

How Good is my Learning Algorithm? Building Cross-Validation Confidence Intervals for Test Error

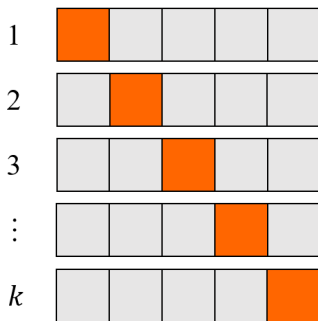
Alexandre Bayle

Department of Statistics, Harvard University

Joint work with Pierre Bayle, Lucas Janson, Lester Mackey

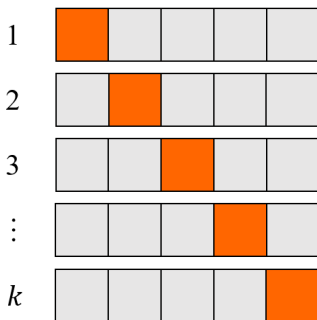
Stat 310 – April 3, 2024

Cross-validation (CV)



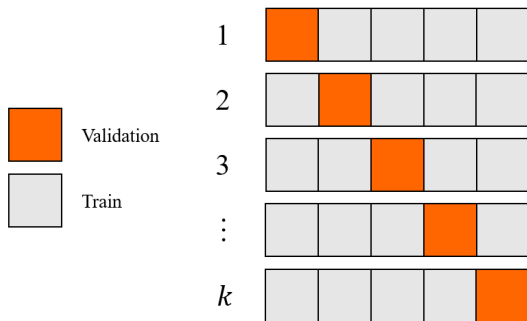
- Divide data into k folds

Cross-validation (CV)



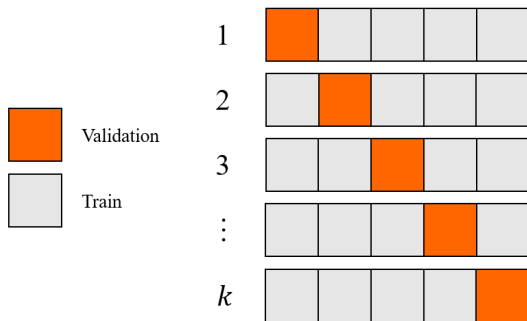
- Divide data into k folds
- Fit k prediction rules, each with one fold held out

Cross-validation (CV)



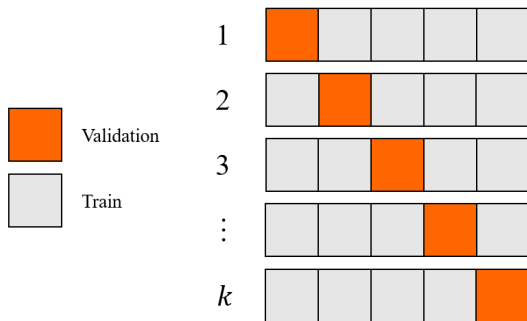
- Divide data into k folds
- Fit k prediction rules, each with one fold held out
- Evaluate each prediction rule on its held-out set

Cross-validation (CV)



- Divide data into k folds
- Fit k prediction rules, each with one fold held out
- Evaluate each prediction rule on its held-out set
- Average the k error estimates

Cross-validation (CV)



- Divide data into k folds
- Fit k prediction rules, each with one fold held out
- Evaluate each prediction rule on its held-out set
- Average the k error estimates

Pros:

- Unbiased for test error
- Lower variance than single train-test split

**How good is my
learning algorithm?**

How good is my learning algorithm?

Need: Test error confidence intervals to quantify uncertainty

How good is my learning algorithm?

Need: Test error confidence intervals to quantify uncertainty

Problem: Existing intervals often invalid & CV distribution is complex

Is algorithm A actually better than algorithm B?

Is algorithm A actually better than algorithm B?

