to ordering of training points}}\right)$$

For any subset $S \subset \{1, \dots, n\}$,

$$\hat{\mu}_{-S} = \mathcal{A}\left((X_i, Y_i) : i \in \{1, \dots, n\} \setminus S\right)$$

Some options...

Naive method:

$$C(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm Q_{1-\alpha}(|Y_i - \hat{\mu}(X_i)|)$$

- Computational cost: one regression $\hat{\mu}$
- Poor coverage in practice (overfitting)

Some options...

Holdout method:

Choose holdout set $S_{\text{hold}} \subset \{1, \dots, n\}$

$$C(X_{n+1}) = \hat{\mu}_{-S_{\text{hold}}}(X_{n+1}) \pm Q_{1-\alpha} \left(|Y_i - \hat{\mu}_{-S_{\text{hold}}}(X_i)| : i \in S_{\text{hold}} \right)$$

- Computational cost: one regression $\hat{\mu}_{-S_{\text{hold}}}$
- Wider intervals (reduced sample size $n - |S_{\text{hold}}| < n$)

¹Papadopoulos 2008, Vovk 2012, Lei et al. 2018

Some options...

Holdout method:

Choose holdout set $S_{\text{hold}} \subset \{1, \dots, n\}$

$$C(X_{n+1}) = \hat{\mu}_{-S_{\text{hold}}}(X_{n+1}) \pm Q_{1-\alpha} \left(|Y_i - \hat{\mu}_{-S_{\text{hold}}}(X_i)| : i \in S_{\text{hold}} \right)$$

- Computational cost: one regression $\hat{\mu}_{-S_{\text{hold}}}$
- Wider intervals (reduced sample size $n - |S_{\text{hold}}| < n$)
- Distribution-free coverage due to exchangeability¹
of $|Y_i - \hat{\mu}_{-S_{\text{hold}}}(X_i)|$, $i \in S_{\text{hold}} \cup \{n+1\}$

¹Papadopoulos 2008, Vovk 2012, Lei et al. 2018

Some options...

Cross-validation:

Split data into $S_1 \cup \dots \cup S_K$

For each $i \in S_k$, $R_i^{\text{CV}} = |Y_i - \hat{\mu}_{-S_k}(X_i)|$

$$C(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm Q_{1-\alpha} \left(R_i^{\text{CV}} \right)$$

— Computational cost: $K + 1$ regressions

Some options...

Jackknife a.k.a. leave-one-out cross-validation ($K = n$)

Residuals $R_i^{\text{LOO}} = |Y_i - \hat{\mu}_{-i}(X_i)|$

$$C(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm Q_{1-\alpha} \left(R_i^{\text{LOO}} \right)$$

Some options...

Jackknife a.k.a. leave-one-out cross-validation ($K = n$)

Residuals $R_i^{\text{LOO}} = |Y_i - \hat{\mu}_{-i}(X_i)|$

$$C(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm Q_{1-\alpha}(R_i^{\text{LOO}})$$

— Predictive coverage holds under assumptions²

Asymptotic stability: $\hat{\mu}(X_{n+1}) \approx \hat{\mu}_{-i}(X_{n+1})$

²Steinberger & Leeb 2018

Some options...

