

Generalised Bayes with intractable distances

Outline

1. Generalised Bayesian beliefs
2. Loss-based generalised Bayes
3. Computational challenges & Examples
4. Novel asymptotic theory

Today's take-away in a nutshell

1. Generalised Bayes with distances (between model & data-generating process) is attractive
2. Challenge: often, such distances are intractable & require simulation (KDE, MMD, DPD)
 1. Turner and Sederberg (2014) - KDE
 2. Cherief-Abdellatif and Alquier (2020) - MMD
 3. Jewson et al. (2018) - Beta Divergence
 4. Pacchiardi et al. (2022) - Kernel Scores
 5. Legramanti et al. (2023) - IPMs
3. Show that this requires analysing novel posterior object
4. Need to make sure we are (asymptotically) targeting the correct object

Standard Bayesian beliefs

Standard Bayes posterior

Empirical measure, $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

$$\pi(\theta | \widehat{\text{KL}}_n) = \frac{\prod_{i=1}^n p_{\theta}(x_i) \pi(\theta)}{\int \prod_{i=1}^n p_{\theta}(x_i) \pi(\theta) d\theta} = \operatorname{argmin}_{q \in \mathcal{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q} \left[n \cdot \widehat{\text{KL}}_n(P_n, p_{\theta}) \right] + \text{KL}(q, \pi) \right\}$$

$$\widehat{\text{KL}}_n(P_n, p_{\theta}) = -\frac{1}{n} \sum_{i=1}^n \log p_{\theta}(x_i) + C$$

Loss-based generalised Bayesian beliefs

Loss-based generalised posterior

$$\pi(\theta | \widehat{D}_n) = \frac{\exp\{-n \cdot \mathbf{D}(P_n, P_\theta)\} \pi(\theta)}{\int \exp\{-n \cdot \mathbf{D}(P_n, P_\theta)\} \pi(\theta) d\theta} = \operatorname{argmin}_{q \in \mathcal{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim q} \left[n \cdot \mathbf{D}(P_n, P_\theta) \right] + \operatorname{KL}(q, \pi) \right\}$$

Estimator for some loss function \mathbf{D}

Why distance-based generalised Bayesian beliefs? (Robustness)

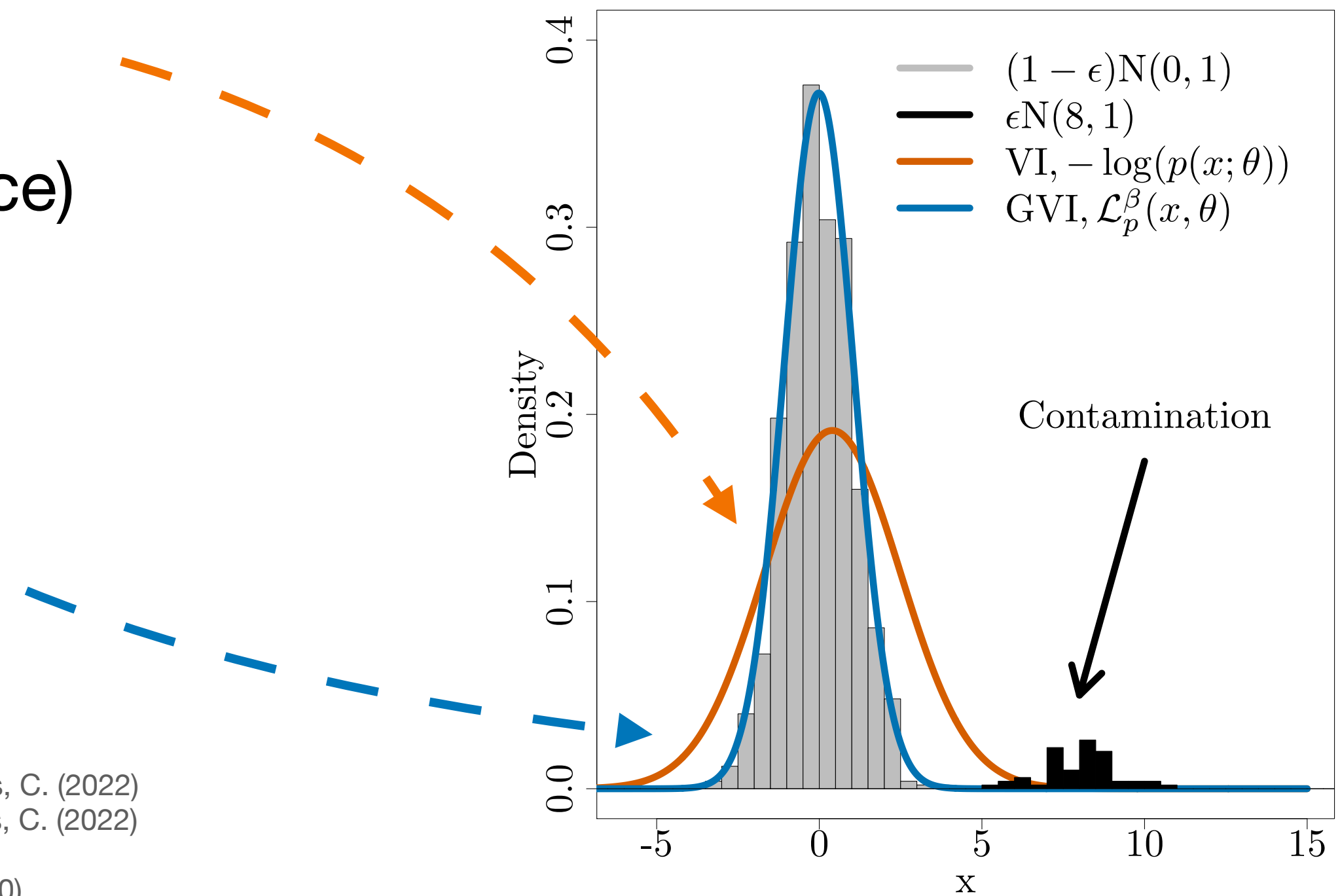
$$\pi(\theta | D_n) = \frac{\exp\{-n \cdot D(P_n, p_\theta)\} \pi(\theta)}{\int \exp\{-n \cdot D(P_n, p_\theta)\} \pi(\theta) d\theta}$$

1. **Standard Bayes**: brittle/unreliable/non-robust under model misspecification (because KL is not a robust measure of distance)
2. But we want to keep Bayesian advantages (UQ, expert input via priors)

Solution: belief distributions on the parameter through **robust distances**

Generalised Bayesian inference for discrete intractable likelihood, Matsubara, T., Knoblauch, J., Briol, F.X., and Oates, C. (2022)
Robust generalised Bayesian inference for intractable likelihoods, Matsubara, T., Knoblauch, J., Briol, F.X., and Oates, C. (2022)
Generalised Bayesian likelihood-free inference using scoring rule estimators, Pacchiardi, L. & Dutta, R. (2021)
MMD-Bayes: Robust Bayesian estimation via maximum mean discrepancy, Cherrief-Abdellatif, B.E. & Alquier, P. (2020)
Principles of Bayesian Inference using General Divergence Criteria, Jewson, J., Smith, J., & Holmes, C., (2016)
Robust Bayes estimation using the density power divergence, Ghosh, A. & Basu, A. (2016)
Bayesian model robustness via disparities, Hooker, G. & Vidyashankar, A. N. (2014)

...



When are distance-based generalised Bayesian beliefs difficult?

$$\pi(\theta | \widehat{\text{KL}}_n) = \frac{\prod_{i=1}^n p_{\theta}(x_i) \pi(\theta)}{\int \prod_{i=1}^n p_{\theta}(x_i) \pi(\theta) d\theta}$$



Need to evaluate $p_{\theta}(x)$ for $x \in \{x_i : 1 \leq i \leq n\}$

Why loss-based generalised Bayesian beliefs?

$$\pi(\theta | \widehat{\text{KL}}_n) = \frac{\prod_{i=1}^n p_{\theta}(x_i) \pi(\theta)}{\int \prod_{i=1}^n p_{\theta}(x_i) \pi(\theta) d\theta}$$



Need to evaluate $p_{\theta}(x)$ for $x \in \{x_i : 1 \leq i \leq n\}$

$$\pi(\theta | D_n) = \frac{\exp\{-n \cdot D(P_n, p_{\theta})\} \pi(\theta)}{\int \exp\{-n \cdot D(P_n, p_{\theta})\} \pi(\theta) d\theta}$$