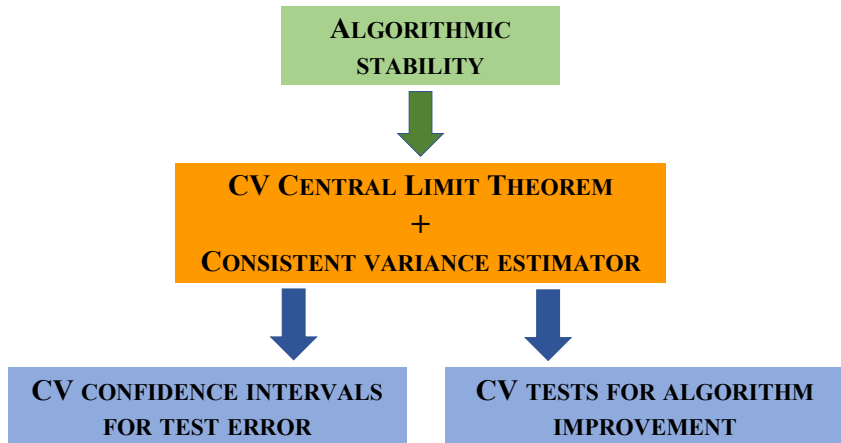
Need: Trustworthy hypothesis tests of algorithm improvement

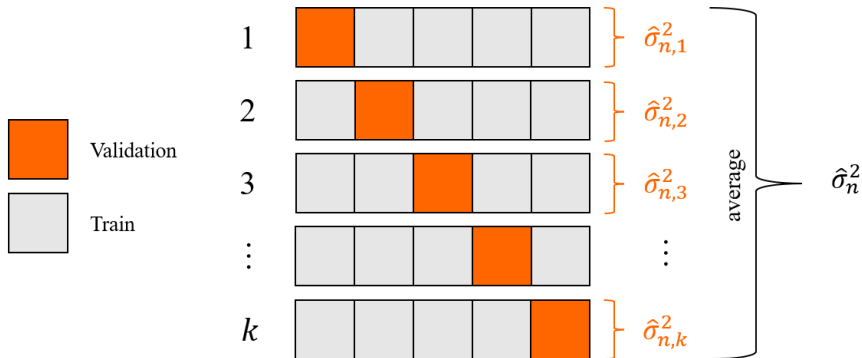
Is algorithm A actually better than algorithm B?

Need: Trustworthy hypothesis tests of algorithm improvement

Problem: Standard tests (like the cross-validated t -test, the repeated train-validation t -test, and the 5×2 -fold CV test) **do not appropriately account for dependence and have no correctness guarantees**

Our Contributions





$$Z = \frac{\sqrt{n} \cdot \hat{R}_n}{\hat{\sigma}_n}$$

(Bayle, Bayle, Janson, and Mackey, 2020)

Hold-out CLT



Validation



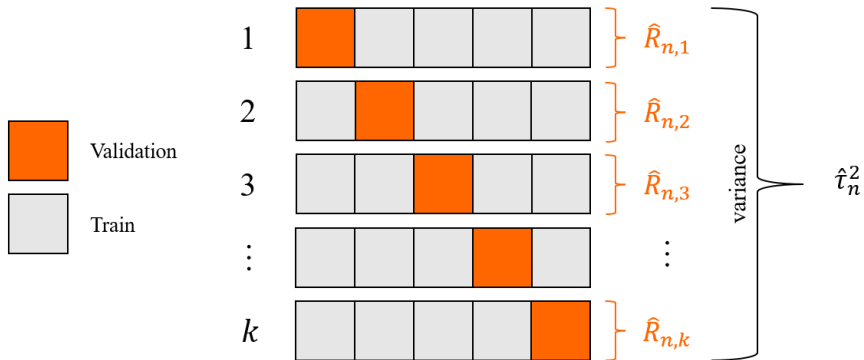
Train

Single split



} $\hat{\sigma}_n^2$
 \hat{R}_n

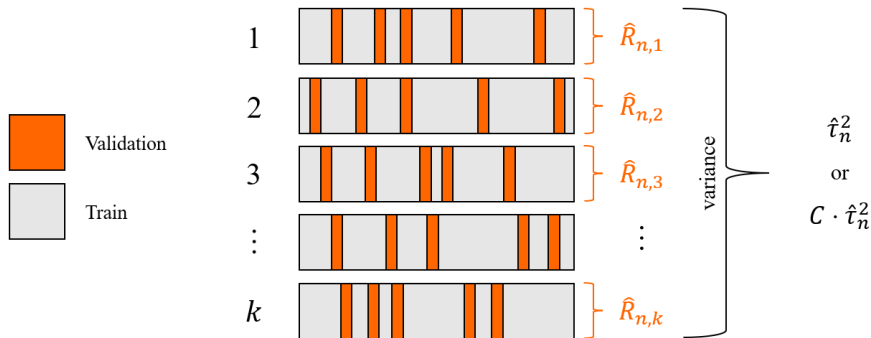
$$Z = \frac{\sqrt{n/k} \cdot \hat{R}_n}{\hat{\sigma}_n}$$



$$t = \frac{\frac{1}{k} \sum_{i=1}^k \hat{R}_{n,i}}{\hat{t}_n / \sqrt{k}}$$

(Dietterich, 1998)