Cross-conformal prediction:³

Split data into $S_1 \cup \dots \cup S_K$,

& find all y overlapping $\geq n - (1 - \alpha)(n + 1)$

of the intervals $\hat{\mu}_{-S_k}(X_{n+1}) \pm R_i^{\text{CV}}$

— Computational cost: $K + 1$ regressions

— Distribution-free theory

Coverage $\geq 1 - 2\alpha - 2K/n$

³Vovk 2015, Vovk et al 2018

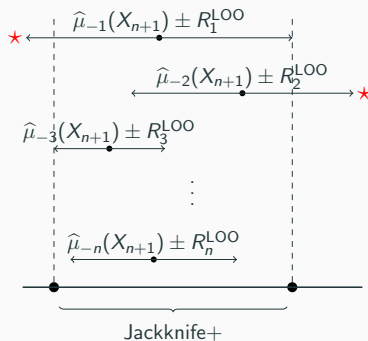
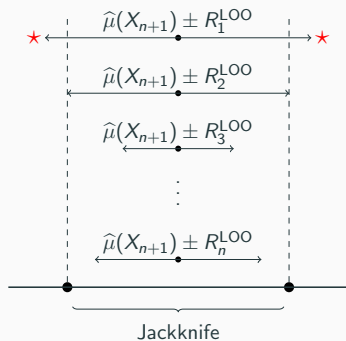
Jackknife+:

$$C(X_{n+1}) = \left[Q_{\alpha} \left(\hat{\mu}_{-i}(X_{n+1}) - R_i^{\text{LOO}} \right), Q_{1-\alpha} \left(\hat{\mu}_{-i}(X_{n+1}) + R_i^{\text{LOO}} \right) \right]$$

Compare to jackknife:

$$C(X_{n+1}) = \left[Q_{\alpha} \left(\hat{\mu}(X_{n+1}) - R_i^{\text{LOO}} \right), Q_{1-\alpha} \left(\hat{\mu}(X_{n+1}) + R_i^{\text{LOO}} \right) \right]$$

Our methods



Extension to CV+:

$$C(X_{n+1}) = \left[Q_{\alpha} \left(\hat{\mu}_{-S_k}(X_{n+1}) - R_i^{\text{CV}} \right), Q_{1-\alpha} \left(\hat{\mu}_{-S_k}(X_{n+1}) + R_i^{\text{CV}} \right) \right]$$

CV+ interval \supseteq Vovk's cross-conformal prediction set

Theorem: For any distrib. P and any \mathcal{A} , jackknife+ satisfies

$$\mathbb{P}\{Y_{n+1} \in C(X_{n+1})\} \geq 1 - 2\alpha.$$

Theorem: For any distrib. P and any \mathcal{A} , jackknife+ satisfies

$$\mathbb{P}\{Y_{n+1} \in C(X_{n+1})\} \geq 1 - 2\alpha.$$

Theorem: For any distrib. P and any \mathcal{A} , K -fold CV+ satisfies

$$\mathbb{P}\{Y_{n+1} \in C(X_{n+1})\} \geq \begin{cases} 1 - 2\alpha - 1/K & \text{(new)} \\ 1 - 2\alpha - 2K/n & \text{(Vovk et al)} \end{cases}$$

$$\rightsquigarrow \geq 1 - 2\alpha - \sqrt{2/n}.$$

Proof sketch (jackknife+)

Define:

- Regression $\tilde{\mu}_{-\{i,j\}}$ for each $i, j \in \{1, \dots, n+1\}$
- Residuals $R_{ij} = |Y_i - \tilde{\mu}_{-\{i,j\}}(X_i)|$
- Comparison matrix for R_{ij} vs R_{ji} :

$$A_{ij} = \mathbf{1}\{R_{ij} > R_{ji}\} \quad (\text{and } A_{ii} = 0)$$

Proof sketch (jackknife+)

Define:

- Regression $\tilde{\mu}_{-\{i,j\}}$ for each $i, j \in \{1, \dots, n+1\}$
- Residuals $R_{ij} = |Y_i - \tilde{\mu}_{-\{i,j\}}(X_i)|$
- Comparison matrix for R_{ij} vs R_{ji} :

$$A_{ij} = \mathbf{1}\{R_{ij} > R_{ji}\} \quad (\text{and } A_{ii} = 0)$$

Verify: if $Y_{n+1} \notin C(X_{n+1})$ then

$$\sum_{i=1}^n \mathbf{1}\left\{ |Y_{n+1} - \hat{\mu}_{-i}(X_{n+1})| > |Y_i - \hat{\mu}_{-i}(X_i)| \right\} \geq (1 - \alpha)(n + 1)$$