Need to evaluate/approximate $D(P_n, p_{\theta})$



E.g., MMD^2 for simulators,
 KSD^2 for unnormalised models

Computational challenges (MMD)

$$\pi(\theta | D_n) \propto \exp\{-n \cdot D(P_n, P_\theta)\} \pi(\theta)$$



Actual target: $\text{MMD}^2(P_n, P_\theta) = \mathbb{E}_{X \sim P_n, X' \sim P_n} [k(X, X')] - 2\mathbb{E}_{X \sim P_n, Y \sim P_\theta} [k(X, Y)] + \mathbb{E}_{Y \sim P_\theta, Y' \sim P_\theta} [k(Y, Y')]$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(x_i, x_j) - \frac{2}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim P_\theta} [k(x_i, Y)] + \mathbb{E}_{Y \sim P_\theta, Y' \sim P_\theta} [k(Y, Y')]$$

Intractable expectations

Computation: challenges (MMD/kernel score example)

$$\pi(\theta | D_n) \propto \exp\{-n \cdot D(P_n, P_\theta)\} \pi(\theta)$$



Actual target: $\text{MMD}^2(P_n, P_\theta) = \mathbb{E}_{X \sim P_n, X' \sim P_n} [k(X, X')] - 2\mathbb{E}_{X \sim P_n, Y \sim P_\theta} [k(X, Y)] + \mathbb{E}_{Y \sim P_\theta, Y' \sim P_\theta} [k(Y, Y')]$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(x_i, x_j) - \frac{2}{n} \sum_{i=1}^n \mathbb{E}_{Y \sim P_\theta} [k(x_i, Y)] + \mathbb{E}_{Y \sim P_\theta, Y' \sim P_\theta} [k(Y, Y')]$$

Intractable expectations

What we compute: $\text{MMD}^2(P_n, P_{m,\theta}) = \mathbb{E}_{X \sim P_n, X' \sim P_n} [k(X, X')] - 2\mathbb{E}_{X \sim P_n, Y \sim P_{m,\theta}} [k(X, Y)] + \mathbb{E}_{Y \sim P_{m,\theta}, Y' \sim P_{m,\theta}} [k(Y, Y')]$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(x_i, x_j) - 2\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m k(x_i, y_j) + \frac{1}{m^2} \sum_{j=1}^m \sum_{j=1}^m k(y_i, y_j)$$

Computation: challenges

Some other examples for D with this problem:

$$W_p^p(P_n, P_\theta) = \left[\int_{x \in \mathcal{X}} \left| F_{P_\theta}^{-1}(x) - F_n^{-1}(x) \right| dx \right]^{1/p} \quad \text{if } \mathcal{X} \subset \mathbb{R}$$

F_n^{-1} = empirical quantile function

$F_{P_\theta}^{-1}$ = model quantile function

Computation: challenges

Some other examples for D with this problem:

$$W_p^p(P_n, P_\theta) = \left[\int_{x \in \mathcal{X}} \left| F_{P_\theta}^{-1}(x) - F_n^{-1}(x) \right| dx \right]^{1/p} \quad \text{if } \mathcal{X} \subset \mathbb{R}$$

F_n^{-1} = empirical quantile function
 $F_{P_\theta}^{-1}$ = model quantile function

$$KS(P_n, P_\theta) = \sup_{x \in \mathcal{X}} \left| F_{P_\theta}(x) - F_n(x) \right|$$

F_n = empirical cdf
 F_{P_θ} = model cdf

$$CvM(P_n, P_\theta) = \int_{x \in \mathcal{X}} \left| F_{P_\theta}(x) - F_n(x) \right|^2 dF_n(x)$$

Computation: challenges

Some other examples for D with this problem:

$$W_p^p(P_n, P_\theta) = \left[\int_{x \in \mathcal{X}} \left| F_{P_\theta}^{-1}(x) - F_n^{-1}(x) \right| dx \right]^{1/p} \quad \text{if } \mathcal{X} \subset \mathbb{R}$$

$$KS(P_n, P_\theta) = \sup_{x \in \mathcal{X}} \left| F_{P_\theta}(x) - F_n(x) \right|$$

$$CvM(P_n, P_\theta) = \int_{x \in \mathcal{X}} \left| F_{P_\theta}(x) - F_n(x) \right|^2 dF_n(x)$$

F_n^{-1} = empirical quantile function
 $F_{P_\theta}^{-1}$ = model quantile function

F_n = empirical cdf
 F_{P_θ} = model cdf

Also: general IPMs, $\beta/\alpha/\gamma$ -divergences, ..., any intractable models as well...

Computation: challenges

What we actually want:

$$\pi(\theta | D_n) \propto \exp\{-n \cdot D(P_n, p_\theta)\} \pi(\theta)$$

What we can compute:

$$\pi(\theta | \widehat{D}_{m,n}, y_{1:m}) \propto \exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\} \pi(\theta)$$

Estimator $\widehat{D} \approx D$

Sample-based approximation of p_θ ,

Two pieces: n-sample data points, m-simulated data points

$$P_{m,\theta} = \frac{1}{m} \sum_{j=1}^m \delta_{y_j} \quad y_{1:m} \stackrel{i.i.d.}{\sim} p_\theta$$

Computation: challenges

What we actually want:

$$\pi(\theta | D_n) \propto \exp\{-n \cdot D(P_n, p_\theta)\} \pi(\theta)$$

What we (can) compute:

$$\pi(\theta | \widehat{D}_{m,n}, y_{1:m}) \propto \exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\} \pi(\theta)$$