(Corrected) repeated train-validation t



$$t = \frac{\frac{1}{k} \sum_{i=1}^k \hat{R}_{n,i}}{\hat{t}_n / \sqrt{k}} \quad C = 1 + k \frac{n_{\text{test}}}{n_{\text{train}}}$$

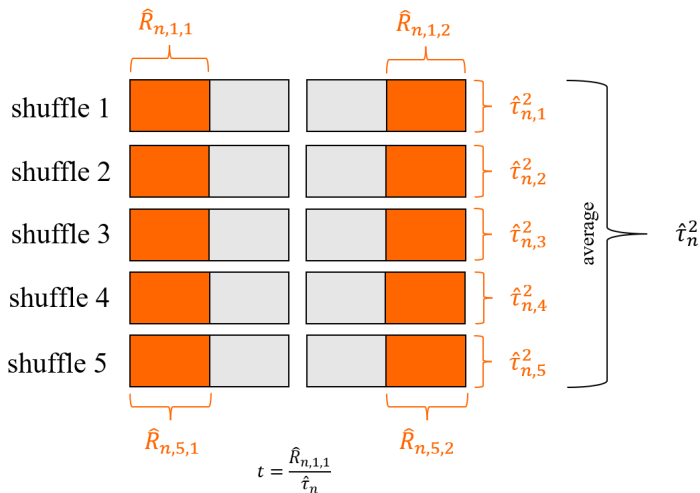
(Nadeau and Bengio, 2003)

5x2 CV



Validation

Train



(Dietterich, 1998)

Notations and Problem Setup

- **Datapoints** Z_1, \dots, Z_n
 - Often each $Z_i = (X_i, Y_i)$ with covariates X_i and response Y_i
 - For any vector B of indices, Z_B denotes the corresponding vector of datapoints

Notations and Problem Setup

- **Datapoints** Z_1, \dots, Z_n
 - Often each $Z_i = (X_i, Y_i)$ with covariates X_i and response Y_i
 - For any vector B of indices, Z_B denotes the corresponding vector of datapoints
- **Loss function** $h_n(Z_i, Z_B)$: error when training on Z_B and testing on Z_i

Notations and Problem Setup

- **Datapoints** Z_1, \dots, Z_n
 - Often each $Z_i = (X_i, Y_i)$ with covariates X_i and response Y_i
 - For any vector B of indices, Z_B denotes the corresponding vector of datapoints
- **Loss function** $h_n(Z_i, Z_B)$: error when training on Z_B and testing on Z_i
 - **Regression:** $h_n(Z_i, Z_B) = (Y_i - \hat{f}(X_i; Z_B))^2$ for $\hat{f}(\cdot; Z_B)$ trained on Z_B

Notations and Problem Setup

- **Datapoints** Z_1, \dots, Z_n
 - Often each $Z_i = (X_i, Y_i)$ with covariates X_i and response Y_i
 - For any vector B of indices, Z_B denotes the corresponding vector of datapoints
- **Loss function** $h_n(Z_i, Z_B)$: error when training on Z_B and testing on Z_i
 - **Regression:** $h_n(Z_i, Z_B) = (Y_i - \hat{f}(X_i; Z_B))^2$ for $\hat{f}(\cdot; Z_B)$ trained on Z_B
 - **Classification:** $h_n(Z_i, Z_B) = \mathbb{1}[Y_i \neq \hat{f}(X_i; Z_B)]$

Notations and Problem Setup

- **Datapoints** Z_1, \dots, Z_n
 - Often each $Z_i = (X_i, Y_i)$ with covariates X_i and response Y_i
 - For any vector B of indices, Z_B denotes the corresponding vector of datapoints
- **Loss function** $h_n(Z_i, Z_B)$: error when training on Z_B and testing on Z_i
 - **Regression:** $h_n(Z_i, Z_B) = (Y_i - \hat{f}(X_i; Z_B))^2$ for $\hat{f}(\cdot; Z_B)$ trained on Z_B
 - **Classification:** $h_n(Z_i, Z_B) = \mathbb{1}[Y_i \neq \hat{f}(X_i; Z_B)]$
- **Validation sets** $\{B'_j\}_{j=1}^k$ and associated **training sets** $\{B_j\}_{j=1}^k$
 - Validation sets partition datapoint indices $\{1, \dots, n\}$ into k folds
 - k can be fixed or grow with n

