

Bootstrap validity

Arshak Minasyan

February 5, 2017

1 Introduction

intro and literature review

2 Bootstrap validity

In this section we will discuss the framework of Generalized Linear Models (GLMs) and show that the distribution of $\tilde{\theta} - \theta^*$ can be mimicked using the counterpart from bootstrap world $\tilde{\theta}^b - \tilde{\theta}$, i.e.

$$\mathcal{L}^b(\tilde{\theta}^b - \tilde{\theta} \mid \mathbf{Y}) \approx \mathcal{L}(\tilde{\theta} - \theta^*). \quad (1)$$

2.1 Generalized Linear regression

In this section we derive the necessary terms in general case (without decomposing the initial variable vector v ; for the ease of representation we denote it θ) needed for further analysis for the special class of distributions called exponential class of distributions.

Further we will assume \mathcal{P} to be an exponential family (EF), which has a number of good properties, namely, the log-likelihood function could be written in this way $\ell(v, y) = vy - g(v)$ and

$$L(\theta) = \sum_{i=1}^n \ell(Y_i, f(X_i)) = \sum_{i=1}^n \{Y_i \Psi_i^T \theta - g(\Psi_i^T \theta)\} = \mathbf{Y}^T \Psi \theta - A(\theta) \quad (2)$$

with

$$A(\theta) \stackrel{\text{def}}{=} \sum_{i=1}^n g(\Psi_i^T \theta) \quad (3)$$

Define

$$D^2 \stackrel{\text{def}}{=} \nabla^2 A(\theta) = \sum_{i=1}^n \Psi_i \Psi_i^T g''(\Psi_i^T \theta). \quad (4)$$

We also define the Fisher information matrix

$$\mathbb{F} = D^2 = -\nabla^2 \mathbb{E}L(\theta^*) = \sum_{i=1}^n \Psi_i \Psi_i^T g''(\Psi_i^T \theta^*) = \Phi g''(\Psi^T \theta^*) \Psi^T. \quad (5)$$

We subtract the deterministic part from $L(\theta)$ and get $\zeta(\theta) = L(\theta) - \mathbb{E}L(\theta)$, then

$$\nabla \zeta(\theta) = \sum_{i=1}^n \epsilon_i \Psi_i = \Psi \epsilon, \quad V^2 = \mathbb{V}ar(\nabla \zeta(\theta)) = \Psi \mathbb{V}ar(\epsilon) \Psi^T, \quad (6)$$

where $\epsilon_i = Y_i - \mathbb{E}Y_i$.

The Fisher expansion reads as

$$\|D(\tilde{\theta} - \theta^*) - \xi\| \leq \diamond(x) \quad (7)$$

on a dominating (elliptic) set of probability at least $1 - e^x$, where ξ is defined as follows

$$\xi \stackrel{\text{def}}{=} D^{-1} \nabla L(\theta^*) = D^{-1} \nabla \zeta(\theta^*) = D^{-1} \Psi \epsilon. \quad (8)$$

The latter means that ξ is simply the linear combination of errors ϵ_i . We achieve this result thanks to the structure of likelihood for EFC which stochastic part is linear in \mathbf{Y} , the only source of randomness.

2.2 Bootstrap counterpart

In the bootstrap world we have

$$L^b(\theta) = \sum_{i=1}^n \ell_i(\theta) w_i^b = \sum_{i=1}^n (Y_i \Psi_i^T \theta - g(\Psi_i^T \theta)) w_i^b = \theta^T \Psi \mathcal{W}^b \mathbf{Y} - A^b(\theta) \quad (9)$$

with

$$A^b(\theta) \stackrel{\text{def}}{=} \sum_{i=1}^n g(\Psi_i^T \theta) w_i^b. \quad (10)$$

Note that

$$\mathbb{E}^b A^b(\theta) = A. \quad (11)$$

We define the counterpart of zeta function in bootstrap world as follows

$$\zeta^b(\theta) = \nabla L^b(\theta) - \mathbb{E}^b \nabla L^b(\theta) \quad (12)$$