Estimator $\widehat{D} \approx D$

Sample-based approximation of p_θ ,

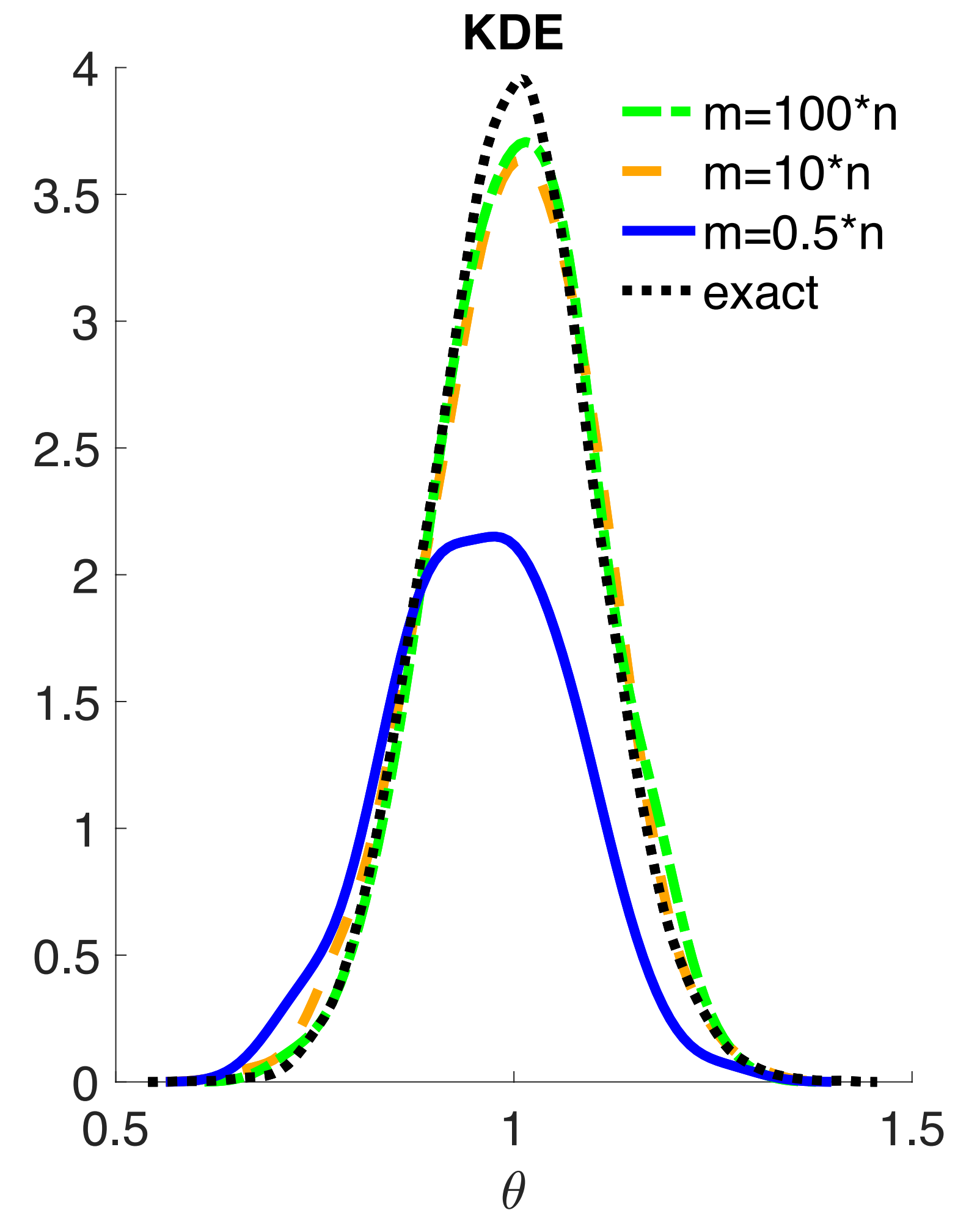
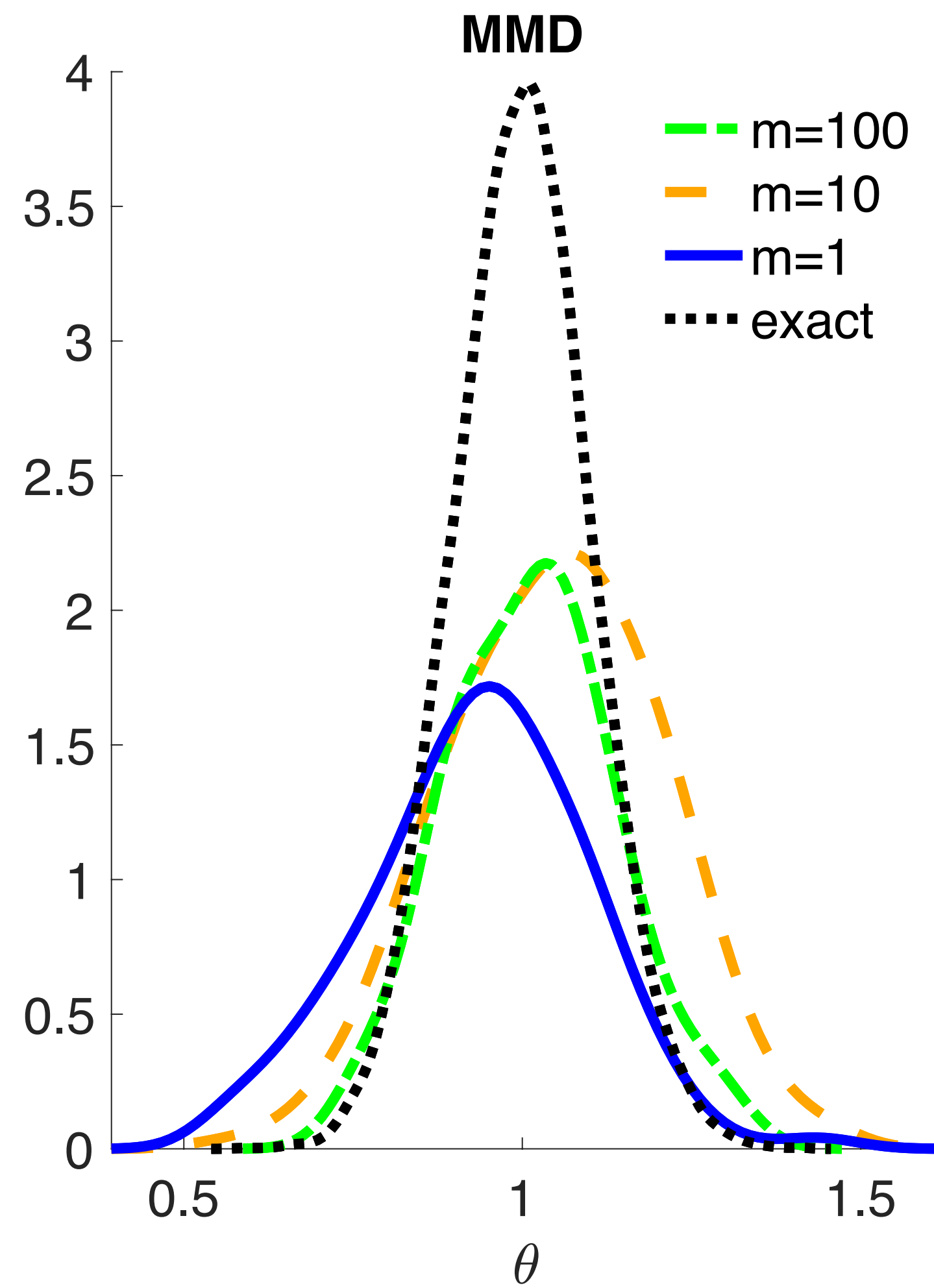
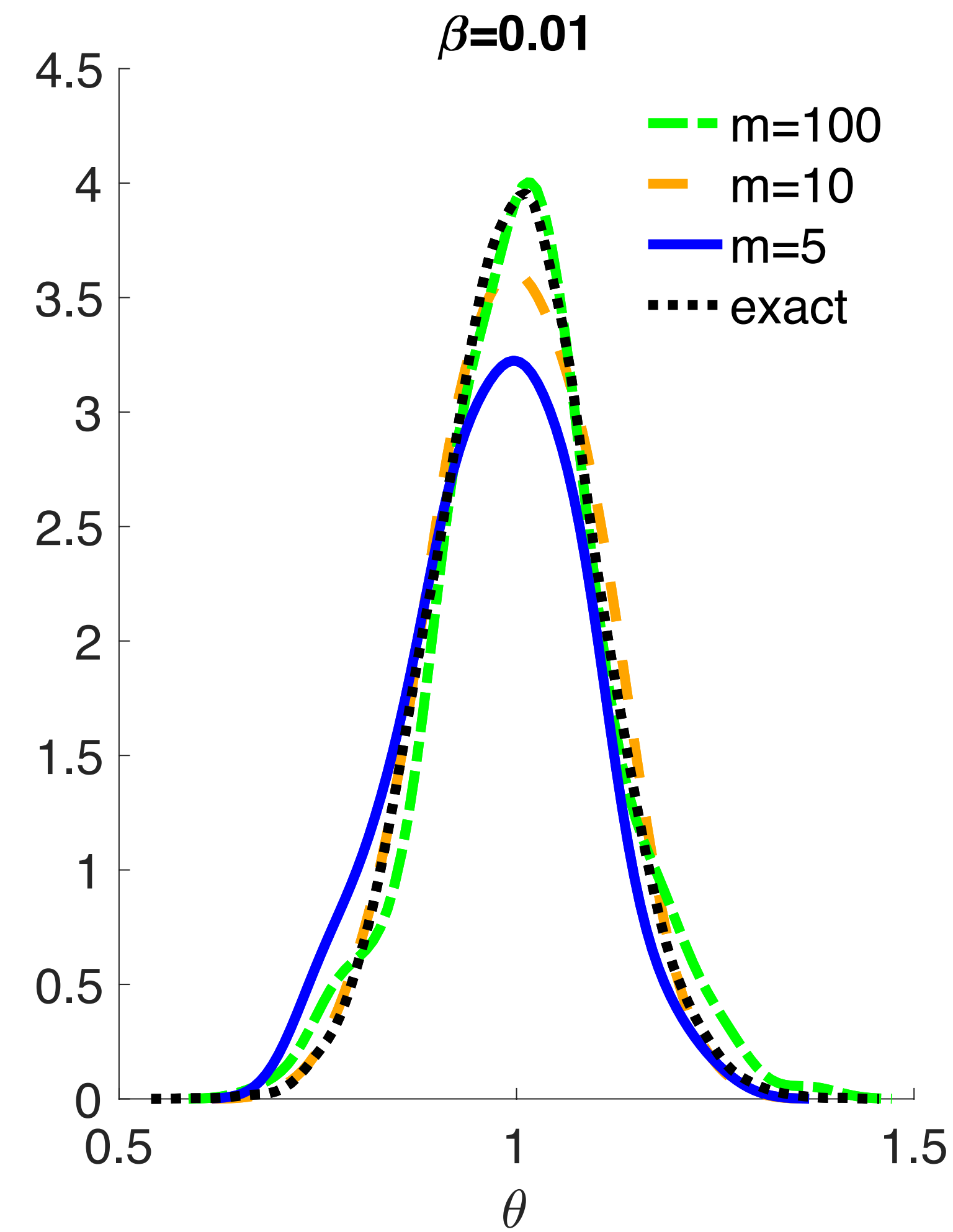
$$P_{m,\theta} = \frac{1}{m} \sum_{j=1}^m \delta_{y_j} \quad y_{1:m} \stackrel{i.i.d.}{\sim} p_\theta$$

Research questions we tackle:

- (1) How different is $\pi(\theta | D_{m,n}, y_{1:m})$ from $\pi(\theta | \widehat{D}_{m,n}, y_{1:m})$?
- (2) How does $\widehat{D}_{m,n}$ dictate behaviour of and $\pi(\theta | \widehat{D}_{m,n}, y_{1:m})$?

