



Pacini, D. (2020). Proximal Statistic: Asymptotic Normality. *Statistics and Probability Letters*, 167, Article 108896.  
<https://doi.org/10.1016/j.spl.2020.108896>

Peer reviewed version

License (if available):  
CC BY-NC-ND

Link to published version (if available):  
[10.1016/j.spl.2020.108896](https://doi.org/10.1016/j.spl.2020.108896)

[Link to publication record on the Bristol Research Portal](#)  
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Elsevier at <https://doi.org/10.1016/j.spl.2020.108896>. Please refer to any applicable terms of use of the publisher.

## University of Bristol – Bristol Research Portal

### General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available: <http://www.bristol.ac.uk/red/research-policy/pure/user-guides/brp-terms/>

# Proximal Statistic: Asymptotic Normality

David Pacini - University of Bristol

July 2020

*Abstract.* This note introduces an asymptotically normal statistic for the value function of a convex stochastic minimization programme, which may have more than one minimizer. The statistic uses a recursive estimator, based on the proximal algorithm, of one of the minimizers.

Keywords: Set Identification; Proximal algorithm

## 1. The Problem: Nonunique Global Minimizer

Consider the problem of constructing, from a random sample  $\{z_i\}_{i=1}^n$ , an asymptotically normal statistic for  $\varphi_o$  defined as

$$\varphi_o := \inf_{q \in Q} F(q, P_o),$$

where  $Q$  is a known set,  $P_o$  is the unknown distribution of the random vector  $z_i$  taking values in  $\mathbb{R}^L$ , and  $F(q, P_o) := \int f(q, z) dP_o(z)$  for a known real-valued function  $q \mapsto f(q, z_i)$ . This problem arises, for instance, in the context of selecting parametric statistical models, where  $F(q, P_o)$  is the expectation of minus the individual log-likelihood with parameters  $q$  (see e.g., Vuong, 1989), or in regression models with interval-censored variables, where  $F(q, P_o)$  is the support function of a convex set in the direction  $q$  (see e.g., Beresteanu, Molchanov, and Molinari, 2011). In these applications, the asymptotic normality requirement serves to facilitate inference. When  $\arg \min_q F(q, P_o)$  exists and is unique, a solution is the plug-in statistic  $\hat{\varphi}_n := F(\hat{q}_n, P_n)$ , where  $\hat{q}_n \in \arg \min_{q \in Q} F(q, P_n)$  and  $P_n$  is the empirical distribution function. Under an *envelope* and a *Lipschitz-continuity* condition on  $f$ , it is known (see e.g., Shapiro, Dentcheva, and Ruszczyński, 2014, Theorem 5.7) that the sequence  $n^{1/2}(\hat{\varphi}_n - \varphi_o)$  converges in distribution (denoted  $\rightsquigarrow$ ) to a normal random variable  $N(0, \text{avar}(\hat{\varphi}_n))$  with mean zero and variance  $\text{avar}(\hat{\varphi}_n) := E[[f(q_\star, z_i) - \varphi_o]^2]$  for  $q_\star \in \arg \min_{q \in Q} F(q, P_o)$ . When  $\arg \min_q F(q, P_o)$  is not unique, it is also known (see e.g., Shapiro et al., 2014, Theorem 5.7) that  $n^{1/2}(\hat{\varphi}_n - \varphi_o) \rightsquigarrow \mathbb{G}_\star := \inf_{q \in Q_\star} \mathbb{G}_q$ , where  $Q_\star := \arg \min_{q \in Q} F(q, P_o)$  and  $q \mapsto \mathbb{G}_q$  is a Gaussian process. The plug-in statistic is no longer a solution to the problem of interest because  $\mathbb{G}_\star$  is not normal.

This note considers the case when  $Q_\star$  may not be a singleton. It investigates the following statistic.

**Definition (Proximal Statistic).** Let  $\{z_i\}_{i=1}^n$  be i.i.d.  $P_o$ . Let  $Q$  denote the closed unit ball in  $\mathbb{R}^L$  and let  $\mathcal{Z} \subseteq \mathbb{R}^M$  denote the support of the random vector  $z_i$ . Let  $f : Q \times \mathcal{Z} \mapsto \bar{\mathbb{R}}$  be a random convex function such

that: (A.i) There exists a measurable function  $e : \mathbb{R}^L \mapsto \mathbb{R}$ , not depending on  $q$ , such that  $E[e(z_i)]^2 < \infty$  and  $\sup_{q \in Q} |f(q, z_i)| \leq e(z_i)$  a.e.  $z_i$ ; (A.ii) There exists a measurable function  $m : \mathbb{R}^L \mapsto \mathbb{R}_+$  such that  $E[m(z_i)^2] < \infty$  and  $|f(q, z_i) - f(\tilde{q}, z_i)| \leq m(z_i)\|q - \tilde{q}\|$  a.e.  $z_i$  for all  $q, \tilde{q} \in Q$ , where  $\|\cdot\|$  is the Euclidean norm. Define the proximal function

$$prox_P(v) := \arg \min_{q \in Q} F(q, P) + \frac{1}{2}\|q - v\|^2.$$

For  $n > 8$ , define  $k_n := \lceil n^{1/3} \rceil$ . Let  $\hat{q}_{k_n}$  denote the last element in the sequence  $\{\hat{q}_k\}_{k=2}^{k_n}$  defined recursively by

$$\hat{q}_{k+1} := (1 - k^{-1})prox_n(\hat{q}_k), \text{ where } prox_n(\hat{q}_k) := prox_{P_n}(\hat{q}_k) \quad (1)$$

for an arbitrary starting value  $\hat{q}_2 \in Q$ . The *proximal statistic* is  $\hat{\varphi}_{k_n} := F(\hat{q}_{k_n}, P_n)$ .  $\square$