Proof sketch (jackknife+)

Define:

- Regression $\tilde{\mu}_{-\{i,j\}}$ for each $i, j \in \{1, \dots, n+1\}$
- Residuals $R_{ij} = |Y_i - \tilde{\mu}_{-\{i,j\}}(X_i)|$
- Comparison matrix for R_{ij} vs R_{ji} :

$$A_{ij} = \mathbf{1}\{R_{ij} > R_{ji}\} \quad (\text{and } A_{ii} = 0)$$

Verify: if $Y_{n+1} \notin C(X_{n+1})$ then

$$\sum_{i=1}^n \mathbf{1}\left\{ \overbrace{|Y_{n+1} - \hat{\mu}_{-i}(X_{n+1})|}^{R_{n+1,i}} > \overbrace{|Y_i - \hat{\mu}_{-i}(X_i)|}^{R_{i,n+1}} \right\} \geq (1 - \alpha)(n + 1)$$

Proof sketch (jackknife+)

Define:

- Regression $\tilde{\mu}_{-i,j}$ for each $i, j \in \{1, \dots, n+1\}$
- Residuals $R_{ij} = |Y_i - \tilde{\mu}_{-i,j}(X_i)|$
- Comparison matrix for R_{ij} vs R_{ji} :

$$A_{ij} = \mathbf{1}\{R_{ij} > R_{ji}\} \quad (\text{and } A_{ii} = 0)$$

Verify: if $Y_{n+1} \notin C(X_{n+1})$ then

$$\sum_{i=1}^n \underbrace{\mathbf{1}\left\{ \overbrace{|Y_{n+1} - \hat{\mu}_{-i}(X_{n+1})|}^{R_{n+1,i}} > \overbrace{|Y_i - \hat{\mu}_{-i}(X_i)|}^{R_{i,n+1}} \right\}}_{A_{n+1,i}} \geq (1 - \alpha)(n + 1)$$

Proof sketch (jackknife+)

Summary so far:

$$Y_{n+1} \notin C(X_{n+1}) \quad \Rightarrow \quad \sum_i A_{n+1,i} \geq (1 - \alpha)(n + 1)$$

Proof sketch (jackknife+)

Summary so far:

$$Y_{n+1} \notin C(X_{n+1}) \quad \Rightarrow \quad \sum_i A_{n+1,i} \geq (1 - \alpha)(n + 1)$$

Data points are exchangeable, so

$$\mathbb{P} \left\{ \sum_i A_{n+1,i} \geq (1 - \alpha)(n + 1) \right\} = \frac{\mathbb{E} [\# \text{ rows } j \text{ with } \sum_i A_{ji} \geq (1 - \alpha)(n + 1)]}{n + 1}$$

Proof sketch (jackknife+)

Summary so far:

$$Y_{n+1} \notin C(X_{n+1}) \Rightarrow \sum_i A_{n+1,i} \geq (1 - \alpha)(n + 1)$$

Data points are exchangeable, so

$$\mathbb{P} \left\{ \sum_i A_{n+1,i} \geq (1 - \alpha)(n + 1) \right\} = \frac{\mathbb{E} [\# \text{ rows } j \text{ with } \sum_i A_{ji} \geq (1 - \alpha)(n + 1)]}{n + 1}$$

Deterministically, at most $2\alpha(n + 1)$ such rows $\Rightarrow \mathbb{P} \{ \dots \} \leq 2\alpha$

Proof sketch (jackknife+)

Worst case for # rows j with $\sum_i A_{ji} \geq (1 - \alpha)(n + 1)$:

50% 1's	all 1's	} $2\alpha(n + 1)$ rows
all 0's	all 0's	
		} $(1 - 2\alpha)(n + 1)$ rows

$$\left[\min_i \hat{\mu}_{-i}(X_{n+1}) - Q_{1-\alpha}\left(R_i^{\text{LOO}}\right), \max_i \hat{\mu}_{-i}(X_{n+1}) + Q_{1-\alpha}\left(R_i^{\text{LOO}}\right) \right]$$