Preliminary answers to (1) and (2)

Data generated iid from $N(\theta, 1)$, $\theta = 1$



Computation: accounting for model samples

Question (1): How different is $\pi(\theta | D_n)$ from $\pi(\theta | \widehat{D}_{m,n}, y_{1:m})$?

$$\pi(\theta | \widehat{D}_{m,n}, y_{1:m}) \propto \exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta}, y_{1:m})\} \pi(\theta)$$



Sample-based approximation of p_θ ,

$$P_{m,\theta} = \frac{1}{m} \sum_{j=1}^m \delta_{y_j} \quad y_{1:m} \stackrel{i.i.d.}{\sim} p_\theta$$



But hold on! $y_{1:m}$ itself is random! (not fixed/conditioned on!)

Computation: accounting for model samples

Question (1): How different is $\pi(\cdot | D_n)$ from $\pi(\cdot | \widehat{D}_{m,n})$?

$$\pi(\theta | \widehat{D}_{m,n}, y_{1:m}) \propto \exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\} \pi(\theta)$$

$\pi(\theta | \widehat{D}_{m,n}, y_{1:m})$ constructed from draws $y_{1:m} \stackrel{i.i.d.}{\sim} P_\theta$
is itself random, and an (unbiased) estimate of $\bar{\pi}(\theta | \widehat{D}_{m,n})$

$$\bar{\pi}(\theta | \widehat{D}_{m,n}) \propto \pi(\theta) \cdot \mathbb{E}_{y_{1:m} \sim P_\theta} \left[\exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\} \right] \quad (\neq \pi(\theta | D_n))$$

Computation: accounting for model samples

Question (1): How different is $\pi(\cdot | D_n)$ from $\pi(\cdot | \widehat{D}_{m,n})$?

$$\pi(\theta | \widehat{D}_{m,n}) \propto \exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\} \pi(\theta)$$

$\pi(\theta | \widehat{D}_{m,n})$ constructed from draws $y_{1:m} \stackrel{i.i.d.}{\sim} p_\theta$
is itself random, and an (unbiased) estimate of $\bar{\pi}(\theta | \widehat{D}_{m,n})$

$$\bar{\pi}(\theta | \widehat{D}_{m,n}) \propto \pi(\theta) \cdot \underbrace{\mathbb{E}_{y_{1:m} \sim P_\theta} \left[\exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\} \right]}_{= \text{like an intractable likelihood!}} \quad (\neq \pi(\theta | D_n))$$

Computation: accounting for model samples

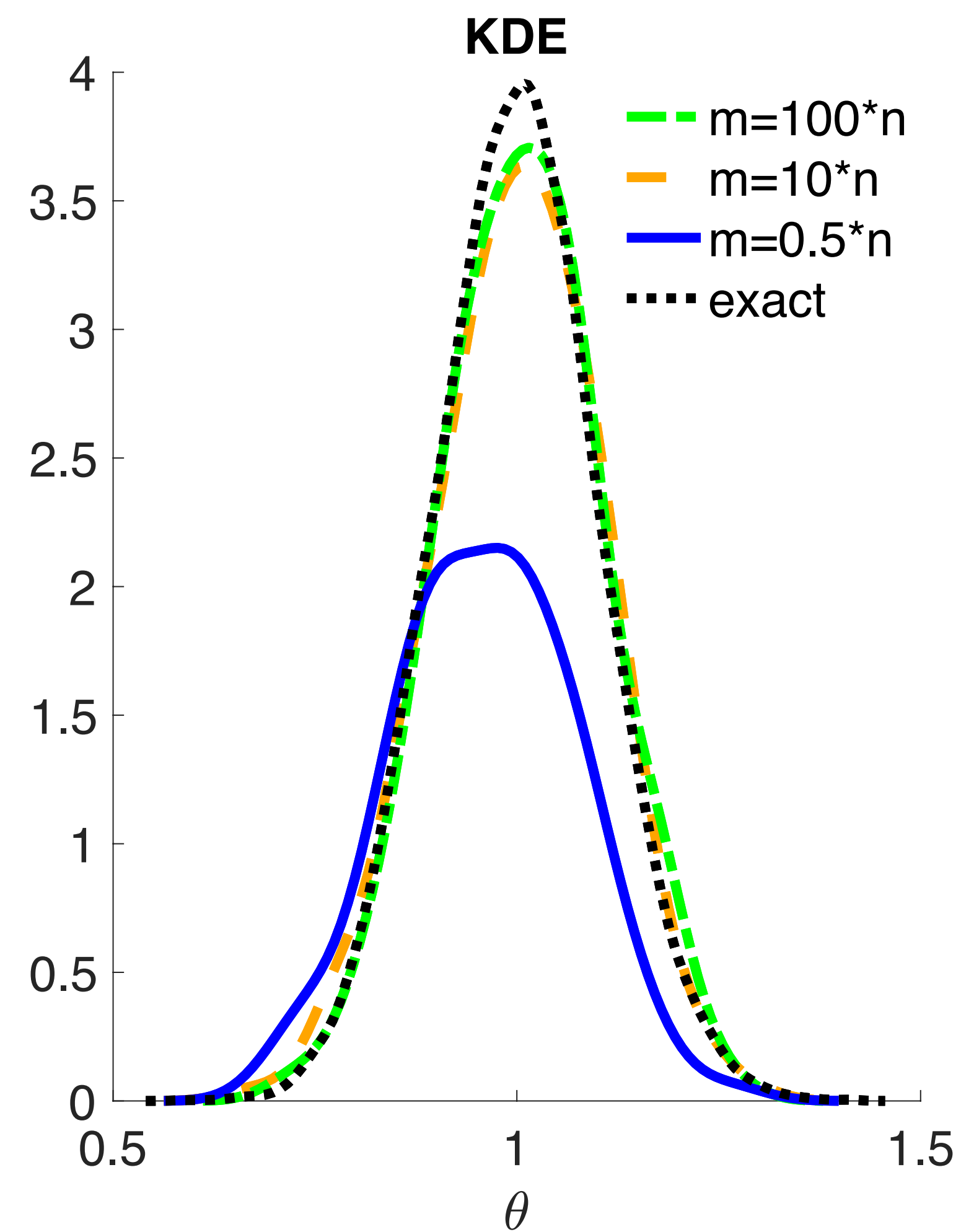
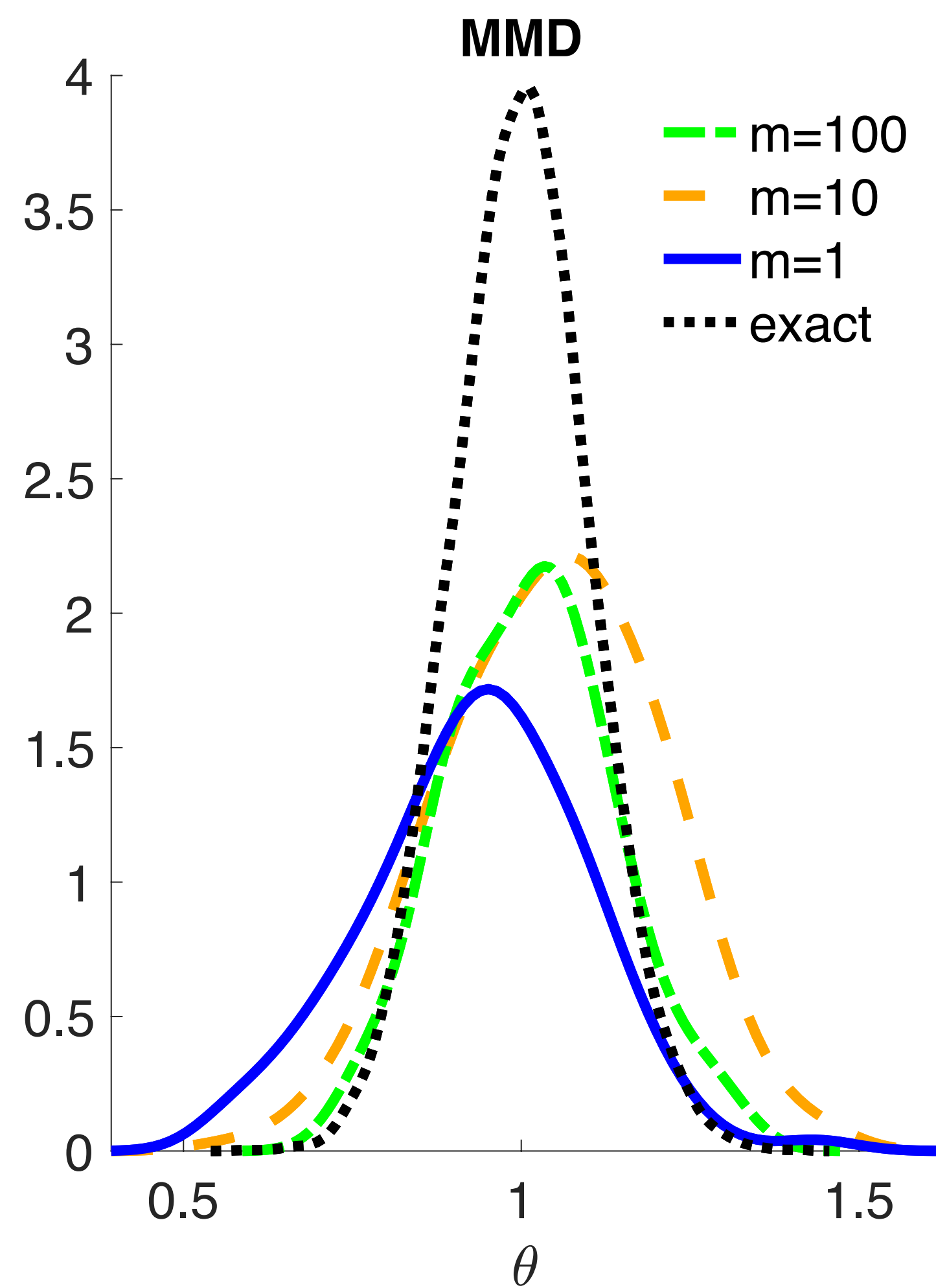
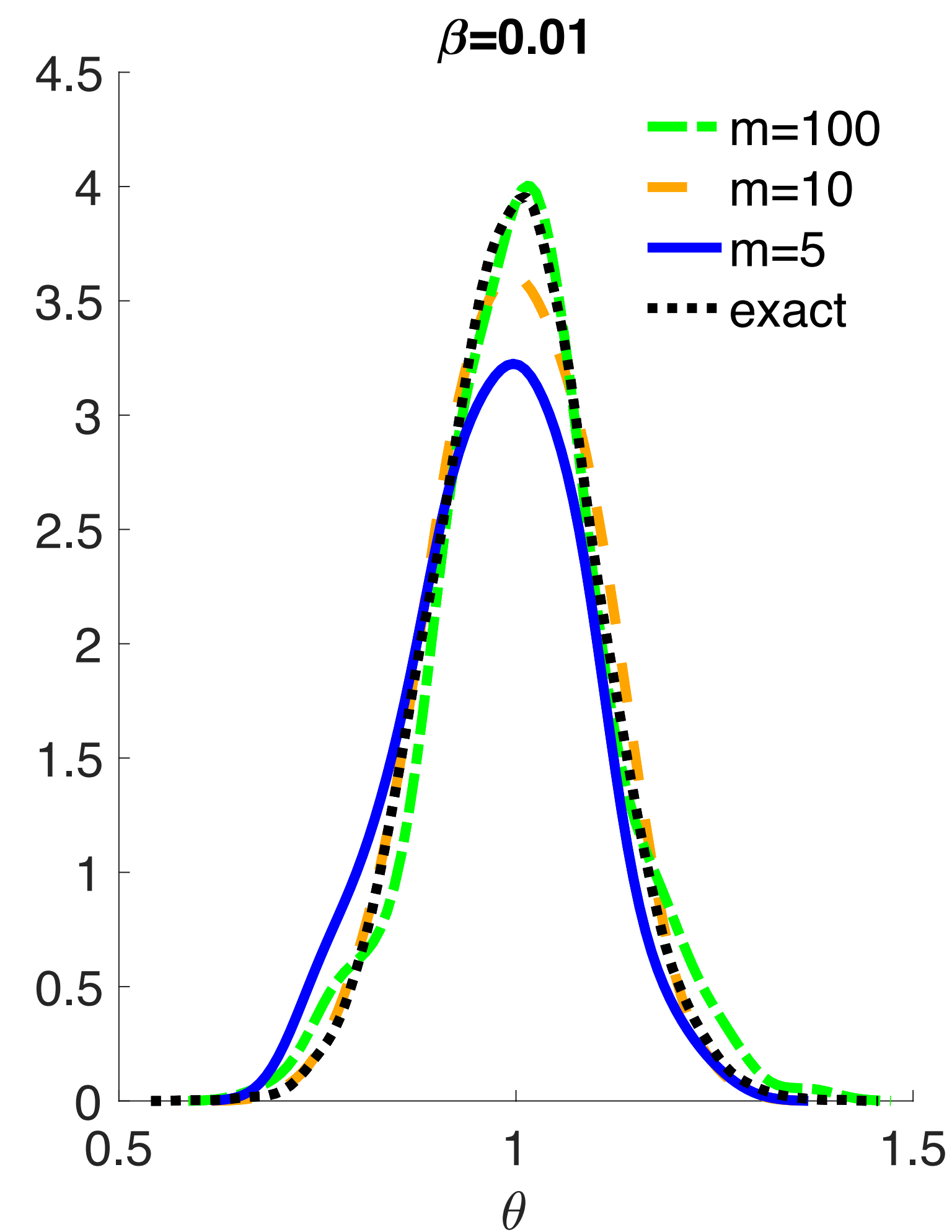
Question (1): How different is $\pi(\cdot | D_n)$ from $\pi(\cdot | \widehat{D}_{m,n})$?

