

R Programming and Data Analysis

Evaluating Variability and Uncertainty

Introduction

1. Studying properties by simulation
 - a. Sampling distribution of MLE
 - b. Empirical coverage study for a credible interval
2. Dimension Reduction Application
3. Cross-validation
4. Bootstrapping

Sampling Distribution of an Estimator

Sampling Distribution of an Estimator

- The Beta-Binomial (BB) distribution is an extension of Binomial which allows for extra variation.
- $Z \sim \text{Binomial}(m, p)$ measures the number of successes out of m independent success/failure trials, each having success probability p .
- Suppose

$$Y \sim \text{Binomial}(m, p)$$

$$p \sim \text{Beta}(\pi\rho^2(1 - \rho^2), (1 - \pi)\rho^2(1 - \rho^2)),$$

where $a = \pi(1 - \rho^2)/\rho^2$ and $b = (1 - \pi)(1 - \rho^2)/\rho^2$.

- Marginally, $Y \stackrel{\text{iid}}{\sim} \text{BB}(\pi, \rho)$, with density

$$f(y \mid m, \pi, \rho) = \frac{\Gamma(m+1)}{\Gamma(y+1)\Gamma(m-y+1)} \frac{\Gamma(a+y)\Gamma(b+m-y)\Gamma(a+b)}{\Gamma(a+b+m)\Gamma(a)\Gamma(b)},$$

$$E(Y) = m\pi$$

$$\text{Var}(Y) = m\pi(1 - \pi)\{1 + \rho(m - 1)\}$$

Sampling Distribution of an Estimator

- Given fixed m_1, \dots, m_n and a random sample y_1, \dots, y_n from $\text{BB}(\pi, \rho)$, the maximum likelihood estimator (MLE) is obtained by

$$\hat{\theta} = \underset{(\pi, \rho)}{\operatorname{argmax}} \left\{ \log \left[\prod_{i=1}^n f(y_i \mid m_i, \pi, \rho) \right] \right\}.$$

- Large sample theory says that $\hat{\theta} \sim \text{N}(\theta, \mathcal{I}^{-1}(\theta))$, approximately, for large n .
- Problem:** Use simulation study to determine the distribution of $\hat{\theta}$ for a fixed n .

Sampling Distribution of an Estimator

- Use `optim` to maximize the likelihood numerically.
- π and ρ are probabilities, therefore must be constrained to $(0, 1)$.
 - a. `optim` supports simple bound constraints
 - b. `constrOptim` supports linear inequality constraints
 - c. We will transform the problem to an unconstrained optimization.
- Let $G(x) = 1/\{1 + e^{-x}\}$ denote the cdf of the standard logistic distribution. G is a (one-to-one and onto) mapping from \mathbb{R} to $(0, 1)$.
- Transformed problem is

$$\hat{\phi} = \operatorname{argmax}_{(\phi_1, \phi_2)} \left\{ \log \left[\prod_{i=1}^n f(y_i \mid m_i, \pi = G(\phi_1), \rho = G(\phi_2)) \right] \right\}.$$

- MLE for the original space is $\hat{\theta} = (G(\hat{\phi}_1), G(\hat{\phi}_2))$.

Sampling Distribution of an Estimator

... Demonstration ...
(See `sampling.Rmd`)

Empirical Coverage Study for a Credible Interval

Empirical Coverage Study for a Credible Interval

- In a Bayesian analysis, we assume a model for the observed data \mathbf{y} and a distribution for $\boldsymbol{\theta} \in \Theta$ which describes our prior belief,

$$\mathbf{y} \sim f(\mathbf{y} \mid \boldsymbol{\theta}), \quad \boldsymbol{\theta} \sim f(\boldsymbol{\theta}).$$

- All inference on $\boldsymbol{\theta}$ is then based on the posterior distribution

$$f(\boldsymbol{\theta} \mid \mathbf{y}) = \frac{f(\mathbf{y} \mid \boldsymbol{\theta})f(\boldsymbol{\theta})}{\int f(\mathbf{y} \mid \boldsymbol{\vartheta})f(\boldsymbol{\vartheta})d\boldsymbol{\vartheta}}.$$