Notations and Problem Setup

- **Datapoints** Z_1, \dots, Z_n
 - Often each $Z_i = (X_i, Y_i)$ with covariates X_i and response Y_i
 - For any vector B of indices, Z_B denotes the corresponding vector of datapoints
- **Loss function** $h_n(Z_i, Z_B)$: error when training on Z_B and testing on Z_i
 - **Regression:** $h_n(Z_i, Z_B) = (Y_i - \hat{f}(X_i; Z_B))^2$ for $\hat{f}(\cdot; Z_B)$ trained on Z_B
 - **Classification:** $h_n(Z_i, Z_B) = \mathbb{1}[Y_i \neq \hat{f}(X_i; Z_B)]$
- **Validation sets** $\{B'_j\}_{j=1}^k$ and associated **training sets** $\{B_j\}_{j=1}^k$
 - Validation sets partition datapoint indices $\{1, \dots, n\}$ into k folds
 - k can be fixed or grow with n

Cross-validation (CV) error

$$\hat{R}_n = \frac{1}{n} \sum_{j=1}^k \sum_{i \in B'_j} h_n(Z_i, Z_{B_j})$$

Why CV Error?

Cross-validation error: $\hat{R}_n = \frac{1}{n} \sum_{j=1}^k \sum_{i \in B'_j} h_n(\mathbf{Z}_i, \mathbf{Z}_{B_j})$

Why CV Error?

Cross-validation error: $\hat{R}_n = \frac{1}{n} \sum_{j=1}^k \sum_{i \in B'_j} h_n(\mathbf{Z}_i, \mathbf{Z}_{B_j})$

k -fold test error: $R_n = \frac{1}{n} \sum_{j=1}^k \sum_{i \in B'_j} \mathbb{E}[h_n(\mathbf{Z}_i, \mathbf{Z}_{B_j}) \mid \mathbf{Z}_{B_j}]$
 $= \frac{1}{k} \sum_{j=1}^k \mathbb{E}[h_n(\mathbf{Z}_0, \mathbf{Z}_{B_j}) \mid \mathbf{Z}_{B_j}]$

- Average test error of the k prediction rules $\hat{f}(\cdot; \mathbf{Z}_{B_j})$
- Common inferential target

Why CV Error?

Cross-validation error: $\widehat{R}_n = \frac{1}{n} \sum_{j=1}^k \sum_{i \in B'_j} h_n(\mathbf{Z}_i, \mathbf{Z}_{B_j})$

k -fold test error: $R_n = \frac{1}{n} \sum_{j=1}^k \sum_{i \in B'_j} \mathbb{E}[h_n(\mathbf{Z}_i, \mathbf{Z}_{B_j}) \mid \mathbf{Z}_{B_j}]$
 $= \frac{1}{k} \sum_{j=1}^k \mathbb{E}[h_n(\mathbf{Z}_0, \mathbf{Z}_{B_j}) \mid \mathbf{Z}_{B_j}]$

- Average test error of the k prediction rules $\widehat{f}(\cdot; \mathbf{Z}_{B_j})$
- Common inferential target

Goal: Establish a [central limit theorem](#) for $\widehat{R}_n - R_n$

Application: Confidence Intervals for Test Error

Problem

Construct an asymptotically-exact $(1 - \alpha)$ -confidence interval for k -fold test error R_n

Application: Confidence Intervals for Test Error

Problem

Construct an asymptotically-exact $(1 - \alpha)$ -confidence interval for k -fold test error R_n

Solution: CV Confidence Interval for Test Error

If we have a CLT and a variance estimator $\hat{\sigma}_n^2$ that satisfies relative error consistency ($\hat{\sigma}_n^2/\sigma_n^2 \xrightarrow{p} 1$), then the interval