$$\begin{aligned}\bar{\pi}(\theta | \widehat{D}_{m,n}) &\propto \pi(\theta) \cdot \mathbb{E}_{y_{1:m} \sim p_\theta} \left[\exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\} \right] \\ &= \pi(\theta) \cdot \int_{\mathcal{X}^m} \left(\prod_{j=1}^m p_\theta(y_j) \right) \cdot \underbrace{\exp\{-n \cdot \widehat{D}(P_n, P_{m,\theta})\}}_{= \text{like a kernel in ABC/BSL}} dy_{1:m} \\ &\neq \pi(\theta | D_n)\end{aligned}$$

However: clearly $\bar{\pi}_n^D \neq \pi_n^{D,\text{ideal}}$; so can we quantify/bound the difference?

Preliminary answers to (1) and (2)

Data generated iid from $N(\theta, 1)$, $\theta = 1$



Computation: accounting for model samples

Theorem 1:

Under mild regularity conditions on the moments of $\widehat{D}(P_n, P_{m,\theta})$ (in $y_{1:m}$) and $\pi_n^{D,\text{ideal}}$, then for all $\gamma \in [0,2]$ and as $n, m \rightarrow \infty$ with $m \gg n$,

$$\int_{\Theta} \|\theta\|^\gamma \left| \pi(\theta | D_n) - \bar{\pi}(\theta | \widehat{D}_{m,n}) \right| d\theta = O_p(\max\{m^{-\kappa_1}, m^{-\kappa_2}\}),$$

which immediately implies that

$$\text{TVD}(\pi_n^{D,\text{ideal}}, \bar{\pi}_n^D) = O_p(\max\{m^{-\kappa_1}, m^{-\kappa_2}\}), \quad \left\| \mathbb{E}_{\theta \sim \pi_n^{D,\text{ideal}}}[\theta] - \mathbb{E}_{\theta \sim \bar{\pi}_n^D}[\theta] \right\| = O_p(\max\{m^{-\kappa_1}, m^{-\kappa_2}\}).$$

Computation: accounting for model samples

Theorem 1:

Under mild regularity conditions on the moments of $\widehat{D}(P_n, P_{m,\theta})$ (in $y_{1:m}$) and $\pi_n^{D,\text{ideal}}$, then for all $\gamma \in [0,2]$ and as $n, m \rightarrow \infty$ with $m \gg n$,

$$\int_{\Theta} \|\theta\|^\gamma \left| \pi(\theta | D_n) - \bar{\pi}(\theta | \widehat{D}_{m,n}) \right| d\theta = O_p(\max\{m^{-\kappa_1}, m^{-\kappa_2}\}),$$