Empirical Coverage Study for a Credible Interval

- If $Y_i \stackrel{\text{ind}}{\sim} \text{Binomial}(m_i, p)$ and $p \sim \text{Beta}(a, b)$, for given a and b , then

$$\begin{aligned} f(p \mid \mathbf{y}) &\propto \left\{ \prod_{i=1}^n \binom{m_i}{y_i} p^{y_i} (1-p)^{m_i-y_i} \right\} \frac{p^{a-1} (1-p)^{b-1}}{B(a, b)} \\ &\propto p^{\sum_{i=1}^n y_i + a - 1} (1-p)^{\sum_{i=1}^n m_i - \sum_{i=1}^n y_i + b - 1}, \end{aligned}$$

so that $p \mid \mathbf{y} \sim \text{Beta}(\sum_{i=1}^n y_i + a, \sum_{i=1}^n m_i - \sum_{i=1}^n y_i + b)$.

- A level $1 - \alpha$ credible interval for p uses the $\alpha/2$ and $1 - \alpha/2$ quantiles.
- **Problem:** What is the expected width and coverage probability?

Sampling Distribution of an Estimator

... Demonstration ...
(See `credible.Rmd`)

Dimension Reduction Application

Dimension Reduction Application

- We will now consider some dimension reduction methods to use with the diamonds dataset (in the ggplot2 package).
- Xia et al. (2002) propose several methods for dimension reduction in regression models when the form of the regression function is not known.
- Suppose $\{(\mathbf{x}_i, y_i) : i = 1, \dots, n\}$ is an iid sample from

$$Y = \tilde{g}(\mathbf{X}) + \epsilon, \quad E(\epsilon \mid \mathbf{X}) = 0, \quad \mathbf{X} \in \mathbb{R}^p$$

where \tilde{g} is unknown regression function.

- It is often conceivable that, for some $d < p$,

$$\begin{aligned}\tilde{g}(\mathbf{x}) &= g(\mathbf{B}^\top \mathbf{x}) \\ &= g(\mathbf{b}_1^\top \mathbf{x}, \dots, \mathbf{b}_d^\top \mathbf{x}), \quad \mathbf{B} = (\mathbf{b}_1 \cdots \mathbf{b}_d) \in \mathbb{R}^{p \times d}.\end{aligned}$$

Here, \mathbf{B} is the basis of a d -dimensional subspace of \mathbb{R}^p which preserves the regression relationship between y and \mathbf{x} .

Dimension Reduction Application

- The Minimum Average Variance Estimation (MAVE) method (Xia et al., 2002) estimates \mathbf{B} by minimizing the objective function

$$Q(\mathbf{B}) = \sum_{j=1}^n \min_{\alpha_j, \gamma_j} \left\{ \sum_{i=1}^n (y_i - \alpha_j - \gamma_j^\top \mathbf{B}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 w_{ij} \right\}$$

- The model is based on Taylor series expansions

$$\begin{aligned} g(\mathbf{B}^\top \mathbf{x}_i) &\approx \underbrace{g(\mathbf{B}^\top \mathbf{x}_0)}_{\alpha_0} + \underbrace{[\nabla g(\mathbf{B}^\top \mathbf{x}_0)]^\top}_{\gamma_0^\top} (\mathbf{x}_i - \mathbf{x}_0), \quad i = 1, \dots, n, \\ &= \alpha_0 + \gamma_0^\top (\mathbf{x}_i - \mathbf{x}_0), \end{aligned}$$

around a given point \mathbf{x}_0 .