$$C_\alpha = \hat{R}_n \pm q_{1-\alpha/2} \hat{\sigma}_n / \sqrt{n}$$

satisfies

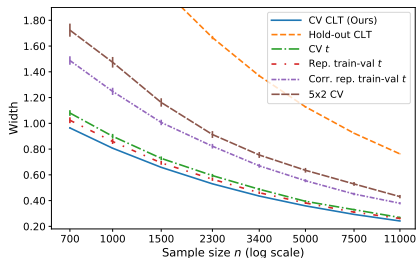
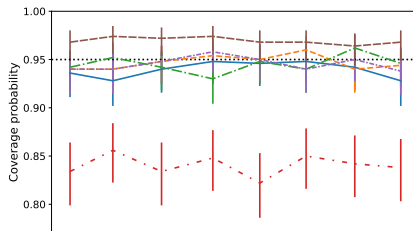
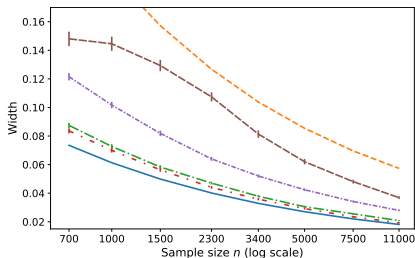
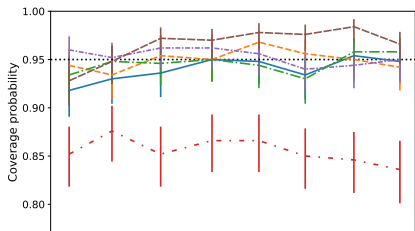
$$\lim_{n \rightarrow \infty} \mathbb{P}(R_n \in C_\alpha) = 1 - \alpha$$

where $q_{1-\alpha/2}$ is the $(1 - \alpha/2)$ -quantile of a standard normal distribution

Confidence Intervals for Test Error, $k=10$

$$C_\alpha = \widehat{R}_n \pm q_{1-\alpha/2} \widehat{\sigma}_n / \sqrt{n} \quad \text{with} \quad 1 - \alpha = 0.95$$

Our CV CLT procedure: valid coverage, smallest width



ℓ^2 -regularized logistic regression (Higgs)

Random forest regression (FlightDelays)

Application: CIs for Test Error Difference

Problem

Construct an asymptotically-exact $(1 - \alpha)$ -confidence interval for the difference in k -fold test errors

Application: CIs for Test Error Difference

Problem

Construct an asymptotically-exact $(1 - \alpha)$ -confidence interval for the difference in k -fold test errors

Solution: CV Confidence Interval for Test Error Difference

For a target loss function ℓ , define the \mathcal{A}_1 - \mathcal{A}_2 loss difference

$$h_n(Z_0, Z_B) = \ell(Y_0, \hat{f}_1(X_0; Z_B)) - \ell(Y_0, \hat{f}_2(X_0; Z_B)),$$

if we have a CLT and a variance estimator $\hat{\sigma}_n^2$ that satisfies relative error consistency ($\hat{\sigma}_n^2/\sigma_n^2 \xrightarrow{P} 1$), then the interval

$$C_\alpha = \hat{R}_n^{(1)} - \hat{R}_n^{(2)} \pm q_{1-\alpha/2} \hat{\sigma}_n / \sqrt{n}$$

satisfies

$$\lim_{n \rightarrow \infty} \mathbb{P}(R_n^{(1)} - R_n^{(2)} \in C_\alpha) = 1 - \alpha$$

Application: Tests for Algorithm Improvement

Problem

Construct an asymptotically-exact level α test of whether \mathcal{A}_1 has smaller k -fold test error than \mathcal{A}_2

Application: Tests for Algorithm Improvement

Problem

Construct an asymptotically-exact level α test of whether \mathcal{A}_1 has smaller k -fold test error than \mathcal{A}_2

Solution: CV Test for Improved Test Error

For a target loss function ℓ , define the \mathcal{A}_1 - \mathcal{A}_2 loss difference

$$h_n(Z_0, Z_B) = \ell(Y_0, \hat{f}_1(X_0; Z_B)) - \ell(Y_0, \hat{f}_2(X_0; Z_B)),$$

and consider testing $H_0 : R_n \geq 0$ (\mathcal{A}_1 not better) against $H_1 : R_n < 0$ (\mathcal{A}_1 is better). If we have a CLT and a variance estimator $\hat{\sigma}_n^2$ that satisfies relative error consistency ($\hat{\sigma}_n^2/\sigma_n^2 \xrightarrow{P} 1$), then the test

$$\text{REJECT } H_0 \Leftrightarrow \hat{R}_n < q_\alpha \hat{\sigma}_n / \sqrt{n}$$