Simple algebra and the fact that $\nabla \mathbb{E}^b L^b(\tilde{\theta}) = 0$ brings us to the following expression

$$\zeta^b(\theta^*) = \sum_{i=1}^n [Y_i \Psi_i^T - \Psi_i g'(\Psi_i^T \theta^*)] \epsilon_i^b = Y^T \Psi \mathcal{E}^b - \nabla A(\theta^*) \mathcal{E}^b, \quad (13)$$

where $\epsilon_i^b = w_i^b - 1$. The Fisher expansion in bootstrap world has the following form:

$$\|D(\tilde{\theta}^b - \tilde{\theta}) - \xi^b\| \leq \diamond^b(x) \quad (14)$$

One can see that the standartized score in real world is

$$\xi = D^{-1} \Psi \epsilon \quad (15)$$

and the standartized score in bootstrap world is

$$\xi^b = D^{-1} [\mathbf{Y}^T \Psi - \nabla A(\theta^*)] \mathcal{E}^b | \mathbf{Y} \sim \mathcal{N}(0, V^2), \quad (16)$$

where $V^2 \stackrel{\text{def}}{=} \text{Var}(\xi^b)$. The latter is based on the fact that the random components with respect to the measure \mathbb{P}^b are standard normal random variables $\epsilon_i^b | \mathbf{Y} \sim \mathcal{N}(0, 1)$ by construction.

In order to proceed with bootstrap validity we need to compare these two standartized scores.

2.3 Gaussian approximation (GAR) for weighted sum

So far we have used the Fisher expansion in order to approximate the standardized score $D(\tilde{\theta} - \theta^*)$ with normalized score ξ with respect to the measure \mathbb{P} . Analogically we have approximated $D(\tilde{\theta}^b - \tilde{\theta})$ with ξ^b , the counterparts of standardized and normalized scores with respect to \mathbb{P}^b .

ξ^b is indeed Gaussian by construction, hence there is no need to approximate it, while the vector $\xi = D^{-1}\Psi\epsilon$ is the linear combination of the original errors $\epsilon_i \stackrel{\text{def}}{=} Y_i - \mathbb{E}Y_i$. So, denote the Gaussian approximation of ξ as $\bar{\xi}$, which is a centered normal random variable with variance $\text{Var}(\xi)$, i.e.

$$\bar{\xi} \sim \mathcal{N}(0, \text{Var}(\xi)), \quad (17)$$

where $\text{Var}(\xi) = D^{-1}\Psi\text{Var}(\epsilon)\Psi^T D^{-1}$.

To pose the problem in more understandable way we recall that our errors ϵ_i s are independent with zero mean and variance of v_i^2 , i.e.

$$\mathbb{E}\epsilon_i = 0 \quad \text{Var}\epsilon_i = v_i^2. \quad (18)$$

We introduce the independent and normal random variables $\tilde{\epsilon}_i$ with $\text{Var}\tilde{\epsilon}_i = v_i^2$ and define $\bar{\xi}$ as follows

$$\bar{\xi} \stackrel{\text{def}}{=} D^{-1}\Psi\tilde{\epsilon}, \quad (19)$$

where $\tilde{\epsilon} \sim \mathcal{N}(0, \text{diag}(v_1^2, \dots, v_n^2))$. It is easy to check that the covariance matrices of ξ and $\bar{\xi}$ coincide.

Define also S_k as follows

$$S_k \stackrel{\text{def}}{=} \sum_{i=1}^{k-1} c_k \epsilon_k + \sum_{i=k+1}^n c_k \tilde{\epsilon}_k \quad \text{for } k = 2, \dots, n, \quad (20)$$

where c_k are the coefficients of error terms, and

$$S = \sum_{i=1}^n c_k \epsilon_i, \quad \tilde{S} = \sum_{i=1}^n c_k \tilde{\epsilon}_i. \quad (21)$$

Let f is a smooth function with

$$\left| f(x+d) - f(x) - d^T f'(x) - \frac{d^T f''(x) d}{2} \right| \leq C_f \|d\|_\infty^3. \quad (22)$$

Then, using the telescopic principle one can see that

$$f(S) - f(\tilde{S}) = \sum_{k=1}^n (f(S_k + c_k \epsilon_k) - f(S_k + c_k \tilde{\epsilon}_k)), \quad (23)$$

since $S_k + \epsilon_k = S_{k+1} + \tilde{\epsilon}_{k+1}$ for $k = 1, \dots, n-1$.

