

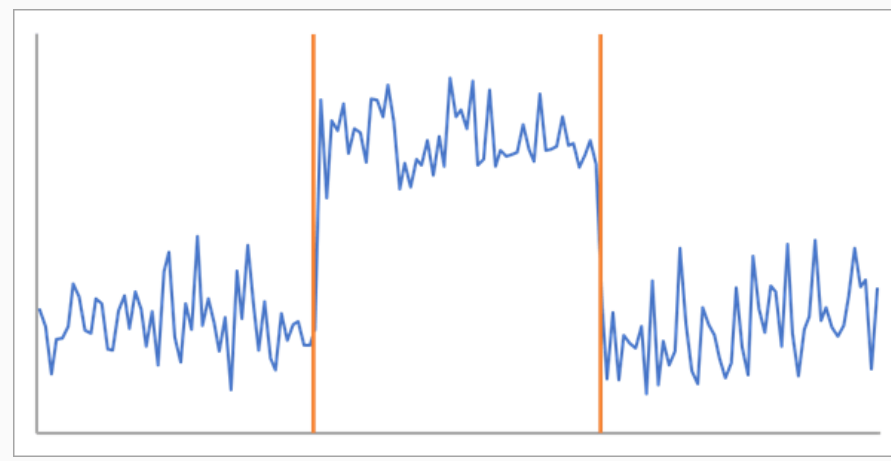
Motivations

- Observe i.i.d sequence $(X_t)_{t \in \mathbb{N}}$ drawn from parametric distributions p_θ .
- Estimate θ and confidence set $\hat{\Theta}_t^\delta = \hat{\Theta}_t^\delta(X_1, \dots, X_t)$.
- Time-uniformity:

$$\mathbb{P}(\forall t \in \mathbb{N}, \theta \in \hat{\Theta}_t^\delta) \geq 1 - \delta \quad \text{or} \quad \mathbb{P}(\theta \in \hat{\Theta}_\tau^\delta) \geq 1 - \delta \text{ for } \tau \text{ stopping time.}$$
- Typical applications:



Stochastic bandits



Change-point detection

Setting: exponential families

Parametric family

$$p_\theta(x) = h(x) \exp(\langle \theta, F(x) \rangle - \mathcal{L}(\theta))$$

- $\Theta \subseteq \mathbb{R}^d$: open set,
- $F(x)$: feature function,
- $\mathcal{L}(\theta)$: log-partition function (convex, assume $\det \nabla^2 \mathcal{L}(\theta) > 0$ for all $\theta \in \Theta$).

Bregman divergence

$$\begin{aligned} \mathcal{B}_\mathcal{L}(\theta', \theta) &= \mathcal{L}(\theta') - \mathcal{L}(\theta) - \langle \theta' - \theta, \nabla \mathcal{L}(\theta) \rangle \\ &= KL(p_\theta \| p_{\theta'}). \end{aligned}$$

Standard parameter estimate

$$\hat{\theta}_t = \nabla \mathcal{L}^{-1}(\hat{\mu}_t), \text{ where } \hat{\mu}_t = \frac{1}{t} \sum_{s=1}^t F(X_s).$$

Example: Gaussian $\mathcal{N}(\mu, \sigma^2)$ with known variance

$$\theta = \mu, \Theta = \mathbb{R}, F(x) = \frac{x}{\sigma}, \mathcal{L}(\theta) = \frac{\theta^2}{2\sigma^2}, \mathcal{B}_\mathcal{L}(\theta', \theta) = \frac{(\theta' - \theta)^2}{2\sigma^2}.$$

Lemma 1 (log-Laplace control) For $\theta \in \Theta$ and λ s.t $\theta + \lambda \in \Theta$,

$$\log \mathbb{E}_\theta [e^{\langle \lambda, F(X) - \mathbb{E}_\theta[F(X)] \rangle}] = \mathcal{B}_\mathcal{L}(\theta + \lambda, \theta).$$

Lemma 2 (Bregman duality) For any $\alpha \in [0, 1]$,

$$\mathcal{B}_{\mathcal{L}, \theta'}^*(\alpha(\nabla \mathcal{L}(\theta) - \nabla \mathcal{L}(\theta'))) = \mathcal{B}_\mathcal{L}(\theta', \theta_\alpha),$$

where

$$\begin{aligned} \theta_\alpha &= \nabla \mathcal{L}^{-1}(\alpha \nabla \mathcal{L}(\theta) + (1 - \alpha) \nabla \mathcal{L}(\theta')), \\ \mathcal{B}_{\mathcal{L}, \theta'}^*(x) &= \sup_{\lambda} \langle \lambda, x \rangle - \mathcal{B}_\mathcal{L}(\theta' + \lambda, \theta'). \end{aligned}$$