has asymptotic level α for q_α the α -quantile of a standard normal distribution

Tests for Algorithm Improvement, $k=10$, $\alpha=0.05$

Our CV CLT procedure: valid size, most powerful

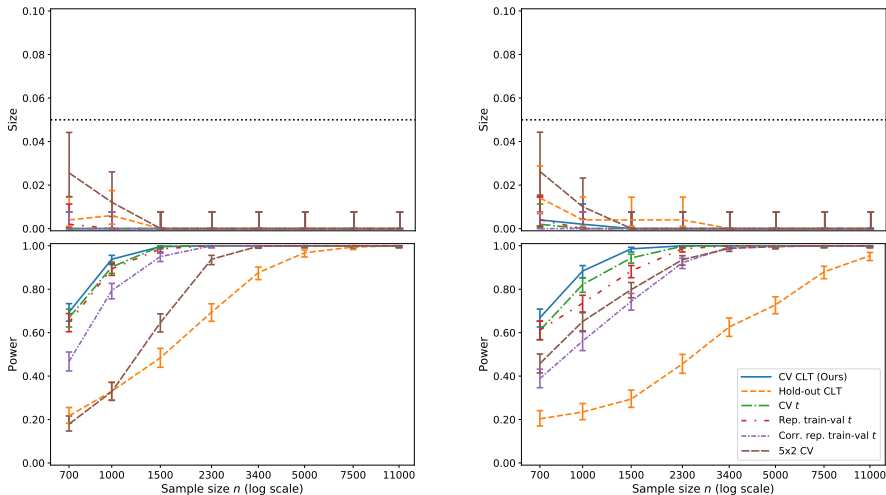


Figure: Size when testing $H_1 : \text{Err}(\mathcal{A}_1) < \text{Err}(\mathcal{A}_2)$ (top) and power when testing $H_1 : \text{Err}(\mathcal{A}_2) < \text{Err}(\mathcal{A}_1)$ (bottom) of level-0.05 tests for improved test error.
Left (classif.): $\mathcal{A}_1 = \text{logistic}$, $\mathcal{A}_2 = \text{NN}$. **Right** (reg.): $\mathcal{A}_1 = \text{RF}$, $\mathcal{A}_2 = \text{ridge}$.

Stability

How much does the performance of a learned prediction rule change when one point in the training set is changed?

Stability

How much does the performance of a learned prediction rule change when one point in the training set is changed?

- Uniform stability: worst-case change in loss h_n

Stability

How much does the performance of a learned prediction rule change when one point in the training set is changed?

- Uniform stability: worst-case change in loss h_n
- Mean-square stability: mean-square change in loss h_n

Stability

How much does the performance of a learned prediction rule change when one point in the training set is changed?

- Uniform stability: worst-case change in loss h_n
- Mean-square stability: mean-square change in loss h_n
- **Loss stability**
 - Mean-square change in loss *difference*
$$h_n(Z_0, Z_{1:m}) - \mathbb{E}[h_n(Z_0, Z_{1:m}) \mid Z_{1:m}]$$

Stability

How much does the performance of a learned prediction rule change when one point in the training set is changed?

- Uniform stability: worst-case change in loss h_n
- Mean-square stability: mean-square change in loss h_n
- **Loss stability**
 - Mean-square change in loss *difference*

$$h_n(Z_0, Z_{1:m}) - \mathbb{E}[h_n(Z_0, Z_{1:m}) \mid Z_{1:m}]$$

Note: $\gamma_{\text{loss}}(h_n) \leq \gamma_{\text{ms}}(h_n)$ [Kumar et al., 2013]

[Bousquet and Elisseeff, 2002, Kale et al., 2011, Kumar et al., 2013, Celisse and Guedj, 2016, ...]

Asymptotic Normality of CV

CV Central Limit Theorem (Bayle, Bayle, Janson, and Mackey, 2020)

Suppose Z_0, Z_1, \dots, Z_n are i.i.d., and define the expected loss function

$$\bar{h}_n(Z_0) = \mathbb{E}[h_n(Z_0, Z_{1:n(1-1/k)}) \mid Z_0] \quad \text{with} \quad \sigma_n^2 = \text{Var}(\bar{h}_n(Z_0)).$$

Asymptotic Normality of CV

CV Central Limit Theorem (Bayle, Bayle, Janson, and Mackey, 2020)

Suppose Z_0, Z_1, \dots, Z_n are i.i.d., and define the expected loss function

$$\bar{h}_n(Z_0) = \mathbb{E}[h_n(Z_0, Z_{1:n(1-1/k)}) \mid Z_0] \quad \text{with} \quad \sigma_n^2 = \text{Var}(\bar{h}_n(Z_0)).$$