Now, using Taylor expansion for every k and condition (22) we get $\forall k$

$$\left| f(S_k + c_k \epsilon_k) - f(S_k + c_k \tilde{\epsilon}_k) - c_i^T f'(S_k)(\epsilon_k - \tilde{\epsilon}_k) - c_i^T f''(S_k) c_i \frac{\epsilon_k^2 - \tilde{\epsilon}_k^2}{2} \right| \leq C_f \|c_i\|_\infty^3 \cdot (|\epsilon_k|^3 + |\tilde{\epsilon}_k|^3).$$

Taking the expected value inside the module and using the fact that $\mathbb{E}(\epsilon_i - \tilde{\epsilon}_i) = \mathbb{E}(\epsilon_i^2 - \tilde{\epsilon}_i^2) = 0$ along with independence of S_k and ϵ_k and $\tilde{\epsilon}_k$ one gets

$$|\mathbb{E}f(S) - \mathbb{E}f(\tilde{S})| = \left| \mathbb{E} \left[\sum_{k=1}^n (f(S_k + \epsilon_k) - f(S_k + \tilde{\epsilon}_k)) \right] \right| \leq C_f \|c_i\|_\infty^3 \cdot \mathbb{E} \sum_{k=1}^n (|\epsilon_k|^3 + |\tilde{\epsilon}_k|^3).$$

Let

$$\delta_n = \|c_i\|_\infty^3 \mathbb{E} \sum_{k=1}^n (|\epsilon_k|^3 + |\tilde{\epsilon}_k|^3) = \|c_i\|_\infty^3 \mathbb{E} \sum_{k=1}^n \left(|\epsilon_k|^3 + 2\sqrt{\frac{2}{\pi}} v_i^3 \right), \quad (24)$$

which yields

$$|\mathbb{E}f(S) - \mathbb{E}f(\tilde{S})| \leq C_f \cdot \delta_n. \quad (25)$$

Remark 2.1 *The inequality (25) can be easily extended in this form*

$$|\mathbb{E}f(S - z) - \mathbb{E}f(\tilde{S} - z)| \leq C_f \cdot \delta_n. \quad (26)$$

The basic idea for getting the distance between two distributions from here is to approximate the discontinuous function $f(x) = \mathbb{I}(x \geq q)$ by a smooth function $f(x)$ and then to apply the Lindeberg telescopic sum device.

if we construct such function, then the the theorem will look like something like this:

Theorem 2.1 *Let ϵ_i be independent zero mean with finite third moments and $\text{Var}\epsilon_i = v_i^2$ for $i = 1, \dots, n$. Define*

$$S = \sum_{i=1}^n c_i \epsilon_i \quad \text{and} \quad \tilde{S} = \sum_{i=1}^n c_i \tilde{\epsilon}_i, \quad (27)$$

where $\tilde{\epsilon}_i \sim \mathcal{N}(0, v_i^2)$. Then, for any $\Delta > 0$ it holds

$$\left| \mathbb{P}(S \geq q) - \mathbb{P}(\tilde{S} \geq q + \Delta) \right| \leq \diamond(\Delta, \delta_n, x, p), \quad (28)$$

where δ_n is given by (24), x is our tolerance level and p is the dimension of c_i .

2.4 Gaussian comparison & Pinsker's inequality

As we have already showed that ξ can be replaced by $\bar{\xi}$, we now are ready to compare two Gaussian random vectors, namely

$$\bar{\xi} \sim \mathcal{N}(0, \text{Var}(\xi)), \quad \text{and} \quad \xi^b \sim \mathcal{N}(0, V^2). \quad (29)$$

It is sufficient to compare the covariance matrices of these two distributions. The main idea here relies under the Pinsker's inequality.

Consider two random Gaussian vectors with zero mean, namely $\xi \sim \mathcal{N}(0, \Sigma)$ and $\xi^b \sim \mathcal{N}(0, \Sigma^b)$. Let's discuss this in the most general setting. Let T be a function from \mathbb{R}^p to \mathbb{R}^q and $\mathbf{X} = T(\xi)$ and $\mathbf{Y} = T(\xi^b)$. We are interested in bounding the distance between distributions of \mathbf{X} and \mathbf{Y} .

The next theorem bounds from above the Kullback-Leibler divergence between two normal distributions.