Theorem: For any distrib. P and any \mathcal{A} , jackknife-minmax satisfies

$$\mathbb{P}\{Y_{n+1} \in C(X_{n+1})\} \geq 1 - \alpha.$$

Results so far

Method	Assumption-free theory	Typical empirical coverage
Naive	No guarantee	$< 1 - \alpha$
Jackknife	No guarantee	$\approx 1 - \alpha$
Jackknife+	$1 - 2\alpha$ coverage	$\approx 1 - \alpha$
Jackknife-minmax	$1 - \alpha$ coverage	$> 1 - \alpha$

Theorem: There exists a data distribution P and a regression algorithm \mathcal{A} s.t.

1. For the naive method, coverage = 0
2. For jackknife, coverage = 0

Furthermore if $\alpha \leq \frac{1}{2}$, there exist P and \mathcal{A} s.t.

3. For jackknife+, coverage $\leq 1 - 2\alpha + o(1)$

Out-of-sample stability:

$$\mathbb{P} \left\{ \left| \hat{\mu}(X_{n+1}) - \hat{\mu}_{-i}(X_{n+1}) \right| \leq \epsilon \right\} \geq 1 - \nu$$

In-sample stability:

$$\mathbb{P} \left\{ \left| \hat{\mu}(X_i) - \hat{\mu}_{-i}(X_i) \right| \leq \epsilon \right\} \geq 1 - \nu$$

Stability \Rightarrow generalization bounds⁴, predictive coverage⁵

⁴Bousquet & Elisseeff 2002

⁵Steinberger & Leeb 2018

Out-of-sample stability:

$$\mathbb{P} \left\{ \left| \hat{\mu}(X_{n+1}) - \hat{\mu}_{-i}(X_{n+1}) \right| \leq \epsilon \right\} \geq 1 - \nu$$

In-sample stability:

$$\mathbb{P} \left\{ \left| \hat{\mu}(X_i) - \hat{\mu}_{-i}(X_i) \right| \leq \epsilon \right\} \geq 1 - \nu$$

Stability \Rightarrow generalization bounds⁴, predictive coverage⁵

- Example: K -nearest-neighbors satisfies out-of-sample stability with $\epsilon = 0$ and $\nu = K/n$

⁴Bousquet & Elisseeff 2002

⁵Steinberger & Leeb 2018

Theorem: With out-of-sample stability, for jackknife,

$$\mathbb{P} \{ \text{dist}(Y_{n+1}, C(X_{n+1})) \leq \epsilon \} \geq 1 - \alpha - 2\sqrt{\nu}$$

For jackknife+,

$$\mathbb{P} \{ \text{dist}(Y_{n+1}, C(X_{n+1})) \leq 2\epsilon \} \geq 1 - \alpha - 4\sqrt{\nu}$$

If we also assume in-sample stability, then for the naive method,

$$\mathbb{P} \{ \text{dist}(Y_{n+1}, C(X_{n+1})) \leq 2\epsilon \} \geq 1 - \alpha - 4\sqrt{\nu}$$

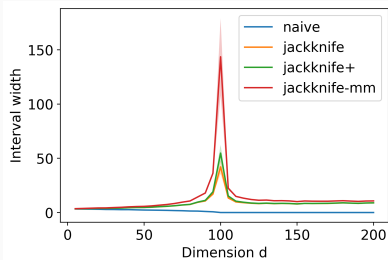
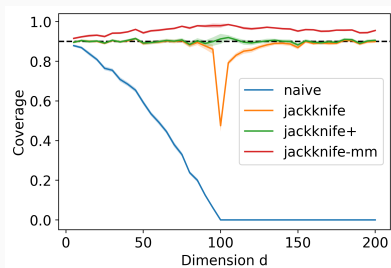
Summary of theory