Bias Term

Variance Term

which immediately implies that

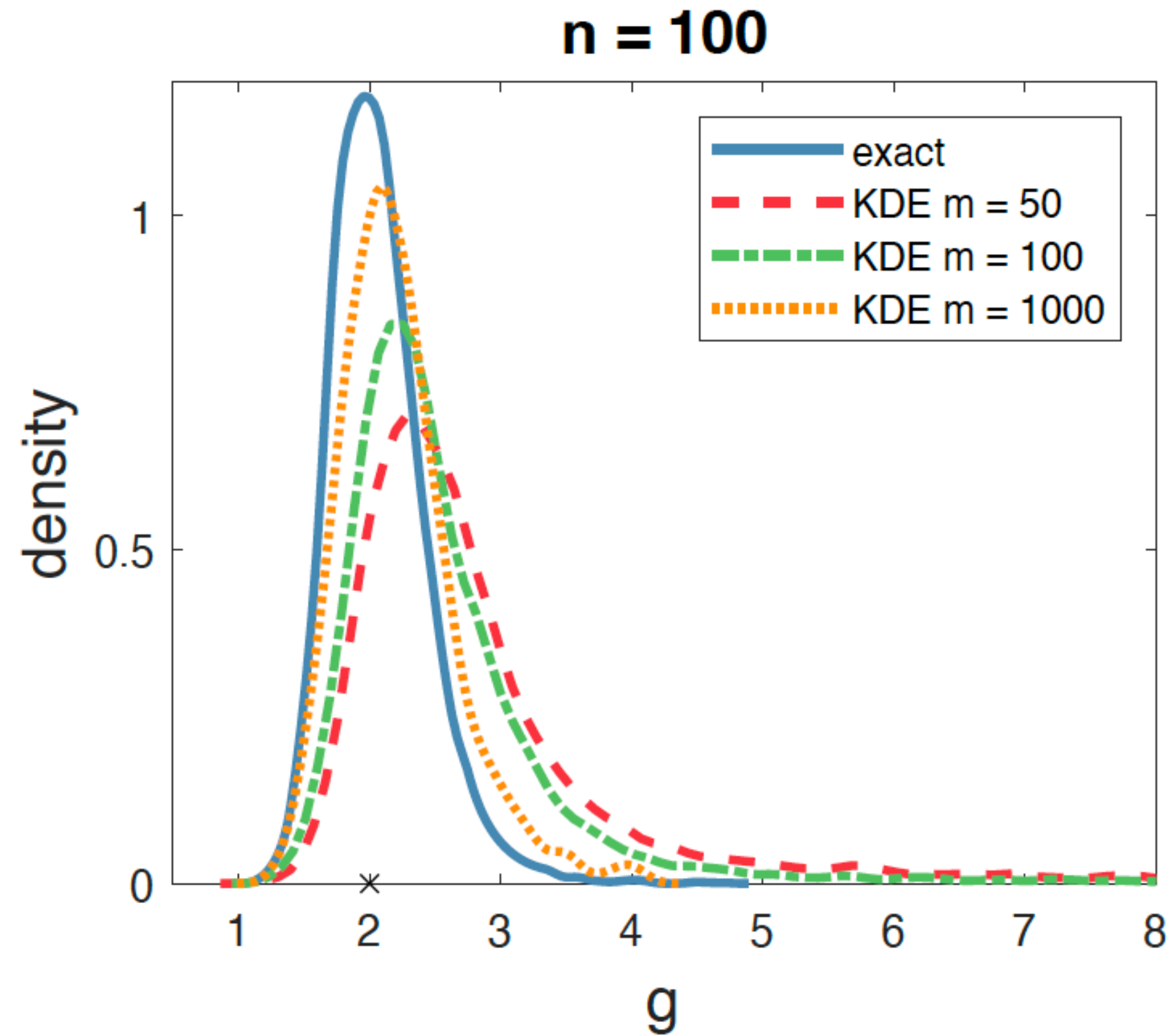
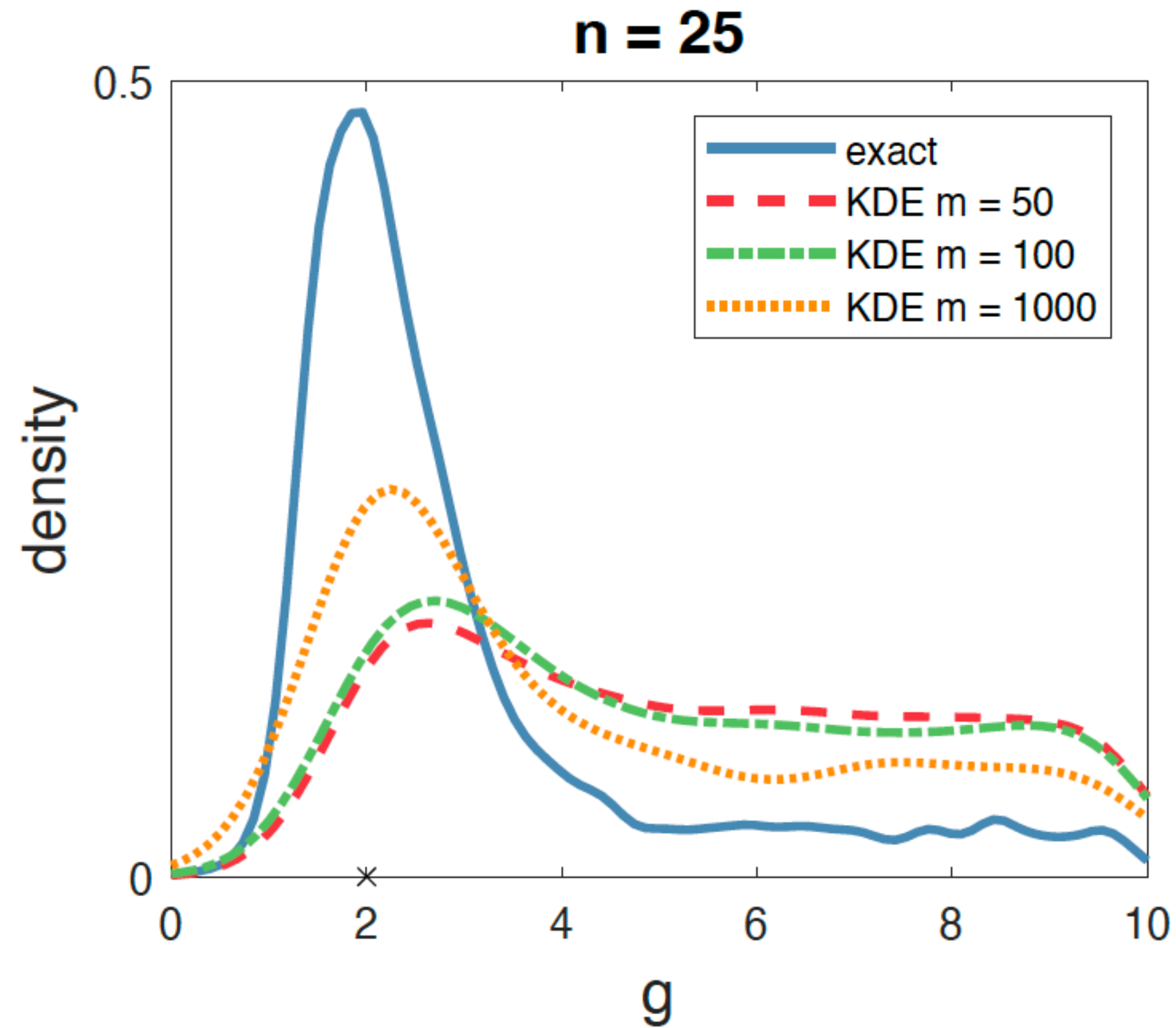
$$\text{TVD}(\pi_n^{D,\text{ideal}}, \bar{\pi}_n^D) = O_p(\max\{m^{-\kappa_1}, m^{-\kappa_2}\}), \quad \left\| \mathbb{E}_{\theta \sim \pi_n^{D,\text{ideal}}}[\theta] - \mathbb{E}_{\theta \sim \bar{\pi}_n^D}[\theta] \right\| = O_p(\max\{m^{-\kappa_1}, m^{-\kappa_2}\}).$$

Tells us that implicitly targeting $\bar{\pi}_n^D$ instead of $\pi_n^{D,\text{ideal}}$ is (generally) fine
(bias vanishes at linear rate for larger model-samples)

Asymptotic & choice of distance

Question (2): How does D dictate behaviour of $\bar{\pi}_n^D$ and $\pi_n^{D,ideal}$?