Theorem 2.2 *Let $\mathbb{P} = \mathcal{N}(\mu, \Sigma)$ and $\mathbb{P}^b = \mathcal{N}(\mu^b, \Sigma^b)$ for some non-degenerate symmetric and positive definite matrices. Assume that*

$$\|\Sigma^{-1/2}\Sigma^b\Sigma^{-1/2} - I_p\|_{op} \leq \frac{1}{2} \quad (30)$$

and

$$\text{tr}(\Sigma^{-1/2}\Sigma^b\Sigma^{-1/2} - I_p)^2 \leq \Delta^2, \quad (31)$$

then

$$\mathcal{D}(\mathbb{P} \parallel \mathbb{P}^b) = -\mathbb{E} \log \frac{d\mathbb{P}^b}{d\mathbb{P}} \leq \frac{\Delta^2}{2} + \frac{1}{2}(\mu - \mu^b)^T \Sigma^{-1} \Sigma^b \Sigma^{-1} (\mu - \mu^b). \quad (32)$$

Moreover, for any measurable set $A \subset \mathbb{R}^p$ it holds

$$|\mathbb{P}(A) - \mathbb{P}^b(A)| \leq \sqrt{\frac{\mathcal{D}(\mathbb{P} \parallel \mathbb{P}^b)}{2}}. \quad (33)$$

Proof. By the following change of variable $u = \Sigma^{-1/2}(x - \mu)$ one can achieve the standard normal Gaussian vector \mathbb{P} , while for \mathbb{P}^b we have the following expression

$$\mathbb{P}^b = \mathcal{N}\left(\Sigma^{-1/2}(\mu^b - \mu), \Sigma^{-1/2}\Sigma^b\Sigma^{-1/2}\right) = \mathcal{N}(b, B), \quad (34)$$

where

$$b \stackrel{\text{def}}{=} \Sigma^{-1/2}(\mu^b - \mu), \quad B \stackrel{\text{def}}{=} \Sigma^{-1/2}\Sigma^b\Sigma^{-1/2} \quad (35)$$

Compute

$$2 \log \frac{d\mathbb{P}^b}{d\mathbb{P}}(\gamma) = \log \det(B) - (\gamma - b)^T B (\gamma - b) + \|\gamma\|^2, \quad (36)$$

where γ is standard normal random vector. Then taking the expectation of (36) one can easily get

$$\mathcal{D}(\mathbb{P} \parallel \mathbb{P}^b) = -\mathbb{E} \log \frac{d\mathbb{P}^b}{d\mathbb{P}} = -\frac{1}{2} [\log \det(B) - \text{tr}(B - I_p) - b^T B b]. \quad (37)$$

Let λ_j is the j th eigenvalue of matrix $B - I_p$ then (37) could be rewritten in this way

$$\mathcal{D}(\mathbb{P} \parallel \mathbb{P}^b) = \frac{1}{2} \left[\sum_{j=1}^p \left\{ -\log(\lambda_j + 1) + \lambda_j \right\} + b^T B b \right]. \quad (38)$$

Then, we note that for $x \in (0, 1/2)$ it holds

$$f(x) = x^2 - x + \log(x + 1) \geq 0, \quad \forall x \in \left(0, \frac{1}{2}\right). \quad (39)$$

The simple algebra verifies this. Hence,

$$\frac{1}{2} \left[\sum_{j=1}^p \left\{ -\log(\lambda_j + 1) + \lambda_j \right\} + b^T B b \right] \leq \frac{1}{2} \left[\sum_{j=1}^p \lambda_j^2 + b^T B b \right] \leq \frac{1}{2} b^T B b + \frac{1}{2} \Delta^2, \quad (40)$$

where the latter is obtained using the condition (31). Then substituting back to the original variables we get

$$\mathcal{D}(\mathbb{P} \parallel \mathbb{P}^b) = -\mathbb{E} \log \frac{d\mathbb{P}^b}{d\mathbb{P}} \leq \frac{\Delta^2}{2} + \frac{1}{2} (\mu - \mu^b)^T \Sigma^{-1} \Sigma^b \Sigma^{-1} (\mu - \mu^b) \quad (41)$$

as required.

Applying Pinsker's inequality we obtain

$$\sup_A |\mathbb{P}(A) - \mathbb{P}^b(A)| \leq \sqrt{\frac{\mathcal{D}(\mathbb{P} \parallel \mathbb{P}^b)}{2}} \leq \frac{1}{2} \sqrt{\Delta^2 + (\mu - \mu^b)^T \Sigma^{-1} \Sigma^b \Sigma^{-1} (\mu - \mu^b)} \quad (42)$$

■

Remark 2.2 Note that $1/2$ from the left hand side of (30) can be replaced by some constant ϵ , which will lead to the similar result parametrized on ϵ .