Time-uniform Bregman deviation

Main theorem

Regularized parameter estimate

$$\hat{\theta}_{t,c}(\theta) = (\nabla \mathcal{L})^{-1} \left(\frac{t}{t+c} \hat{\mu}_t + \frac{c}{t+c} \nabla \mathcal{L}(\theta) \right),$$

Bregman information gain

$$\gamma_{t,c}(\theta) = \log \left(\frac{\int_{\Theta} \exp(-c \mathcal{B}_\mathcal{L}(\theta', \theta)) d\theta'}{\int_{\Theta} \exp(-(t+c) \mathcal{B}_\mathcal{L}(\theta', \hat{\theta}_{t,c}(\theta))) d\theta'} \right),$$

Time-uniform deviation

$$\mathbb{P} \left(\exists t \in \mathbb{N}, (t+c) \mathcal{B}_\mathcal{L}(\theta, \hat{\theta}_{t,c}(\theta)) \geq \log \frac{1}{\delta} + \gamma_{t,c}(\theta) \right) \leq \delta.$$

Remarks

- Valid for generic families, not just one-dimensional.
- $\gamma_{t,c}(\theta) = \frac{\dim \Theta}{2} \log(1 + \frac{t}{c}) + \mathcal{O}(1) \implies$ asymptotic confidence radius is $\propto \sqrt{\frac{\log t}{t}}$.
- Implicit confidence set in θ , but easy to compute numerically.
- Explicit instantiation to many classical families:
 \hookrightarrow Gaussian, Bernoulli, Exponential, Gamma, Weibull, Pareto, Poisson, χ^2 .

GLR test in exponential families

Change of measure detection: distribution of X_u is $p_{\theta(u)}$

$$\mathcal{H}_0 \text{ (null): } \exists \theta_0 \in \Theta, \forall u \in \mathbb{N}, \theta(u) = \theta_0 \quad \text{(no change),}$$

$$\mathcal{H}_1 \text{ (alt.): } \exists s \in \mathbb{N}, \theta_1, \theta_2 \in \Theta, \forall u \in \mathbb{N}, \theta(u) = \theta_1 \mathbb{1}_{u \leq s} + \theta_2 \mathbb{1}_{u > s} \quad \text{(change).}$$

Scan statistic

$$\hat{\theta}_{a:b} = \nabla \mathcal{L}^{-1} \left(\frac{1}{b-a+1} \sum_{s=a}^b F(X_u) \right).$$

Generalized Likelihood Ratio

$$\begin{aligned} G_{1:s:t} &= \inf_{\theta_0} \sup_{\theta_1, \theta_2} \log \left(\frac{\prod_{u=1}^s p_{\theta_1}(X_u) \prod_{u=s+1}^t p_{\theta_2}(X_u)}{\prod_{u=1}^t p_{\theta_0}(X_u)} \right) \\ &= \inf_{\theta_0} s \mathcal{B}_\mathcal{L}(\theta_0, \hat{\theta}_{1:s}) + (t-s) \mathcal{B}_\mathcal{L}(\theta_0, \hat{\theta}_{s+1:t}). \end{aligned}$$

Doubly time-uniform deviation $g(t) = (t+1) \log^2(t+1) / \log(2)$,

$$\mathbb{P}_\theta \left(\exists t \in \mathbb{N}, \exists s < t: (t-s+c) \mathcal{B}_\mathcal{L}(\theta, \hat{\theta}_{s+1:t,c}(\theta)) \geq \log \left(\frac{g(t)}{\delta} \right) + \gamma_{s+1:t,c}(\theta) \right) \leq \delta.$$

Regularized GLR test

$$\tau_{c,\delta} = \min \left\{ t \in \mathbb{N}: \exists s < t, \theta \in \Theta: (s+c) \mathcal{B}_\mathcal{L}(\theta, \hat{\theta}_{1:s,c}(\theta)) \geq \log \left(\frac{2}{\delta} \right) + \gamma_{1:s,c}(\theta) \right.$$

$$\left. \text{and } (t-s+c) \mathcal{B}_\mathcal{L}(\theta, \hat{\theta}_{s+1:t,c}(\theta)) \geq \log \left(\frac{2g(t)}{\delta} \right) + \gamma_{s+1:t,c}(\theta) \right\}$$

has a false alarm probability $\leq \delta$.