If **loss stability** = $o(\sigma_n^2/n)$

Asymptotic Normality of CV

CV Central Limit Theorem (Bayle, Bayle, Janson, and Mackey, 2020)

Suppose Z_0, Z_1, \dots, Z_n are i.i.d., and define the expected loss function

$$\bar{h}_n(Z_0) = \mathbb{E}[h_n(Z_0, Z_{1:n(1-1/k)}) \mid Z_0] \quad \text{with} \quad \sigma_n^2 = \text{Var}(\bar{h}_n(Z_0)).$$

If **loss stability** $= o(\sigma_n^2/n)$ and $(\bar{h}_n(Z_0) - \mathbb{E}[\bar{h}_n(Z_0)])^2/\sigma_n^2$ is uniformly integrable

Sufficient condition: $\sup_n \mathbb{E}[|\bar{h}_n(Z_0) - \mathbb{E}[\bar{h}_n(Z_0)]|^\alpha / \sigma_n^\alpha] < \infty$ for some $\alpha > 2$

Asymptotic Normality of CV

CV Central Limit Theorem (Bayle, Bayle, Janson, and Mackey, 2020)

Suppose Z_0, Z_1, \dots, Z_n are i.i.d., and define the expected loss function

$$\bar{h}_n(Z_0) = \mathbb{E}[h_n(Z_0, Z_{1:n(1-1/k)}) \mid Z_0] \quad \text{with} \quad \sigma_n^2 = \text{Var}(\bar{h}_n(Z_0)).$$

If **loss stability** $= o(\sigma_n^2/n)$ and $(\bar{h}_n(Z_0) - \mathbb{E}[\bar{h}_n(Z_0)])^2/\sigma_n^2$ is **uniformly integrable** then

$$\frac{\sqrt{n}}{\sigma_n}(\hat{R}_n - R_n) \xrightarrow{d} \mathcal{N}(0, 1).$$

Sufficient condition: $\sup_n \mathbb{E}[|\bar{h}_n(Z_0) - \mathbb{E}[\bar{h}_n(Z_0)]|^\alpha / \sigma_n^\alpha] < \infty$ for some $\alpha > 2$

Asymptotic Normality of CV

CV Central Limit Theorem (Bayle, Bayle, Janson, and Mackey, 2020)

Suppose Z_0, Z_1, \dots, Z_n are i.i.d., and define the expected loss function

$$\bar{h}_n(Z_0) = \mathbb{E}[h_n(Z_0, Z_{1:n(1-1/k)}) \mid Z_0] \quad \text{with} \quad \sigma_n^2 = \text{Var}(\bar{h}_n(Z_0)).$$

If **loss stability** $= o(\sigma_n^2/n)$ and $(\bar{h}_n(Z_0) - \mathbb{E}[\bar{h}_n(Z_0)])^2/\sigma_n^2$ is **uniformly integrable** then

$$\frac{\sqrt{n}}{\sigma_n}(\hat{R}_n - R_n) \xrightarrow{d} \mathcal{N}(0, 1).$$

Sufficient condition: $\sup_n \mathbb{E}[|\bar{h}_n(Z_0) - \mathbb{E}[\bar{h}_n(Z_0)]|^\alpha / \sigma_n^\alpha] < \infty$ for some $\alpha > 2$

Many learning algorithms enjoy decaying loss stability

(e.g., SGD, ERM, k -NN, decision trees, ensemble methods)

Consistent Variance Estimation

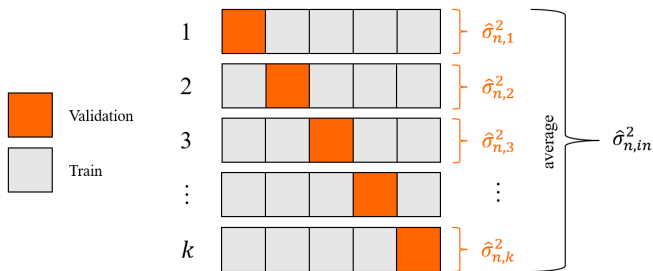
Goal: Find a practical estimator $\hat{\sigma}_n^2$ satisfying $\hat{\sigma}_n^2 / \sigma_n^2 \xrightarrow{p} 1$ under weak conditions.

Consistent Variance Estimation

Goal: Find a practical estimator $\hat{\sigma}_n^2$ satisfying $\hat{\sigma}_n^2 / \sigma_n^2 \xrightarrow{P} 1$ under weak conditions.