The restrictions on  $f(\cdot, \cdot)$  and  $Q$  in the Definition ensure that  $\varphi_o$  is finite and the proximal statistic is measurable.<sup>1</sup> (A.i) and (A.ii) are, respectively, the envelope and Lipschitz-continuity restrictions. The proximal statistic, unlike the plug-in statistic, uses the proximal estimator  $\hat{q}_{k_n}$ . The recursive scheme (1), leading to the proximal estimator, is an application of the Halpern-type proximal point algorithm introduced by Kamikura and Takahashi (2000).<sup>2</sup> The next section establishes sufficient conditions under which  $\hat{\varphi}_{k_n}$  is asymptotically normal and  $\hat{q}_{k_n}$  converges in probability to the minimum-norm argmin.

## 2. Main Result

**Proposition A (Asymptotic Normality).** *Let (A.i)-(A.ii) hold and suppose that: (A.iii)  $\sup_{v \in Q} \|prox_n(v) - prox_o(v)\| = O_{P_o}(n^{-1/2})$ , where  $prox_o(v) := prox_{P_o}(v)$ . Then,  $n^{1/2}(\hat{\varphi}_{k_n} - \varphi_o) \rightsquigarrow N(0, avar(\hat{\varphi}_{k_n}))$ , where  $avar(\hat{\varphi}_{k_n}) := E[[f(q_\star, z_i) - \varphi_o]^2]$  and  $q_\star \in Q_\star$  is the minimum-norm fixed point of  $v \mapsto prox_o(v)$ .*

The proof is below. (A.iii) is a uniform rate-of-convergence restriction on  $prox_n$ . (A.i)-(A.iii) do not restrict  $Q_\star$  to be a singleton.<sup>3</sup> If  $Q_\star$  is a singleton, the proximal and plug-in statistics will have the same asymptotic normal distribution, c.f., Proposition A with Shapiro et al. (2014, Theorem 5.7). The asymptotic normality of the proximal statistic follows from the result that, even when there are multiple minimizers,  $\hat{q}_{k_n}$ , unlike  $\hat{q}_n$ , converges in probability:

**Lemma 1.**  $\|\hat{q}_{k_n} - q_\star\| \xrightarrow{P_o} 0$ .

The proofs of the Lemmas are in the Appendix. The result in Lemma 1 is achieved, using (A.iii), by setting

$k_n = n^{1/3}$ .<sup>4</sup> This result is new and it may be of independent interest from Proposition A. (A.i), (A.ii), and Lemma 1 serve to establish the following two intermediate results.

**Lemma 2.**  $n^{1/2}[F(q_*, P_n) - \varphi_o] \rightsquigarrow N\left(0, E[[f(q_*, z_i) - \varphi_o]^2]\right)$ .

**Lemma 3.**  $n^{1/2}F(q_*, P_n) - n^{1/2}F(\hat{q}_{k_n}, P_n) \xrightarrow{P_o} 0$ .

*Proof of Proposition A.* Rewrite Lemma 2 as

$$X_n := n^{1/2}[F(q_*, P_n) - \varphi_o] \rightsquigarrow X := N\left(0, E[[f(q_*, z_i) - \varphi_o]^2]\right).$$

For  $Y_n := n^{1/2}[F(\hat{q}_{k_n}, P_n) - \varphi_o]$ , Lemma 3 states that  $X_n - Y_n \xrightarrow{P_o} 0$ . Then, from van der Vaart (1998, Theorem 2.7(iv)), it follows that  $Y_n = n^{1/2}(\hat{\varphi}_{k_n} - \varphi_o) \rightsquigarrow X$ .  $\diamond$

If  $q_*$  also belongs to  $Q_{**} := \arg \min_{q \in Q_*} E[[f(q_*, z_i)^2]]$ ,  $\hat{\varphi}_{k_n}$  has minimum asymptotic variance. A sufficient condition for  $\hat{\varphi}_{k_n}$  having this property is

**Corollary.** *If  $q \mapsto E[f(q, z_i)^2]$  is convex, then  $\text{avar}(\hat{\varphi}_{k_n}) = \min_{q_* \in Q_*} E[[f(q_*, z_i) - \varphi_o]^2]$ .*

One could also construct a minimum-asymptotic-variance statistic by using, for  $a_k := 1 - k^{-1}$ , the iteration  $q_{k+1} = (1 - a_k)\hat{q}_{**} + a_k \text{prox}_n(\hat{q}_k)$  for any consistent estimator  $\hat{q}_{**}$  of  $q_{**} \in \arg \min_{q \in Q_*} E[f(q, z_i)^2]$ .