- Given an estimate $\hat{\mathbf{B}}$, a prediction of y_0 for a given \mathbf{x}_0 is

$$\hat{\alpha}_0 = g(\hat{\mathbf{B}}^\top \mathbf{x}_0).$$

Dimension Reduction Application

- MAVE objective function

$$Q(\mathbf{B}) = \sum_{j=1}^n \min_{\alpha_j, \gamma_j} \left\{ \sum_{i=1}^n (y_i - \alpha_j - \gamma_j^\top \mathbf{B}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 w_{ij} \right\}$$

- MAVE uses kernel weights to avoid assumptions on the distribution of \mathbf{X} . Relative to a given \mathbf{x}_0 , and given a bandwidth h , the weights

$$w_{i0} = \frac{K_h(\mathbf{x}_i - \mathbf{x}_0)}{\sum_{\ell=1}^n K_h(\mathbf{x}_\ell - \mathbf{x}_0)}, \quad i = 1, \dots, n,$$

are based on a kernel function $K_h(\mathbf{u}) = p^{-1}K(\mathbf{u}/\sqrt{h})$.

- An example is the Gaussian kernel $K_h(\mathbf{u}) = (2\pi)^{p/2} \exp(-\mathbf{u}^\top \mathbf{u}/2)$.
- Choice of bandwidth h has a major influence on the model.

Dimension Reduction Application

- MAVE objective function

$$Q(\mathbf{B}) = \sum_{j=1}^n \min_{\alpha_j, \gamma_j} \left\{ \sum_{i=1}^n (y_i - \alpha_j - \gamma_j^\top \mathbf{B}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 w_{ij} \right\}$$

- Since \mathbf{B} is the basis of a subspace, $\mathbf{B}^\top \mathbf{B} = \mathbf{I}_{d \times d}$. Furthermore, for any orthogonal matrix \mathbf{A} ,

$$\gamma_j^\top \mathbf{A} \mathbf{A}^\top \mathbf{B}^\top = \gamma_j^\top \mathbf{B}^\top \implies Q(\mathbf{B}\mathbf{A}) = Q(\mathbf{B}).$$

Therefore, $Q(\mathbf{B})$ is an optimization problem over the d -dimensional subspaces in \mathbb{R}^p (the “Grassmann manifold”).

Dimension Reduction Application

- MAVE objective function

$$Q(\mathbf{B}) = \sum_{j=1}^n \min_{\alpha_j, \gamma_j} \left\{ \sum_{i=1}^n (y_i - \alpha_j - \gamma_j^\top \mathbf{B}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 w_{ij} \right\}$$

- Given a particular \mathbf{B} and \mathbf{x}_0 , the inner minimization over

$$Q_{\mathbf{B}}(\alpha_0, \gamma_0) = \sum_{i=1}^n (y_i - \alpha_0 - \gamma_0^\top \mathbf{B}^\top (\mathbf{x}_i - \mathbf{x}_0))^2 w_{i0}$$

corresponds to a local regression estimate at \mathbf{x}_j (Loader, 1999).

- Minimize $Q_{\mathbf{B}}(\alpha_0, \gamma_0)$ in one step via weighted least squares (WLS),

$$\begin{pmatrix} \hat{\alpha}_0 \\ \hat{\gamma}_0 \end{pmatrix} = [\mathbf{Z}^\top \mathbf{W}_0 \mathbf{Z}]^{-1} \mathbf{Z}^\top \mathbf{W}_0 \mathbf{y}, \quad \text{where} \quad \mathbf{Z} = \begin{pmatrix} 1 & (\mathbf{x}_1 - \mathbf{x}_0)^\top \mathbf{B} \\ \vdots & \vdots \\ 1 & (\mathbf{x}_n - \mathbf{x}_0)^\top \mathbf{B} \end{pmatrix},$$

$$\mathbf{W}_0 = \text{Diag}(w_{10}, \dots, w_{n0}), \quad \mathbf{y} = (y_1, \dots, y_n)^\top.$$

Dimension Reduction Application

- Xia et al. (2002) propose an iterative algorithm based on quadratic programming to minimize the objective function

$$Q(\mathbf{B}) = \sum_{j=1}^n \min_{\alpha_j, \gamma_j} \left\{ \sum_{i=1}^n (y_i - \alpha_j - \gamma_j^\top \mathbf{B}^\top (\mathbf{x}_i - \mathbf{x}_j))^2 w_{ij} \right\}$$

- Xia et al. (2002) also propose a simpler method called Outer Product of Gradients (OPG). The rate of consistency of OPG is slower than MAVE, but it works for our purposes and it does not require multiple iterations.
- OPG can provide a good starting value for MAVE.