Within-fold variance estimator $\hat{\sigma}_{n,in}^2$

Computes the variance of $h_n(Z_i, Z_{B_j})$ in each fold and takes the average across folds



Consistent Variance Estimation

Goal: Find a practical estimator $\hat{\sigma}_n^2$ satisfying $\hat{\sigma}_n^2/\sigma_n^2 \xrightarrow{p} 1$ under weak conditions.

All-pairs variance estimator $\hat{\sigma}_{n,out}^2$

$$\hat{\sigma}_{n,out}^2 = \frac{1}{n} \sum_{j=1}^k \sum_{i \in B_j'} (h_n(Z_i, Z_{B_j}) - \hat{R}_n)^2$$

- Computes the empirical variance of $h_n(Z_i, Z_{B_j})$ across all folds
- **Advantage:** can also be used for leave-one-out cross-validation

Consistent Variance Estimation

Within-fold variance estimator $\hat{\sigma}_{n,in}^2$

$$\hat{\sigma}_{n,in}^2 = \frac{1}{k} \sum_{j=1}^k \frac{1}{(n/k)-1} \sum_{i \in B'_j} \left(h_n(Z_i, Z_{B_j}) - \frac{k}{n} \sum_{i' \in B'_j} h_n(Z_{i'}, Z_{B_j}) \right)^2$$

All-pairs variance estimator $\hat{\sigma}_{n,out}^2$

$$\begin{aligned} \hat{\sigma}_{n,out}^2 &= \frac{1}{n} \sum_{j=1}^k \sum_{i \in B'_j} (h_n(Z_i, Z_{B_j}) - \hat{R}_n)^2 \\ &= \frac{1}{k} \sum_{j=1}^k \frac{k}{n} \sum_{i \in B'_j} (h_n(Z_i, Z_{B_j}) - \hat{R}_n)^2 \end{aligned}$$

Consistent Variance Estimation

Low computational cost

$\hat{\sigma}_{n,in}^2$ and $\hat{\sigma}_{n,out}^2$ can be computed in $O(n)$ time, and if loss is binary, in $O(k)$ and $O(1)$ respectively

When h_n is binary, as in the case of 0-1 loss, one can compute

- $\hat{\sigma}_{n,out}^2 = \hat{R}_n(1 - \hat{R}_n)$ in $O(1)$ time given access to the overall cross-validation error \hat{R}_n ,
- $\hat{\sigma}_{n,in}^2 = \frac{1}{k} \sum_{j=1}^k \frac{(n/k)}{(n/k)-1} \hat{R}_{n,j}(1 - \hat{R}_{n,j})$ in $O(k)$ time given access to the k average fold errors $\hat{R}_{n,j} \triangleq \frac{k}{n} \sum_{i \in B'_j} h_n(Z_i, Z_{B_j})$.

Consistent Variance Estimation

Theorem: Consistency of CV Variance Estimators (Bayle et al.)

Under exactly the same conditions given for the CV central limit theorem (**loss stability** = $o(\sigma_n^2/n)$ and **uniform integrability**), we have

$$\widehat{\sigma}_{n,in}^2 / \sigma_n^2 \xrightarrow{L^1} 1.$$

Consistent Variance Estimation

Theorem: Consistency of CV Variance Estimators (Bayle et al.)

Under exactly the same conditions given for the CV central limit theorem (**loss stability** = $o(\sigma_n^2/n)$ and **uniform integrability**), we have

$$\widehat{\sigma}_{n,in}^2 / \sigma_n^2 \xrightarrow{L^1} 1.$$

If, additionally, **mean-square stability** = $o(k\sigma_n^2/n)$, then

$$\widehat{\sigma}_{n,out}^2 / \sigma_n^2 \xrightarrow{L^1} 1.$$

- **Mean-square stability** condition particularly mild for leave-one-out CV ($k = n$), becomes $o(\sigma_n^2)$

Summary

- New CV central limit theorem under algorithmic stability

Summary

- New CV central limit theorem under algorithmic stability
- Consistent estimators of CV variance

Summary

- New CV central limit theorem under algorithmic stability
- Consistent estimators of CV variance
- Asymptotically exact confidence intervals and tests for k -fold test error

Summary

- New CV central limit theorem under algorithmic stability
- Consistent estimators of CV variance
- Asymptotically exact confidence intervals and tests for k -fold test error

Thank you!