Remark 2.3 Note that for the case of $\mu = \mu^b = 0$ and $T(x) = x$ mapping one obtains

$$\sup_{A \in \mathcal{B}(\mathbb{R}^p)} |\mathbb{P}(\xi \in A) - \mathbb{P}^b(\xi^b \in A)| \leq \sqrt{\frac{\mathcal{D}(\mathbb{P} \parallel \mathbb{P}^b)}{2}} \leq \frac{\Delta}{2}, \quad (43)$$

where $\mathcal{B}(\mathbb{R}^p)$ is a set of all Borel sets in \mathbb{R}^p .

From the above theorem it follows that in order to show that the distributions (i.e. covariance matrices) are close to each other, it is sufficient to show that the conditions of the theorem fulfil.

3 Conditions

Describe the conditions that are needed in order to apply the above described steps.

4 Main

bad structure

We will discuss the case when the error distribution comes from exponential family with canonical parameter (EFc).

We define the oracle of described model as $v^* = \arg \max_v \mathbb{E} L(v, \mathbf{Y})$, while the data-driven estimator for the oracle v^* is defined as $\tilde{v} = \arg \max L(v|Y)$. Note that the source of randomness in $L(v, \mathbf{Y})$ is behind \mathbf{Y} , the distribution of which, in general, is unknown.

We introduce the bootstrap counterpart of (log) likelihood function $L(\cdot)$ as follows

$$L^b(v) = \sum_{i=1}^n \ell_i(v|\mathbf{Y}) \cdot w_i^b, \quad (44)$$

where w_i^b are known as bootstrap weights with $\mathbb{E} w_i^b = 1$, $\text{Var } w_i^b = 1$ and $\mathbb{E} \exp(w_i^b) < \infty$ for all $i \in [1, n]$.

One can see that these two log-likelihood function are very similar to each other, but indeed live in different probability spaces, since the log-likelihood in (47) is the function of random data \mathbf{Y} , while the log-likelihood in (44) is the bootstrap counterpart of log-likelihood and considered to have a different source of randomness, namely, the randomness comes from the *bootstrap multipliers*; \mathbf{Y} is fixed. The link between these two probability spaces is given with this simple observation

$$\mathbb{E}^b L^b(v) \equiv L(v) \quad (45)$$

Suppose $\mathbf{Y} = (Y_1, \dots, Y_n)$ is iid data coming from some unknown distribution \mathcal{P} . In general, $Y_i \sim \mathcal{P}_{f(x)}$, where $f(x)$ is the true function that we are trying to approximate. Consider $X_i \in \mathbb{R}^d$ and $Y_i \sim P_i \in (\mathcal{P}_{\mathbf{v}})$, which means that $\exists v_i : P_i = P_{v_i}$. Function $\mathbf{v}(x)$ can be represented in the following way

$$v(x) = \sum_{j=1}^{\infty} \theta_j \psi_j(x).$$

Then, our linear parametric assumption is

$$v(x) = \sum_{j=1}^{p+q} \theta_j \psi_j(x) = \sum_{j=1}^p \theta_j \psi_j(x) + \sum_{j=p+1}^{p+q} \theta_j \psi_j(x) \text{ for given basis } \psi_j(\cdot). \quad (46)$$

Denote $\eta_i = \theta_{p+i}$ and column-vector $\eta = (\eta_1, \dots, \eta_q)^T \in \mathbb{R}^q$. We will mostly discuss the case of finite p and q .

The log-likelihood function is defined as follows

$$L(v; \mathbf{Y}) = \log \frac{d\mathbb{P}_v}{d\mu_0^n}(\mathbf{Y}) = \sum_{i=1}^n v_i Y_i - g(v_i), \quad (47)$$

where v is the parameter of interest. Define

$$\tilde{v} = \arg \max_{v \in \Theta} L(v, \mathbf{Y}) \quad v^* = \arg \max_{v \in \Theta} \mathbb{E}L(v, \mathbf{Y}) \quad (48)$$

to be maximum likelihood estimator and oracle, respectively.