Outer Product of Gradients (OPG) Algorithm

Inputs: dimension d and bandwidth h .

for $j = 1$ to n **do**

 Compute weights $w_{ij} = \frac{K_h(\mathbf{x}_i - \mathbf{x}_j)}{\sum_{\ell=1}^n K_h(\mathbf{x}_\ell - \mathbf{x}_j)}$ for $i = 1, \dots, n$.

 Compute $(\hat{\alpha}_j, \hat{\gamma}_j)$ by minimizing $Q_{\mathbf{B}}(\alpha_j, \gamma_j)$ via WLS.

end for

return $\hat{\mathbf{B}}$ as matrix of the first d eigenvectors of $\hat{\Sigma} = \frac{1}{n} \sum_{j=1}^n \hat{\gamma}_j \hat{\gamma}_j^\top$.

Dimension Reduction Application

- There are two factors that must be selected by the user: the dimension $d \in \{0, 1, \dots, p\}$ of the reduction, and the bandwidth $h > 0$. We will assume $d = 1$, but select h via **cross-validation**.
- The estimated regression function is based on a random sample, and should be expressed with uncertainty. We will use the **bootstrap** to provide upper and lower bounds to go with the estimated function.

Cross-Validation

Cross-Validation

- In statistics and machine learning, we assume that observed data is generated from an underlying process. The process, not necessarily the data, is our primary interest. But we learn about the process through the data.
- Overfitting occurs when a model is trained to capture all features of the observed data, even spurious ones, but fails to capture important features of the process.
- Using the data at hand, we would like to select a good fitting model which does not overfit.
- Cross-validation (Hastie et al., 2013, §7.10) is a method to assess model fit and ensure that we are not overfitting or oversmoothing the process of interest.

Cross-Validation

- To illustrate cross-validation concretely, consider the following classification problem.
- Suppose $\mathbf{x}_i \in \mathbb{R}^d$ and \mathcal{R} is a subregion of \mathbb{R}^d for $i = 1, \dots, n$. We observe the class assignment

$$y_i = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \mathcal{R}, \\ 0 & \text{otherwise.} \end{cases}$$

Having observed $\mathcal{D} = \{(\mathbf{x}_i, y_i) : i = 1, \dots, n\}$, but without explicitly knowing \mathcal{R} , how to assign the class y^* for a new $\mathbf{x}^* \in \mathbb{R}^d$?

- We could allow some randomness in the observed y_i , but this is not necessary for our illustration.
- The N -Nearest Neighbors algorithm finds the closest N points in \mathcal{D} to \mathbf{x}^* and takes y^* to be the one most frequently occurring in this set (Hastie et al., 2013, Chapter 2).
- Taking N too small will overfit, but N too large will oversmooth.

Cross-Validation

- The most basic way to assess overfitting is to partition the sample into three parts: a **training set**, a **validation set**, and a **test set**.
- For each N , fit the model using observations in the training set, and evaluate the model using observations in the validation set.
- Select the N that gave the best fit on the validation set, and let this be the final model.
- The fit of the final model may be evaluated using the test set.
- This method assumes that there are no systematic differences between the three datasets. This can be avoided by partitioning at random.

Cross-Validation

- K -fold cross-validation makes more efficient use of the data at the expense of more computing.
- Partition the observations indices $\mathcal{S} = \{1, \dots, n\}$ into K subsets $\mathcal{S}_1, \dots, \mathcal{S}_K$.