Method	Assumption-free theory	Out-of-sample stability	In-sample and out-of-sample stab.
Naive	No guarantee	No guarantee	$\approx 1 - \alpha$
Jackknife	No guarantee	$\approx 1 - \alpha$	$\approx 1 - \alpha$
Jackknife+	$1 - 2\alpha$	$\approx 1 - \alpha$	$\approx 1 - \alpha$
Jackknife-mm	$1 - \alpha$	$1 - \alpha$	$1 - \alpha$

- $n = 100, d = 5, 10, \dots, 200$
- $X_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$
- $Y_i = X_i^\top \beta + \mathcal{N}(0, 1)$
- Regression method \mathcal{A} :
 - Least squares with min ℓ_2 norm (ridge with penalty 0)
 - Stable⁶ if $d \ll n$ or $d \gg n$, unstable if $d \approx n$

⁶Hastie et al 2019, *Ridgeless Least Squares Interpolation*.

Simulation



Real data

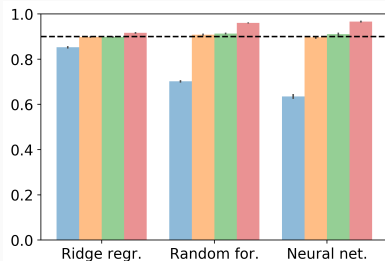
Data set	# samples	# features
Communities & crime	1994	99
BlogFeedback	52397	280
Medical Expenditures Panel	33005	107

- Training data size $n = 200$
- Algorithms: ridge regression, random forests, neural nets

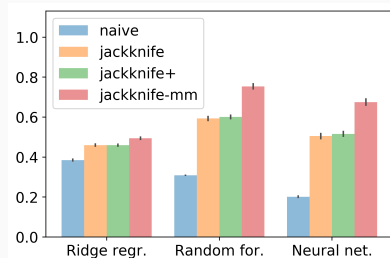
Real data

Communities & crime data:

Coverage



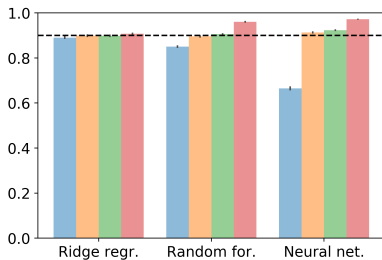
Interval width



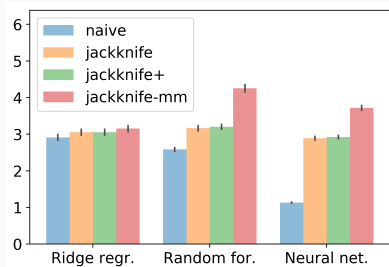
Real data

BlogFeedback data:

Coverage



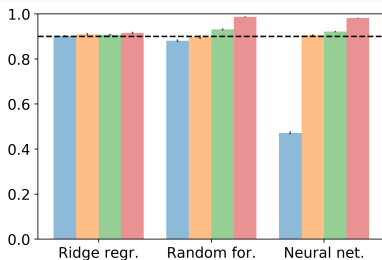
Interval width



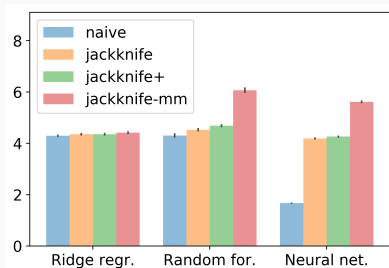
Real data

Medical Expenditure Panel data:

Coverage



Interval width



Summary & open questions

- Jackknife can fail to cover if the regression is unstable
- Jackknife+ accounts for instability in $\hat{\mu}$
to gain assumption-free $1 - 2\alpha$ coverage
- If $\hat{\mu}$ is stable, jackknife+ \approx jackknife (both have $1 - \alpha$ cov.)

- What algorithms/data sets cause jackknife to fail in practice?
- How to use jackknife+ or CV+ for model selection / tuning?