

25. Even More Model Selection (Bootstrap, CV, GCV)

Sunday, March 10, 2024 5:25 PM

Technique #9: Bootstrap Resampling

idea: imitate a validation set

Bootstrap: Given dataset $S = (z_i)_{i=1}^m$ and we have a statistic γ ,

$\gamma(S)$, and we want to estimate distribution of γ (especially the variance)

for $b = 1, \dots, B$

Draw $m' \leq m$ samples from S , with replacement. Call this $S^{(b)}$

Compute $\gamma(S^{(b)})$

$$\text{Then } \widehat{\text{Var}}(\gamma(S)) = \frac{1}{B-1} \sum_{b=1}^B (\gamma(S^{(b)}) - \overline{\gamma(S)}^{\text{sample mean}})^2$$

turns one dataset into many! We're getting something from nothing,
i.e. lifting ourselves by our bootstraps!!

To good to be true? Sort of, not completely.

To use for a criterion to estimate true risk, there's some subtlety on how to do it (otherwise it's like a worse version of cross-validation)

and a 0.632 heuristic. See Hastie, Tibshirani et al. § 7.11

25a. Even More Model Selection (Bootstrap, CV, GCV)a

Sunday, March 10, 2024 5:34 PM

Technique #10: Cross-Validation (CV)

Recall we already talked about (plain) validation

(train/validation/test split, Hoeffding...)

⚠ Don't forget to adjust Hoeffding with $\delta \rightarrow \delta/k$ if you apply to k models simultaneously. (union bound)

Known as the Bonferroni Correction, i.e. divide your p-value by $\frac{\# \text{ of regressors}}{k}$ otherwise you're p-hacking!

(if use 0.05 significance level but run 20 independent test w/ chance results, then 64% chance at least 1 is declared significant)

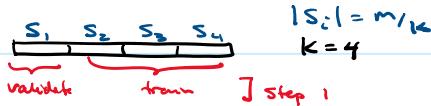
K-fold CV

Better than plain validation

if you don't have much data... but unnecessary if $m \geq 10^5$ or so.

Split data into k -“folds”

for $i = 1, \dots, k$



Train h_i based on training data $S \setminus S_i$

“Validate” h_i using S_i , estimate empirical risk

Return avg. of estimates, $\frac{1}{k} \sum \hat{L}_{S_i}(h_i)$

For a fixed i , this is just plain validation since

$$h_i \perp S_i$$

but since we average over i , ① analysis not easy

② ... but works better than a single plain validation

Typically $k=5$ or 10

$k=1$ doesn't work!

$k=m$ is aka “Leave-one-out” CV ... computationally expensive!

LOO

~ m times more expensive than a single train.

See § 7.10.2 Hastie, Tibshirani et al.

for Right and Wrong ways to do CV

i.e. Feature Selection should be part of CV loop!

Use nice automated CV software frameworks

(Scikit-Learn...) to prevent other kinds of bugs.

25b. Even More Model Selection (Bootstrap, CV, GCV)

Sunday, March 10, 2024 5:52 PM

Technique #11: Generalized CV (GCV)

For Ordinary Least Squares, this is a computationally efficient approximation to LOO CV.

Take linear models $\hat{y} = Py$, P usually parameterized, depends on data X and parameter α

Define $\hat{y}_i = i^{\text{th}}$ output, aka $\hat{y}_i = h_s(x_i)$ estimated model
 and $\hat{y}_i^{(-i)} = i^{\text{th}}$ output if trained on $S^{(-i)} = S \setminus \{x_i\}$ $S = \{x_i\}_{i=1}^m$
 $= h_{S^{(-i)}}(x_i)$ $x_i \in \mathbb{R}^d$

For some models (OLS, cubic smoothing splines)

$$\text{then } \hat{y}_i^{(-i)} - y_i \stackrel{(*)}{=} \frac{\hat{y}_i - y_i}{1 - P_{ii}} \quad P_{ij} = (i, j)^{\text{th}} \text{ entry of } P$$

So LOO CV is

$$CV_{\text{LOO}}(\alpha) = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i^{(-i)} - y_i)^2 \quad \left[\begin{array}{l} \text{Slow! Solve} \\ m \text{ OLS problems} \end{array} \right]$$

Recall
 $P = P(\alpha)$

$$= \frac{1}{m} \sum_{i=1}^m \left(\frac{\hat{y}_i - y_i}{1 - P_{ii}} \right)^2 \quad \left[\begin{array}{l} \text{Fast! Solve} \\ 1 \text{ OLS problem} \end{array} \right]$$

and GCV simplifies even further:

$$GCV(\alpha) = \frac{1}{m} \sum_{i=1}^m \left(\frac{\hat{y}_i - y_i}{1 - \text{tr}(P)/d} \right)^2$$

$$= \underbrace{\left(\frac{1}{1 - \text{tr}(P)/d} \right)^2}_{\text{penalize large effective dim.}} \cdot \underbrace{\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2}_{\text{ERM}}$$

For Tikhonov
 Recall from UPRE section

$$\text{tr}(P_\alpha) = d - \sum_{i=1}^d \frac{\alpha}{\alpha + \lambda_i(X^T X)} \rightarrow 0 \text{ as } \alpha \rightarrow \infty$$

$$\rightarrow d \text{ as } \alpha \rightarrow 0$$

25c. Even More Model Selection (Bootstrap, CV, GCV)