Algorithm 1 Cross-validation for N -Nearest Neighbors

```
for  $N = 1, \dots, n$  do  
  for  $k = 1, \dots, K$  do  
    Predict classes  $\hat{y}_{CV,i}$  for  $i \in \mathcal{S}_k$  using  $\mathcal{S} \setminus \mathcal{S}_k$  as training set.  
  end for  
  Let  $ERROR_{CV(N)} = \sum_{i=1}^n I(y_i \neq \hat{y}_{CV,i})$   
end for
```

- We could also set aside a test set, which would be left out of cross-validation and reserved for evaluating the final model.

Simple Cross-Validation Example

... Demonstration ...
(See `cv.Rmd`)

Cross-validation

- Now consider bandwidth selection for the OPG method, and suppose we have a set of candidate bandwidths \mathcal{H} .
- Selecting h too small will use too much information from immediate neighbors and overfit.
- Selecting h too large will fail to capture important local features and oversmooth.
- It is reasonable to select h by minimizing a prediction error such as

$$\text{SAPE} = \sum_{i=1}^n |y_i - \hat{y}_i|.$$

- We could compute SAPE for each $h \in \mathcal{H}$ using all the observations, and pick an h that gives the smallest value. But this does not reflect if our selected model has overfit.

Cross-validation

Algorithm 2 Cross-validation for OPG

```
for  $h \in \mathcal{H}$  do  
  for  $k = 1, \dots, K$  do  
    Fit model to obs in  $\mathcal{S} \setminus \mathcal{S}_k$  to estimate  $\hat{B}_{CV}$ .  
    Use  $\hat{B}_{CV}$  to get predictions for the obs in  $\mathcal{S}_k$ .  
  end for  
  Compute  $\text{SAPE}_{CV(h)}$   
end for
```

Cross-Validation with OPG

... Demonstration ...
(See `sim.Rmd` and `diamonds.Rmd`)

Bootstrap

Bootstrap

- Bootstrap is a technique to estimate the distribution of a statistic, and thereby obtain estimates for properties such as variance and bias (Efron and Tibshirani, 1994).
- Generally, let $\mathbf{X} = (X_1, \dots, X_n)$ be a random sample and $\hat{\theta}(\mathbf{X})$ be an estimator for a quantity θ of interest.

Algorithm 3 Nonparametric Bootstrap

```
for  $b = 1, \dots, R$  do  
     $s \leftarrow$  simple random sample with replacement of size  $n$  from  $\{1, \dots, n\}$   
     $\mathbf{X}_{\text{boot}} \leftarrow (X_{s_1}, \dots, X_{s_n})$   
     $\hat{\theta}_r \leftarrow \hat{\theta}(\mathbf{X}_{\text{boot}})$   
end for  
return  $\hat{\theta}_1, \dots, \hat{\theta}_R$ 
```

- Bootstrap is very general and can be used in complicated situations where other methods are not available or are too difficult to apply.
- This flexibility comes at a price; bootstrap becomes heavily computational when $\hat{\theta}(\mathbf{X})$ takes effort to compute.

Explanation of Bootstrap

- From lecture Larry Wasserman's lecture notes (<http://www.stat.cmu.edu/~larry/=stat705/Lecture13.pdf>).
- Suppose $\mathbf{X} = (X_1, \dots, X_n)$ is an iid sample from distribution P , and we want to estimate a functional $T_n(P)$.
- Concretely, suppose $T_n(P) = \text{Var}_P(\hat{\theta}(\mathbf{X}))$ for illustration.
- If we knew P , we could approximate $T_n(P)$ by simulation:
 1. Draw $\mathbf{X}^{(1)}$ from P and compute $\hat{\theta}^{(1)} = \hat{\theta}(\mathbf{X}^{(1)})$.
 - \vdots
 - R . Draw $\mathbf{X}^{(R)}$ from P and compute $\hat{\theta}^{(R)} = \hat{\theta}(\mathbf{X}^{(R)})$.Take s^2 to be the sample variance of $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(R)}$.
- By the strong law of large numbers $s^2 \xrightarrow{\text{a.s.}} T_n(P)$.
- In a data analysis situation with an observed x_1, \dots, x_n , P is unknown.