### 3. Illustration: A Model Selection Test under Loss of Point-Identification

This Section illustrates the proximal statistic in the context of extending Vuong (1989) selection test from non-nested point-identifying models to the set-identifying case. Let  $p_o$  denote the density associated to  $P_o$ . For modeling  $p_o$ , consider the families of parametric density functions,  $\mathcal{G} := \{z \mapsto g(z, \theta) : \theta \in \Theta \subset \mathbb{R}^{\dim(\theta)}\}$  and  $\mathcal{H} := \{z \mapsto h(z, \gamma) : \gamma \in \Gamma \subset \mathbb{R}^{\dim(\gamma)}\}$ . The functions  $z \mapsto g(z, \theta)$  and  $z \mapsto h(z, \gamma)$  are known up to the parameters  $\theta$  and  $\gamma$ , respectively. The aim is to choose the model that is 'closest' to  $p_o$ . Define the Kullback-Liebler information criterion as

$$KLIC_o(\mathcal{G}) := \int \ln p_o(z) dP_o(z) - \min_{\theta \in \Theta} G(\theta, P_o),$$

where  $G(\theta, P_o) := \int -\ln g(z, \theta) dP_o(z)$  is the expectation of minus the individual log-likelihood. A similar definition follows for  $KLIC_o(\mathcal{H})$ .  $KLIC_o(\mathcal{G})$  is non-negative and is equal zero if and only if  $p_o(z_i) = g(z_i, \theta_*)$  a.e.  $z_i$  for  $\theta_* \in \arg \min_{\theta \in \Theta} G(\theta, P_o)$ . Define  $\rho_o := KLIC_o(\mathcal{H}) - KLIC_o(\mathcal{G}) = \min_{\theta \in \Theta} G(\theta, P_o) - \min_{\gamma \in \Gamma} H(\gamma, P_o)$ .

Consider the following hypotheses and definitions:

$H_0 : \rho_o = 0$ , meaning that  $\mathcal{G}$  and  $\mathcal{H}$  are *equivalent*.

$H_G : \rho_o > 0$ , meaning that  $\mathcal{G}$  is *better* than  $\mathcal{H}$ .

$H_{\mathcal{H}} : \rho_o < 0$ , meaning that  $\mathcal{G}$  is *worse* than  $\mathcal{H}$ .

These definitions do not require that either model is point-identifying the model's parameter (i.e., there may be  $\theta_o \neq \tilde{\theta}$  such that  $g(z_i, \theta_o) = g(z_i, \tilde{\theta}) = p_o(z_i)$  a.e.  $z_i$ ).

When  $\arg \max_{\theta} G(\theta, P_o)$  and  $\arg \max_{\gamma} H(\gamma, P_o)$  are unique, it is known (see Vuong, 1989, Theorem 5.1) that, if the models are non-nested, the  $n^{1/2}$ -scaled version of the plug-in statistic  $\rho_n := \min_{\theta \in \Theta} G(\theta, P_n) - \min_{\gamma \in \Gamma} H(\gamma, P_n)$  is, under  $H_0$ , asymptotically normal and, under  $H_{\mathcal{G}}$  (res.  $H_{\mathcal{H}}$ ), diverges to  $+\infty(-\infty)$ . When  $\arg \max_{\theta} G(\theta, P_o)$  and/or  $\arg \max_{\gamma} H(\gamma, P_o)$  are not unique, the asymptotic distribution of  $n^{1/2}\rho_n$ , under  $H_0$ , is the difference between the infima of two Gaussian processes.<sup>5</sup> This asymptotic distribution is not normal. To construct an asymptotic normal statistic, let  $\hat{\varphi}_{g,k_n} := G(\hat{\theta}_{k_n}, P_n)$  denote the proximal statistic for  $\varphi_{g_o} := \min_{\theta \in \Theta} G(\theta, P_o)$ . Let  $\text{avar}_n(\hat{\varphi}_{g,k_n})$  denote the plug-in estimator for the asymptotic variance  $\text{avar}(\hat{\varphi}_{g,k_n}) := E[(\ln g(z_i, \theta_*)^2) - E[\ln g(z_i, \theta_*)]^2]$ . Similarly, define  $\hat{\varphi}_{h,k_n}$ ,  $\hat{\gamma}_{k_n}$ ,  $\text{avar}_n(\hat{\varphi}_{h,k_n})$ , and  $\text{acov}_n(\hat{\varphi}_{g,k_n}, \hat{\varphi}_{h,k_n})$ . Define the test statistic  $\hat{\rho}_{k_n} := \hat{\varphi}_{g,k_n} - \hat{\varphi}_{h,k_n}$  and the standard deviation estimator  $\hat{\omega}_n := [\text{avar}_n(\hat{\varphi}_{g,k_n}) + \text{avar}_n(\hat{\varphi}_{h,k_n}) - 2\text{acov}_n(\hat{\varphi}_{g,k_n}, \hat{\varphi}_{h,k_n})]^{1/2}$ .