Friday, March 22, 2024 3:04 PM

Prove:

$$\hat{y}_i^{(-i)} - y_i = \frac{\hat{y}_i - y_i}{1 - P_{ii}} \quad (*)$$

$$X X^T = \sum_{i=1}^m x_i x_i^T \quad \text{so} \dots \hat{w}^{(-i)} = (X^T X + \alpha I)^{-1} \cdot (X^T y - x_i y_i)$$

For $Tikhonov$:

$$\hat{y} = X \cdot \hat{w}, \quad \hat{w} = (X^T X + \alpha I)^{-1} X^T y$$

$$\begin{matrix} d \\ m \end{matrix} \boxed{X} \xrightarrow{d} x_i^T$$

$$P = X G X^T$$

$$P_{ij} = x_i^T G x_j$$

Sherman-Morrison-Woodbury matrix inversion lemma

$$(A + U C V^T)^{-1} = A^{-1} - A^{-1} U (C^{-1} + V^T A^{-1} U)^{-1} V^T A^{-1}$$

$$\text{so} \quad G^{(-i)} = (X^T X + \alpha I - x_i x_i^T)^{-1} = G + \frac{1}{1 - x_i^T G x_i} \cdot G x_i x_i^T G$$

hence

$$\begin{aligned} \hat{w}^{(-i)} &= G^{(-i)} (X^T y - x_i y_i) = (G + \frac{1}{1 - P_{ii}} G x_i x_i^T G) (X^T y - x_i y_i) \\ &= G X^T y - G x_i \cdot y_i + \frac{1}{1 - P_{ii}} G x_i x_i^T G X^T y - \frac{1}{1 - P_{ii}} G x_i x_i^T G x_i \cdot y_i \\ &= \hat{w} + G x_i \left(-y_i + \frac{1}{1 - P_{ii}} x_i^T G X^T y - \frac{1}{1 - P_{ii}} x_i^T G x_i \cdot y_i \right) \\ &= \hat{w} + G x_i \left(-y_i + \frac{P_{ii} y_i}{1 - P_{ii}} + \frac{1}{1 - P_{ii}} \hat{y}_i - \frac{1}{1 - P_{ii}} P_{ii} y_i \right) \\ &= \hat{w} + G x_i \cdot \frac{(\hat{y}_i - y_i)}{1 - P_{ii}} \end{aligned}$$

so

$$\hat{y}_i^{(-i)} = x_i^T \hat{w}^{(-i)} = x_i^T \hat{w} + \underbrace{x_i^T G x_i}_{P_{ii}} \cdot \frac{(\hat{y}_i - y_i)}{1 - P_{ii}} = \hat{y}_i + \frac{P_{ii}}{1 - P_{ii}} (\hat{y}_i - y_i)$$

$$\begin{aligned} \text{so } y_i^{(-i)} - y_i &= \hat{y}_i + \frac{P_{ii}}{1 - P_{ii}} (\hat{y}_i - y_i) - y_i = \frac{1}{1 - P_{ii}} (\hat{y}_i - P_{ii} \hat{y}_i + P_{ii} (y_i - y_i) - y_i + P_{ii} y_i) \\ &= \frac{1}{1 - P_{ii}} (\hat{y}_i - y_i) \quad \square \end{aligned}$$

ASIDE: $(A + \epsilon B)^{-1}$

① Is inverse continuous? Yes, if bounded away from being singular, e.g. $A + \epsilon B$ is invertible for $\epsilon \in (-\delta, \delta)$

② how to compute?
 (2a) Sherman-Morrison-Woodbury
 (2b) Neumann-Series

25d. Even More Model Selection (Bootstrap, CV, GCV)

Friday, March 22, 2024 3:27 PM

Neumann Series

Assume A^{-1} is reference (i.e. A nonsingular)

$$(A - \varepsilon B)^{-1} \rightarrow \textcircled{1} \text{ convert to nice form}$$

$$\begin{aligned} (A - \varepsilon B)x = y &\Leftrightarrow A^{-1}(A - \varepsilon B)x = A^{-1}y \\ &\Leftrightarrow I - \varepsilon A^{-1}Bx = A^{-1}\underbrace{y}_{\tilde{y}} \end{aligned}$$

or, if need to preserve symmetry,

$$I - \varepsilon A^{-1/2}BA^{-1/2}\tilde{x} = A^{-1/2}y$$

\textcircled{2} so wlog, look at perturbations of identity

$$(I - \varepsilon B)^{-1} \quad \text{Scalar case: } \frac{1}{1-x} = 1 + x + x^2 + x^3 + \dots \quad \text{if } |x| < 1$$

by analogy

$$(I - \varepsilon B)^{-1} = \sum_{k=0}^{\infty} (\varepsilon B)^k$$

Neumann Series

$$\left(\text{i.e. } \sum_{k=0}^{\infty} r^k = \frac{1-r^{n+1}}{1-r} \xrightarrow{n \rightarrow \infty} \frac{1}{1-r} \right) \quad \text{Calc 1}$$

Converges (in operator norm), even in Banach space!, if

$$\|\varepsilon B\| < 1 \quad \text{i.e. } \varepsilon < \frac{1}{\|B\|}$$

^T Spectral (operator) norm,

not Euclidean / Frobenius

$$\|B\| = \sup_{\|x\|_2=1} \|Bx\|_2$$

= max singular value