Explanation of Bootstrap

- Bootstrap uses the empirical distribution P_n of x_1, \dots, x_n to estimate P ;

$$P_n(x) = \begin{cases} \frac{1}{n} & \text{if } x \in \{x_1, \dots, x_n\}, \\ 0 & \text{o.w.} \end{cases}$$

This is the same as drawing a simple random sample with replacement of size n from $\{x_1, \dots, x_n\}$.

- In Bootstrap,

1. Draw $\tilde{\mathbf{X}}^{(1)}$ from P_n and compute $\hat{\theta}^{(1)} = \hat{\theta}(\tilde{\mathbf{X}}^{(1)})$.

\vdots

R . Draw $\tilde{\mathbf{X}}^{(R)}$ from P_n and compute $\hat{\theta}^{(R)} = \hat{\theta}(\tilde{\mathbf{X}}^{(R)})$.

Take s_{boot}^2 to be the sample variance of $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(R)}$.

- By the strong law of large numbers $s_{\text{boot}}^2 \xrightarrow{\text{a.s.}} T_n(P_n)$ as $R \rightarrow \infty$.
- $T_n(P_n)$ estimates $T_n(P)$ as $n \rightarrow \infty$, under some regularity conditions.

Simple Bootstrap Example

... Demonstration ...
(See `boot.Rmd`)

Bootstrap

- The basic MAVE and OPG algorithms do not produce estimates of variability, but these can be obtained via the bootstrap.
- Recall that $\alpha_0 = g(\mathbf{B}^\top \mathbf{x}_0)$ is the quantity we wished to estimate. Let us compute $1 - \delta$ level confidence intervals for each \mathbf{x}_0 in a given set \mathcal{X} .

Algorithm 4 Nonparametric Bootstrap with OPG Algorithm

```
for  $r = 1, \dots, R$  do  
   $s \leftarrow$  simple random sample with replacement of size  $n$  from  $\{1, \dots, n\}$   
   $\mathbf{y}_{\text{boot}} \leftarrow (y_{s_1}, \dots, y_{s_n})$   
   $\mathbf{X}_{\text{boot}} \leftarrow (\mathbf{x}_{s_1} \cdots \mathbf{x}_{s_n})^\top$   
   $\hat{B}_r \leftarrow \text{OPG}(\mathbf{y}_{\text{boot}}, \mathbf{X}_{\text{boot}}, h)$   
   $\hat{\alpha}_r \leftarrow$  estimates  $\hat{\alpha}_0$  for each  $\mathbf{x}_0 \in \mathcal{X}$  based on  $\hat{B}_r$   
end for  
return  $\hat{\alpha}_1, \dots, \hat{\alpha}_R$ 
```

- A bootstrap CI for $\mathbf{x}_0 \in \mathcal{X}$ is obtained using the $\delta/2$ and $1 - \delta/2$ quantiles of $\hat{\alpha}_{10}, \dots, \hat{\alpha}_{R0}$.
- Other bootstrap CIs have been proposed (Efron and Tibshirani, 1994).

Bootstrap with OPG

... Demonstration ...
(See `sim.Rmd` and `diamonds.Rmd`)

References I

Bradley Efron and Robert J. Tibshirani. *An Introduction to the Bootstrap*. Chapman and Hall/CRC, 1994.

Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer, 2nd edition, 2013. URL <http://statweb.stanford.edu/~tibs/ElemStatLearn>.

C. Loader. *Local Regression and Likelihood*. Springer, 1999.

Yingcun Xia, Howell Tong, W. K. Li, and Li-Xing Zhu. An adaptive estimation of dimension reduction space. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):363–410, 2002. doi: 10.1111/1467-9868.03411. URL <http://dx.doi.org/10.1111/1467-9868.03411>.