**Proposition B (Model Selection Test for Strictly Non-Nested Models).** *Suppose that  $\{z_i\}_{i=1}^n$  is i.i.d.  $P_o$  and (B.i)  $\Theta$  is a compact convex set and the individual log-likelihood function  $\theta \mapsto \ln g(z_i, \theta)$  is a proper concave function a.e.  $z_i$ ; (B.ii) There exists  $e : \mathbb{R}^L \mapsto \mathbb{R}$  such that  $E[e(z_i)^4] < \infty$  and  $\sup_{\theta \in \Theta} |g(z_i, \theta)| \leq e(z_i)$  a.e.  $z_i$ ; (B.iii) There exists  $m : \mathbb{R}^L \mapsto \mathbb{R}$  such that  $E[m(z_i)^4] < \infty$  and  $|\ln g(z_i, \theta) - \ln g(z_i, \tilde{\theta})| \leq m(z_i)\|\theta - \tilde{\theta}\|$  a.e.  $z_i$ ; (B.iv)  $\sup_{v \in \Theta} \|\text{prox}_{g_n}(v) - \text{prox}_{g_o}(v)\| = O_{P_o}(n^{-1/2})$ , where  $\text{prox}_{g_o}(v) := \arg \min_{\theta \in \Theta} G(\theta, P_o) + 1/2\|\theta - v\|^2$ ; (B.v) If  $\Theta$  and  $g(z_i, \cdot)$  are, respectively, replaced by  $\Gamma$  and  $h(z_i, \cdot)$ , (B.i) to (B.iv) hold; (B.vi)  $\mathcal{G} \cap \mathcal{H} = \emptyset$ . Then, under  $H_0$ ,  $n^{1/2}\hat{\rho}_{k_n}/\hat{\omega}_n \rightsquigarrow N(0, 1)$ ; under  $H_{\mathcal{G}}$ ,  $n^{1/2}\hat{\rho}_{k_n}/\hat{\omega}_n \xrightarrow{P_o} +\infty$ ; and under  $H_{\mathcal{H}}$ ,  $n^{1/2}\hat{\rho}_{k_n}/\hat{\omega}_n \xrightarrow{P_o} -\infty$ .*

Proposition B extends to set-identifying models a result in Vuong (1989, Theorem 5.1). It provides an asymptotically pivotal selection test for the models. One chooses a critical value  $c$  from the standard normal distribution for some significance level. If the realized value  $v$  of the statistic  $n^{1/2}\hat{\rho}_{k_n}/\hat{\omega}_n$  is higher than  $c$ , then one rejects the null hypothesis that the models are equivalent in favor of  $\mathcal{G}$ . If  $v$  is smaller than  $-c$ , then one rejects the null hypothesis that the models are equivalent in favor of  $\mathcal{H}$ . Finally, if the absolute value of  $v$  is smaller than  $c$ , one cannot discriminate between the two models given the data. When both models are point-identifying, this test is asymptotically equivalent to the Vuong test.<sup>6</sup>

We decompose the proof of Proposition B in three Lemmas.

**Lemma 4.**  $n^{1/2}(\hat{\varphi}_{gk_n} - \varphi_{g_o}) \rightsquigarrow N(0, \text{avar}(\hat{\varphi}_{gk_n}))$  and  $n^{1/2}(\hat{\varphi}_{hk_n} - \varphi_{h_o}) \rightsquigarrow N(0, \text{avar}(\hat{\varphi}_{hk_n}))$ .

**Lemma 5.**  $\hat{\omega}_n \xrightarrow{P_o} \omega_o$ .

**Lemma 6.** (B.vi) implies  $\omega_o > 0$ .

**Proof of Proposition B.**  $n^{1/2}(\hat{\varphi}_{gk_n} - \varphi_{g_o}) - n^{1/2}(\hat{\varphi}_{hk_n} - \varphi_{h_o}) = n^{1/2}\hat{\rho}_{k_n} - n^{1/2}\rho_o$ . Under  $H_0$ ,  $n^{1/2}\rho_o = 0$ . Lemma 6 justifies to use Slutsky's Lemma (van der Vaart, 1998, Lemma 3.8 (iii)) to combine Lemmas 4 and 5 to obtain  $n^{1/2}\hat{\rho}_{k_n}/\hat{\omega}_n \rightsquigarrow N(0, 1)$ . The claim for  $n^{1/2}\hat{\rho}_{k_n}/\omega_n$  under  $H_G$  and  $H_{\mathcal{H}}$  follows similarly.  $\diamond$

## Endnotes

<sup>1</sup>For the verification of these properties, see the Appendix. The assumption that  $Q$  is the closed unit ball can be relaxed, at the cost of losing conciseness in the exposition, to  $Q$  being a closed convex subset of  $\mathbb{R}^M$  and proving Lemma 1 below by verifying the conditions in Bauschke and Combettes (2017, Theorem 30.1).