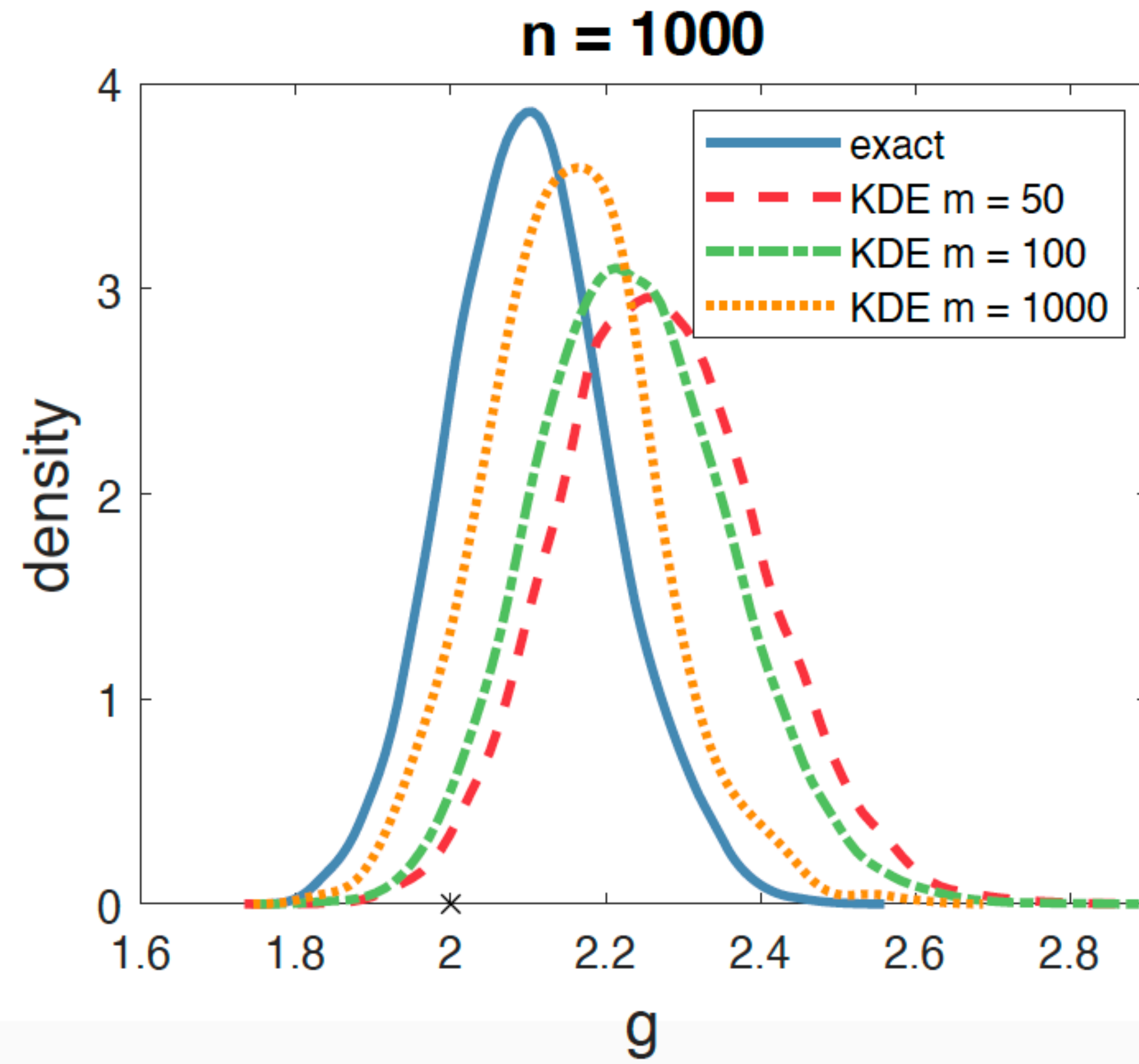
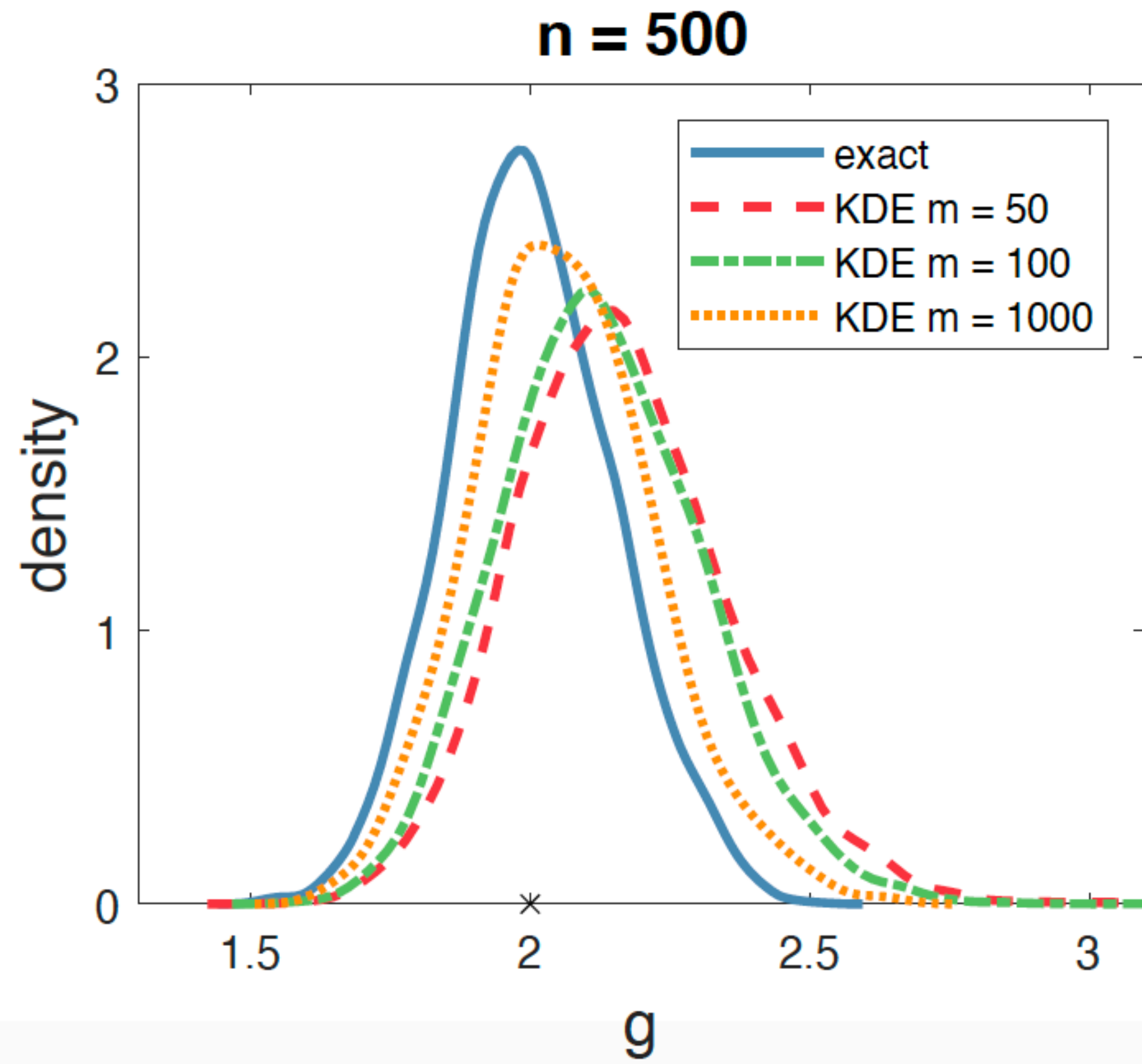
Data generated iid from G-and-K dist



Asymptotic & choice of distance

Question (2): How does D dictate behaviour of $\bar{\pi}_n^D$ and $\pi_n^{D,ideal}$?