<sup>2</sup>For an exposition on the proximal point algorithm, see e.g., Polson, Scott and Willard (2015). The convergence in distribution of  $\hat{\varphi}_{k_n}$  and the convergence in probability of  $\hat{q}_{k_n}$  have so far not been studied.

<sup>3</sup>Possible extensions to Proposition A include studying: (a) the conditions under which the convergence in distribution also holds uniformly; the properties of the proximal statistic when (b)  $q \mapsto f(q, z_i)$  is strongly amenable; (c)  $1/2\|q - \hat{q}_k\|^2$  is replaced by another Bregman divergence; (d)  $Q$  is defined by moment inequality restrictions. These extensions are out of the scope of this note.

<sup>4</sup>Indeed, for any positive integer  $c$  not depending on  $n$ , one can set  $k_n = c + n^{1/3}$ .

<sup>5</sup>This follows from applying Shapiro et al. (2014, Theorem 5.7)

<sup>6</sup>The following extensions to Proposition B are out of the scope of this note. First, the asymptotic approximation in Proposition B is pointwise in  $P_o$ . The development of a uniform asymptotic approximation is needed. Second, one could compare more than two models using multiple testing methods. Third, one could apply Proposition A to a test based on a goodness-of-fit criteria other than the KLIC.

## References

- Andrews, D. (1994): "Empirical Processes in Econometrics", *Handbook of Econometrics*, Vol. 4, Elsevier, Amsterdam.
- Bauschke, H. and P. Combettes (2017): *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer, .
- Beresteanu, A., I. Molchanov and F. Molinari (2011): "Sharp Identification Regions in Models with Convex Moment Predictions", *Econometrica*.
- Halpern, B. (1967) "Fixed Points of Nonexpanding Maps", *Bulletin of the American Mathematical Society*.
- Kamikura, S. and W. Takahashi (2000): "Approximating Solutions of Maximal Monotone Operators in Hilbert Spaces", *Journal of Approximation Theory*.
- Moreau, J.-J. (1965): "Proximite et Dualite dans un Espace Hilbertien", *Bulletin de la Societe Mathematique de France*.

Polson, N., J. Scott and B. Willard (2015): "Proximal Algorithms in Statistics and Machine Learning", *Statistical Science*.

Shapiro, A., D. Dentcheva, and A. Ruszczyński (2014): *Lecture Notes on Stochastic Programming: Modeling and Theory*, Second Edition, SIAM, Philadelphia.

van der Vaart, A. (1998): *Asymptotic Statistics*, Cambridge University Press, Cambridge.

Vuong, Q. (1989): "Likelihood Ratio Tests for Model Selection and Non-Nested Hypotheses", *Econometrica*.

**Acknowledgments.** I would like to thank S. Stouli, A. Santos, P. Lavergne and participants at the First Bristol-Toulouse Econometric Workshop for useful comments.

Appendix (Online Supplementary Material)

We now verify that, under the conditions on  $f(\cdot, \cdot)$  and  $Q$  in the definition of the proximal statistic,  $\hat{\varphi}_{k_n}$  is measurable and  $\varphi_o$  is well-defined, i.e., finite. Under (A.ii),  $q \mapsto f(q, z_i)$  is continuous a.e.  $z_i$ . Since  $f$  is a random function, one has that  $f(q, \cdot)$  is measurable for every  $q$ .  $F(q, P_n)$  is then a Caratheodory function, i.e., a real-valued function such that  $F(q, P_n)$  is measurable for every  $q$  and  $q \mapsto F(q, P_n)$  is continuous a.e.  $z_i$ . It follows (see e.g., Shapiro et al., 2014, Theorem 7.41) that  $F(q, P_n)$  is a random lsc function and so is  $q \mapsto F(q, P_n) + .5\|q - \hat{q}_2\|^2$  because  $q \mapsto .5\|q - \hat{q}_2\|^2$  is continuous. It also follows (see e.g., Shapiro et al., 2014, Theorem 7.42) that  $\hat{q}_3 = (1 - 1/2)prox_{P_n}(\hat{q}_2)$  and  $\hat{\varphi}_3 = F(\hat{q}_3, P_n)$  are both measurable. Hence,  $q \mapsto .5\|q - \hat{q}_3\|^2$  is a random lsc function and so is  $q \mapsto F(q, P_n) + .5\|q - \hat{q}_3\|^2$ . It follows (again from Shapiro et al., 2014, Theorems 7.42 and 7.43) that  $\hat{q}_4 = (1 - 1/3)prox_{P_n}(\hat{q}_3)$  and  $\hat{\varphi}_4 = F(\hat{q}_4, P_n)$  are measurable. By induction, conclude that  $\hat{q}_{k_n}$  and  $\hat{\varphi}_{k_n} := F(\hat{q}_{k_n}, P_n)$  are measurable. Under (A.i) and (A.ii),  $q \mapsto F(q, P_o)$  is a real-valued continuous function (see Shapiro et al., 2014, Theorem 7.48). Since  $Q$  is compact, one is justified then to use the Extreme Value Theorem to claim that  $\varphi_o$  is finite -i.e.  $\min_q F(q, P_o)$  is well-defined. **Proof of Lemma 1.** Assumptions (A.i) and (A.ii) imply (see Shapiro et al., 2014, Theorems 7.48 and 7.49) that  $q \mapsto F(q, P_o)$  is a well-defined and finite valued function. Since we have also assumed that  $q \mapsto f(q, z_i)$  is convex a.e.  $z_i$ , it follows then, from Shapiro et al. (2014, Theorem 7.51), that  $q \mapsto F(q, P_o)$  is a proper convex function. Hence, from Moreau (1965, *Proposition 5.b.*), one has that  $v \mapsto prox_o(v)$  is nonexpansive:

$$\|prox_o(v) - prox_o(\tilde{v})\| \leq \|v - \tilde{v}\| \text{ for any } v, \tilde{v} \in Q. \quad (1.1)$$

Use now the triangle inequality to bound  $\|\hat{q}_{k_n+1} - q_\star\|$  by the sum of a deterministic and a stochastic term  $\|\hat{q}_{k_n+1} - q_\star\| \leq \|q_{k_n+1} - q_\star\| + \|\hat{q}_{k_n+1} - q_{k_n+1}\|$ . Consider first the deterministic term. Let  $a_{k_n} := 1 - k_n^{-1}$ . Define  $q_{k_n+1} := a_{k_n}prox_o(q_{k_n})$  for an arbitrary starting point  $q \in Q$ . Since  $prox_o : Q \mapsto Q$  is nonexpansive (see (1.1))and  $Q$  is the closed unit ball in a Hilbert space,  $a_{k_n} := 1 - k_n^{-1} = 1 - \lceil n^{-1/3} \rceil$  is *acceptable* in the sense of Halpern (1967, Corollary p. 961), viz.

$$\|q_{k_n} - q_\star\| = o(1), \quad (1.2)$$

where  $q_\star$  is the fixed point of  $v \mapsto prox_o(v)$  with the smallest norm. Since the fixed points of  $v \mapsto prox_o(v)$  belong to  $Q_\star$ , a fortiori  $q_\star \in Q_\star$ . One also has  $\|q_{k_n+1} - q_\star\| = o(1)$ .



Consider now the stochastic term. Replacing  $\hat{q}_{k_n+1}$  and  $q_{k_n+1}$  recursively,

$$\|\hat{q}_{k_n+1} - q_{k_n+1}\| = \|a_{k_n} \text{prox}_n(\hat{q}_{k_n}) - a_{k_n} \text{prox}_o(q_{k_n})\|.$$

Add-and-subtract  $\text{prox}_o(\hat{q}_{k_n})$  and use the triangle inequality to get

$$\|\hat{q}_{k_n+1} - q_{k_n+1}\| \leq a_{k_n} \|\text{prox}_n(\hat{q}_{k_n}) - \text{prox}_o(\hat{q}_{k_n})\| + a_{k_n} \|\text{prox}_o(\hat{q}_{k_n}) - \text{prox}_o(q_{k_n})\|.$$

Since  $v \mapsto \text{prox}_o(v)$  is nonexpansive (see (1.1)),

$$\|\hat{q}_{k_n+1} - q_{k_n+1}\| \leq a_{k_n} \|\text{prox}_n(\hat{q}_{k_n}) - \text{prox}_o(\hat{q}_{k_n})\| + a_{k_n} \|\hat{q}_{k_n} - q_{k_n}\|.$$

By recursive substitution,

$$\|\hat{q}_{k_n+1} - q_{k_n+1}\| \leq \frac{a_{k_n}}{1 - a_{k_n}} \|\text{prox}_n(\hat{q}_{k_n}) - \text{prox}_o(\hat{q}_{k_n})\|.$$

Since  $\|\text{prox}_n(\hat{q}_{k_n}) - \text{prox}_o(\hat{q}_{k_n})\| \leq \sup_{q \in Q} \|\text{prox}_n(q) - \text{prox}_o(q)\|$ ,

$$\|\hat{q}_{k_n+1} - q_{k_n+1}\| \leq \frac{a_{k_n}}{1 - a_{k_n}} \sup_{q \in Q} \|\text{prox}_n(q) - \text{prox}_o(q)\|.$$

Since we have assumed that  $n^{1/2}[\text{prox}_n - \text{prox}_o]$  is asymptotically tight (see A.iii), one has

$$\|\hat{q}_{k_n+1} - q_{k_n+1}\| \leq \frac{a_{k_n}}{1 - a_{k_n}} O_{P_o}(n^{-1/2}).$$

Since  $a_{k_n}/(1 - a_{k_n}) = k_n - 1$  and we have assumed  $k_n = o(n^{1/2})$ ,

$$\|\hat{q}_{k_n+1} - q_{k_n+1}\| \leq o(n^{1/2}) O_{P_o}(n^{-1/2}) \leq n^{1/2-1/2} o(1) O_{P_o}(1) \leq o_{P_o}(1).$$

Conclude then  $\|\hat{q}_{k_n} - q_\star\| \leq o(1) + o_{P_o}(1) \leq o_{P_o}(1)$ . △

**Proof of Lemma 2.** For  $\mathcal{F} := \{f(q, \cdot) : q \in Q\}$ , define  $\ell^\infty(\mathcal{F}) := \{f \in \mathcal{F} : \sup_{q \in Q} |f(q, \cdot)| < \infty\}$ .