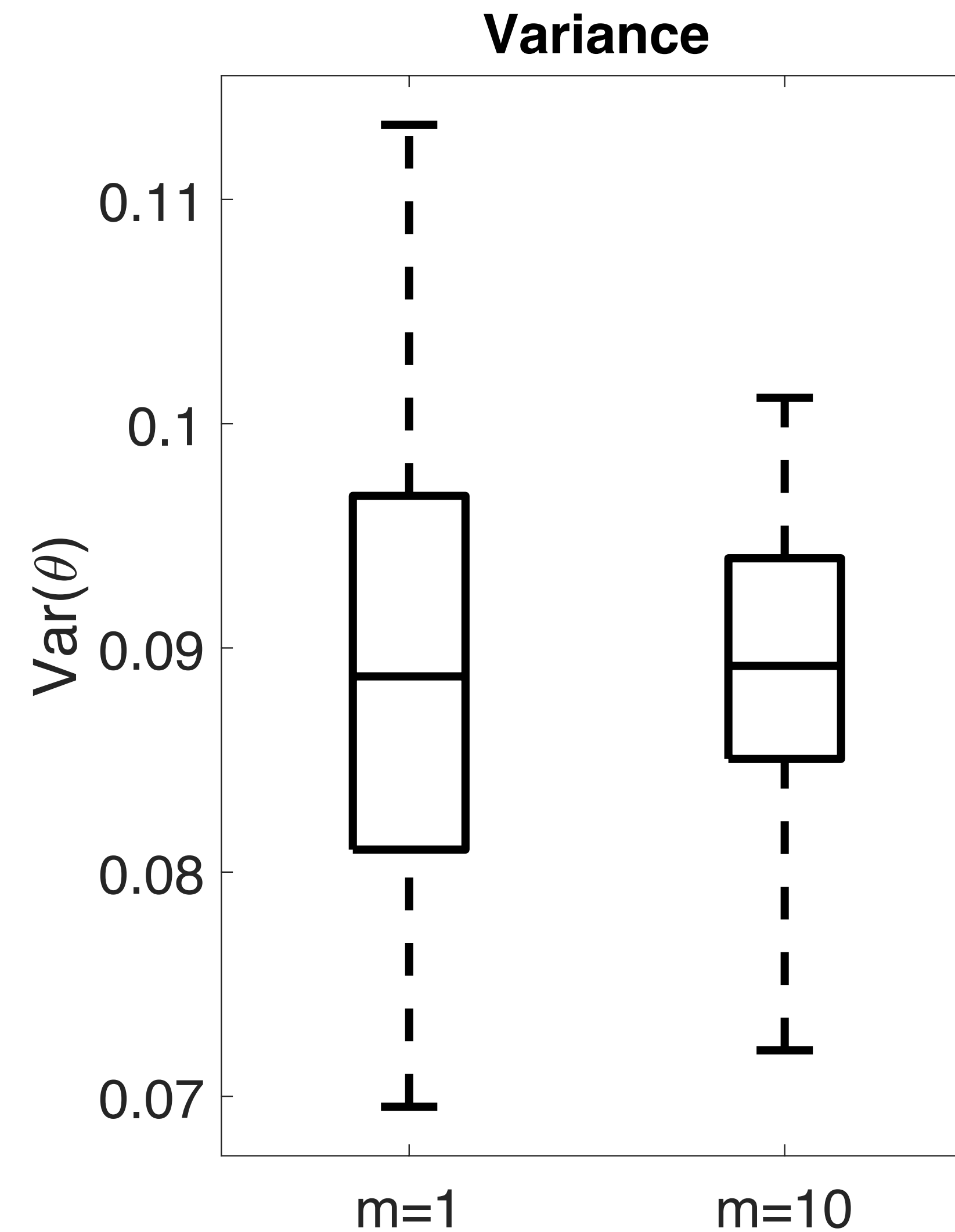
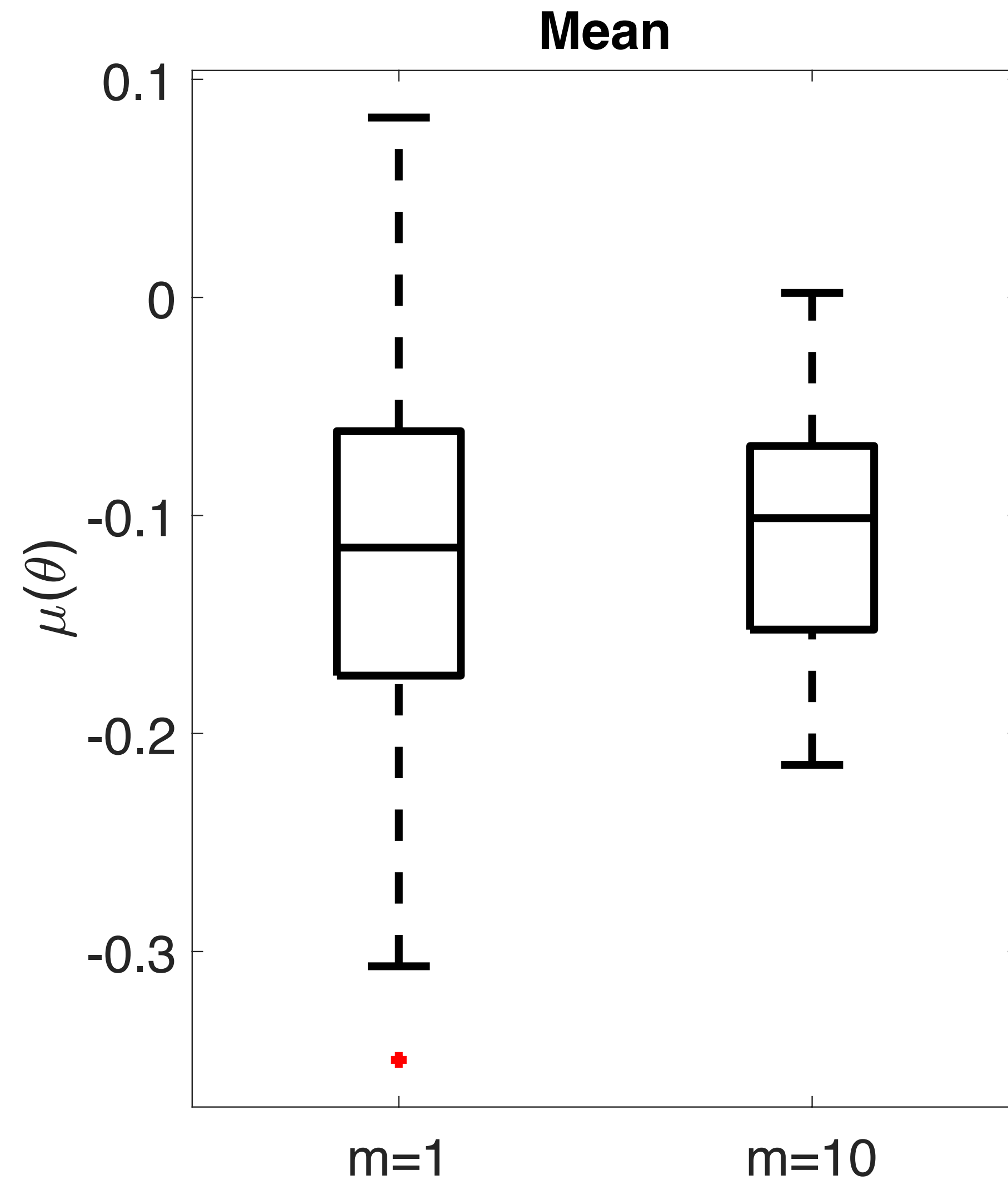
Data generated iid from G-and-K dist



Asymptotics & choice of distance: MMD

Question (2): How does D dictate behaviour of $\bar{\pi}_n^D$?

Data generated iid from t-copula: $\theta = -0.1, \nu = 10$ – fixed



Summary & implications

Research questions we tackled:

(1) How different is $\bar{\pi}_n^D$ from $\pi_n^{D,\text{ideal}}$?

⇒ not very! (at least for n large)

⇒ Standard sampling algorithms will approximate a sensible target!

Summary & implications

Research questions we tackled:

(1) How different is $\bar{\pi}_n^D$ from $\pi_n^{D,\text{ideal}}$?

⇒ not very! (see Theorem 1)

⇒ Standard sampling algorithms will approximate a sensible target!

(2) What behaviour should we expect from $\bar{\pi}_n^D$, ?

⇒ Concentration depends on choice of \mathbf{D} and in particular $\mathbb{E}_{y_{1:m}}[\widehat{\mathbf{D}}_n], (\kappa_1)$ and $\mathbb{V}_{y_{1:m}}[\widehat{\mathbf{D}}_n], (\kappa_2)$

⇒ $\bar{\pi}(\theta \in \cdot \mid \widehat{\mathbf{D}}_{m,n})$ concentration depends on bias & variance!

⇒ choice of $\widehat{\mathbf{D}}_{m,n}$ impacts concentration! → Concentration occurs at slower of

\sqrt{n} – normal parametric rate

$m^{-\kappa_1}$ – Bias of estimator for $\widehat{\mathbf{D}}_{m,n}$

$m^{-\kappa_2}$ – Variance of estimator for $\widehat{\mathbf{D}}_{m,n}$

Summary & implications

Research questions we tackled:

(3) Can this inform/justify methodology?

⇒ Pick \mathbf{D} and $\hat{\mathbf{D}}$ with small bias and low variance: $\text{MMD}^2 \checkmark$, $\text{KDE}^2 \times$

⇒ Given \mathbf{D} and $\hat{\mathbf{D}}_{m,n}$ when no bias $m \asymp \sqrt{n}$, else must choose $m \gg \sqrt{n}$

References

- Generalized Bayesian Likelihood-Free Inference Using Scoring Rules Estimators*, Pacchiardi, L., & Dutta, R. arXiv:2104.03889 (2022)
- A general framework for updating belief distributions*, Bissiri, P.G., Holmes, C., & Walker, S.G. (2016)
- Principles of Bayesian Inference using General Divergence Criteria*, Jewson, J., Smith, J., & Holmes, C., (2016)
- The pseudo-marginal approach for efficient Monte Carlo computations*, Andrieu, C. & Roberts, G. O. (2009)
- Robust Generalised Bayesian Inference for Intractable Likelihoods*, Matsubara, T., Knoblauch, J., Briol, F.-X., & Oates, C. JRSSB (2022)
- Generalised Bayesian Inference for Discrete Intractable Likelihood*, Matsubara, T., Knoblauch, J., Briol, F.-X., & Oates, C. arXiv:2206.08420 (2022)
- An Optimization-centric View on Bayes' Rule: Reviewing and Generalizing Variational Inference*, Knoblauch, J., Jewson, J., & Damoulas, T., JMLR (2022).
- Robust Bayesian Inference for Simulator-based Models via the MMD Posterior Bootstrap*, Dellaporta, C., Knoblauch, J., Damoulas, T., & Briol, F.-X. AISTATS (2022)
- Generalised Posteriors in Approximate Bayesian Computation*, Schmon, S.M., Cannon, P., & Knoblauch, J., AABI (2020)