Let  $H(\epsilon, \mathcal{F}, P)$  denote the cover number of the family of functions  $\mathcal{F}$ .<sup>7</sup> Under (A.ii),  $\mathcal{F}$  is a *type II class* in the sense of Andrews (1994, p. 2270). It follows then from Andrews (1994, Theorem 2) that  $v(z_i) := \max(1, e(z_i), m(z_i))$  is such that  $|f(q, z_i)| \leq v(z_i) \forall f \in \mathcal{F}$  and the uniform entropy integral

$\int_0^1 \sup_{P \in \mathcal{D}} [\ln H(\epsilon(Pv^2)^{1/2}, P, \mathcal{F})]^{1/2} d\epsilon$  satisfies

$$\int_0^1 \sup_{P \in \mathcal{D}} [\ln H(\epsilon(Pv^2)^{1/2}, P, \mathcal{F})]^{1/2} d\epsilon \leq \infty, \quad (2.1)$$

where  $\mathcal{D}$  is the set of all discretely supported probability distributions. Under the bounds  $E[|e(z_i)|^2] < \infty$  in (A.i) and  $E[|m(z_i)|^2] < \infty$  in (A.ii), Jensen's inequality implies

$$P_o v^2 \leq \infty. \quad (2.2)$$

Since  $\mathcal{F}$  is measurable under (A.i) and (A.ii), it follows from (2.1)-(2.2), by van der Vaart (1998, Theorem 19.14), that

$$\mathcal{F} \text{ is } P_o\text{-Donsker}. \quad (2.3)$$

Conclude by restating the definition of  $P_o$ -Donsker class (van der Vaart, 1998, p.269) that  $\mathbb{G}_n f(q) := n^{1/2}[F(q, P_n) - F(q, P_o)] \rightsquigarrow \mathbb{G}f(q)$  in the space  $\ell^\infty(\mathcal{F})$ , where  $q \mapsto \mathbb{G}f(q)$  is a Gaussian process with zero mean and covariance function  $q, \tilde{q} \mapsto E[f(q, z_i)f(\tilde{q}, z_i)] - E[f(q, z_i)]E[f(\tilde{q}, z_i)]$ .  $\triangle$

**Proof of Lemma 3.** We first verify that  $q \mapsto f(q, z_i)$  is square integrable at  $q_*$ :

$$\lim_{q \rightarrow q_*} \int |f(q, z_i) - f(q_*, z)|^2 dP_o(z) = 0. \quad (3.1)$$

For any  $q \in Q$ , (A.ii) implies  $|f(q, z_i) - f(q_*, z_i)|^2 \leq m(z_i)^2 \|q - q_*\|^2$  because  $|f(q, z_i) - f(q_*, z_i)|$  is nonnegative. Taking expectations on both sides

$$\int |f(q, z_i) - f(q_*, z)|^2 dP_o(z) \leq \int m(z)^2 dP_o(z) \|q - q_*\|^2.$$

Under (A.ii),  $\int m(z)^2 dP_o(z) < \infty$ . Hence, (3.1) follows from the last display after taking limits to both sides as  $q \rightarrow q_*$ .

Define  $g : \ell^\infty(\mathcal{F}) \times \mathcal{F} \mapsto \mathbb{R}$  by  $g(h, f) := h(f) - h(f_*)$ , where  $f_* = f(q_*, \cdot)$ . The set  $\mathcal{F}$  is a semi-metric space relative to the  $L_2(P_o)$ -metric. The function  $g$  is continuous with respect to the product semimetric at every point  $(h, f)$  such that  $f \mapsto h(f)$  is continuous. Indeed, if, for any sequence  $\{h_k, f_k\}_k$  in  $\ell^\infty(\mathcal{F}) \times \mathcal{F}$ ,  $\{h_k, f_k\}_k \rightarrow (h, f)$ , then  $h_k \rightarrow h$  uniformly and hence  $h_k(f_k) = h(f_k) + o(1) \rightarrow h(f)$  if  $h$  is continuous at  $f$ .

By van der Vaart (1998, Lemma 28.15), it follows from (2.3) that almost all sample paths of  $\mathbb{G}$  are uniformly continuous on  $\mathcal{F}$ . Thus, the function  $h$  is continuous at  $f_\star \in \mathcal{F}$ .

Set  $f_n := f(\hat{q}_{k_n}, \cdot)$ . Since  $\hat{q}_{k_n} \xrightarrow{P_o} q_\star$  (see Lemma 1), one has, by (3.1), that  $f_n \xrightarrow{P_o} f_\star$  in the metric space  $\mathcal{F}$ . For  $\mathbb{G}_n := n^{1/2}(P_n - P_o)$ , by (2.3),  $\mathbb{G}_n \rightsquigarrow \mathbb{G}$  in the space  $\ell^\infty(\mathcal{F})$ . Hence,

$$(f_n, \mathbb{G}_n) \rightsquigarrow (f_\star, \mathbb{G}) \text{ in the space } \mathcal{F} \times \ell^\infty(\mathcal{F}). \quad (3.2)$$

We have verified that  $g$  is continuous and (3.2) holds. Apply the Continuous Mapping Theorem (van der Vaart, 1998, Theorem 18.11(i)) to obtain

$$\mathbb{G}_n(f_n - f_\star) = g(\mathbb{G}_n, f_n) \rightsquigarrow g(\mathbb{G}, f_\star) = \mathbb{G}f_\star - \mathbb{G}f_\star = 0.$$

Since convergence in probability and convergence in distribution are the same for a degenerate limit (van der Vaart, 1998, Theorem 18.10(iii)),  $\mathbb{G}_n(f_n - f_\star) \xrightarrow{P_o} 0$ . Conclude by replacing  $\mathbb{G}_n(f_n - f_\star)$  by its definition in  $-\mathbb{G}_n(f_n - f_\star) = n^{1/2}F(q_\star, P_n) - n^{1/2}F(\hat{q}_{k_n}, P_n)$ .  $\triangle$

**Proof of Lemma 4.** Under (B.i)-(B.iv), one is justified to set  $q = \theta$ ,  $Q = \Theta$ ,  $f(q, z_i) = -\ln g(z_i, \theta)$ ,  $F(q, P_o) = G(\theta, P_o)$ , etc. It follows then from Proposition A that  $n^{1/2}[\hat{\varphi}_{g, k_n} - \varphi_{g, k_n}] \rightsquigarrow N(0, \text{avar}(\hat{\varphi}_{g, k_n}))$ . A similar reasoning yields  $n^{1/2}(\hat{\varphi}_{hk_n} - \varphi_{ho}) \rightsquigarrow N(0, \text{avar}(\hat{\varphi}_{hk_n}))$ .  $\triangle$

**Proof of Lemma 5.** Define  $\text{avar}_n(\theta) := n^{-1} \sum_{i=1}^n \ln g(z_i, \theta)^2 - [n^{-1} \sum_{i=1}^n \ln g(z_i, \theta)]^2$  and  $\text{avar}(\theta) := E[\ln g(z_i, \theta)^2] - E[\ln g(z_i, \theta)]^2$ . By the triangle inequality with probability approaching one

$$|\text{avar}_n(\hat{\varphi}_{k_n}) - \text{avar}(\hat{\varphi}_{k_n})| \leq |\text{avar}_n(\hat{\theta}_{k_n}) - \text{avar}(\hat{\theta}_{k_n})| + |\text{avar}(\hat{\theta}_{k_n}) - \text{avar}(\theta_\star)|.$$

Consider the first term in the right hand side of this inequality. From (B.ii), (B.iii) and the i.i.d. assumption,  $\sup_\theta |\text{avar}_n(\theta) - \text{avar}(\theta)| = o_{P_o}(1)$ . Hence,  $|\text{avar}_n(\hat{\theta}_{k_n}) - \text{avar}(\hat{\theta}_{k_n})| = o_{P_o}(1)$ . Consider now the second term. From (B.iii),  $\theta \mapsto \text{avar}(\theta)$  is continuous. Since  $\hat{\theta}_{k_n} \xrightarrow{P_o} \theta_\star$ , by the Continuous Mapping Theorem,  $|\text{avar}(\hat{\theta}_{k_n}) - \text{avar}(\theta_\star)| = o_{P_o}(1)$ . It follows then that  $\text{avar}_n(\hat{\varphi}_{k_n}) \xrightarrow{P_o} \text{avar}(\hat{\varphi}_{k_n})$ . A similar result follows for  $\text{avar}_n(\hat{\varphi}_{hk_n})$  and  $\text{acov}_n(\hat{\varphi}_{g, k_n}, \hat{\varphi}_{hk_n})$ . Then, by the Continuous Mapping Theorem,  $\hat{\omega}_n \xrightarrow{P_o} \omega_o$ .  $\triangle$

**Proof of Lemma 6.** It suffices to verify that  $\omega_o = 0$  iff  $g(z_i, \theta_\star) = h(z_i, \gamma_\star)$  for any  $\theta_\star \in \arg \min_{\theta \in \Theta} G(\theta, P_o)$ ,  $\gamma_\star \in \arg \min_{\gamma \in \Gamma} H(\gamma, P_o)$ . Fix  $\theta_\star$  and  $\gamma_\star$ . From the definition of  $\omega_o$ ,  $\omega_o = 0$  iff there exists a constant  $\epsilon$  such that  $g(z_i, \theta_\star) = \epsilon h(z_i, \gamma_\star)$  a.e.  $z_i$ . Since  $z \mapsto g(z, \theta_\star)$  and  $z \mapsto h(z, \gamma_\star)$  are density functions, they integrate to one. It follows then, by integrating both sides of  $g(z, \theta_\star) = \epsilon h(z, \gamma_\star)$  with respect to  $z$ , that  $\epsilon = 1$ .  $\triangle$