Sketch of proof

Martingale construction

Lemma 1 \implies for all suitable λ and an arbitrary $c > 0$,

$$M_t^\lambda = \exp \left(\left\langle \lambda, \sum_{s=1}^t F(X_s) - \mathbb{E}_\theta[F(X)] \right\rangle - t \mathcal{B}_\mathcal{L}(\theta + \lambda, \theta) \right) \text{ defines a } \geq 0 \text{ martingale.}$$

Martingale mixture For $c > 0$,

$$\begin{aligned} q_\theta(\lambda|c) &\propto \exp(\langle \theta + \lambda, c \nabla \mathcal{L}(\theta) \rangle - c \mathcal{L}(\theta)), \\ M_t &= \int M_t^\lambda q_\theta(\lambda|c) d\lambda. \end{aligned}$$

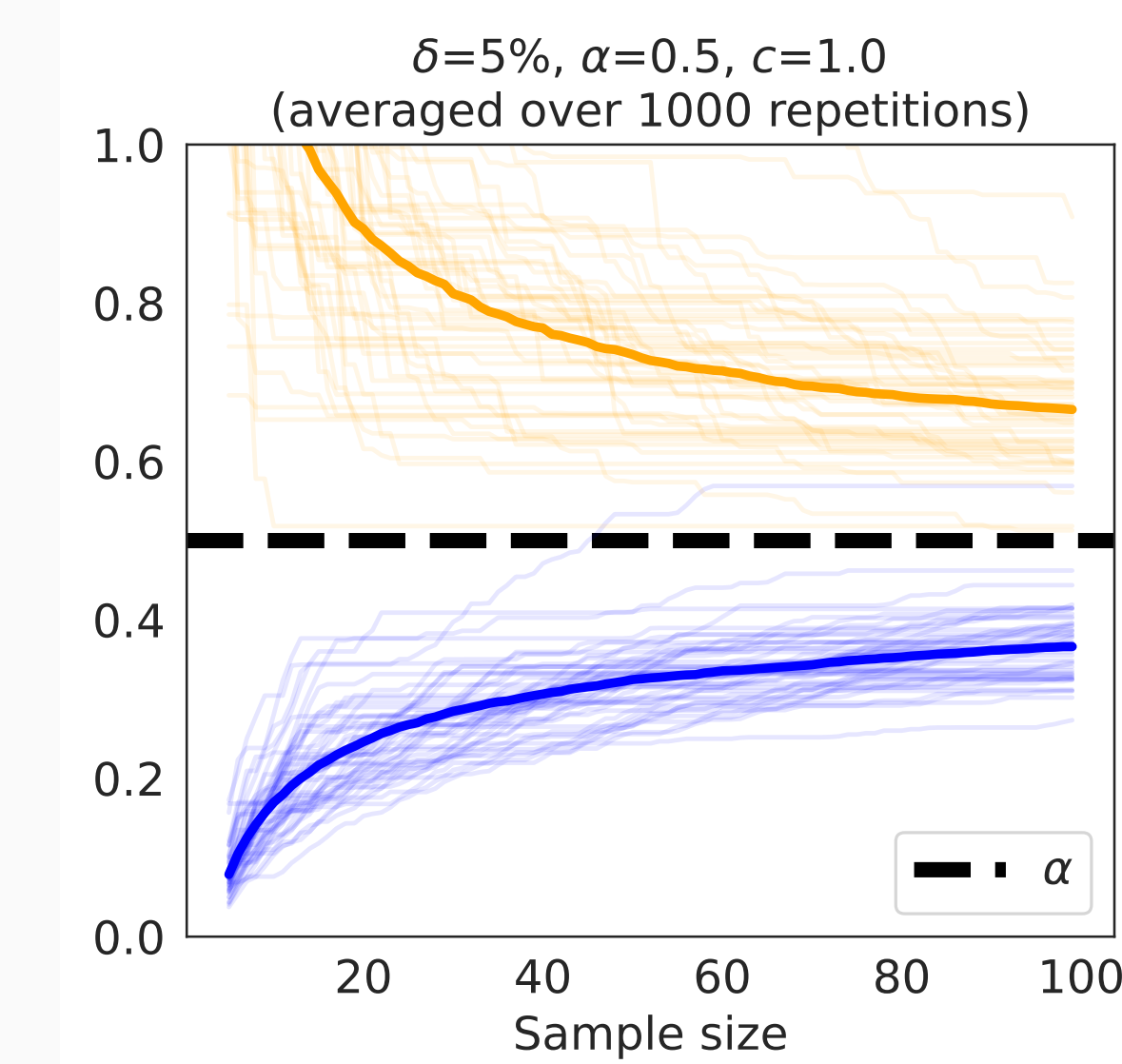
Rewriting

$$\text{Lemma 2 } \implies M_t = \exp \left((t+c) \mathcal{B}_\mathcal{L}(\theta, \hat{\theta}_{t,c}(\theta)) - \gamma_{t,c}(\theta) \right).$$

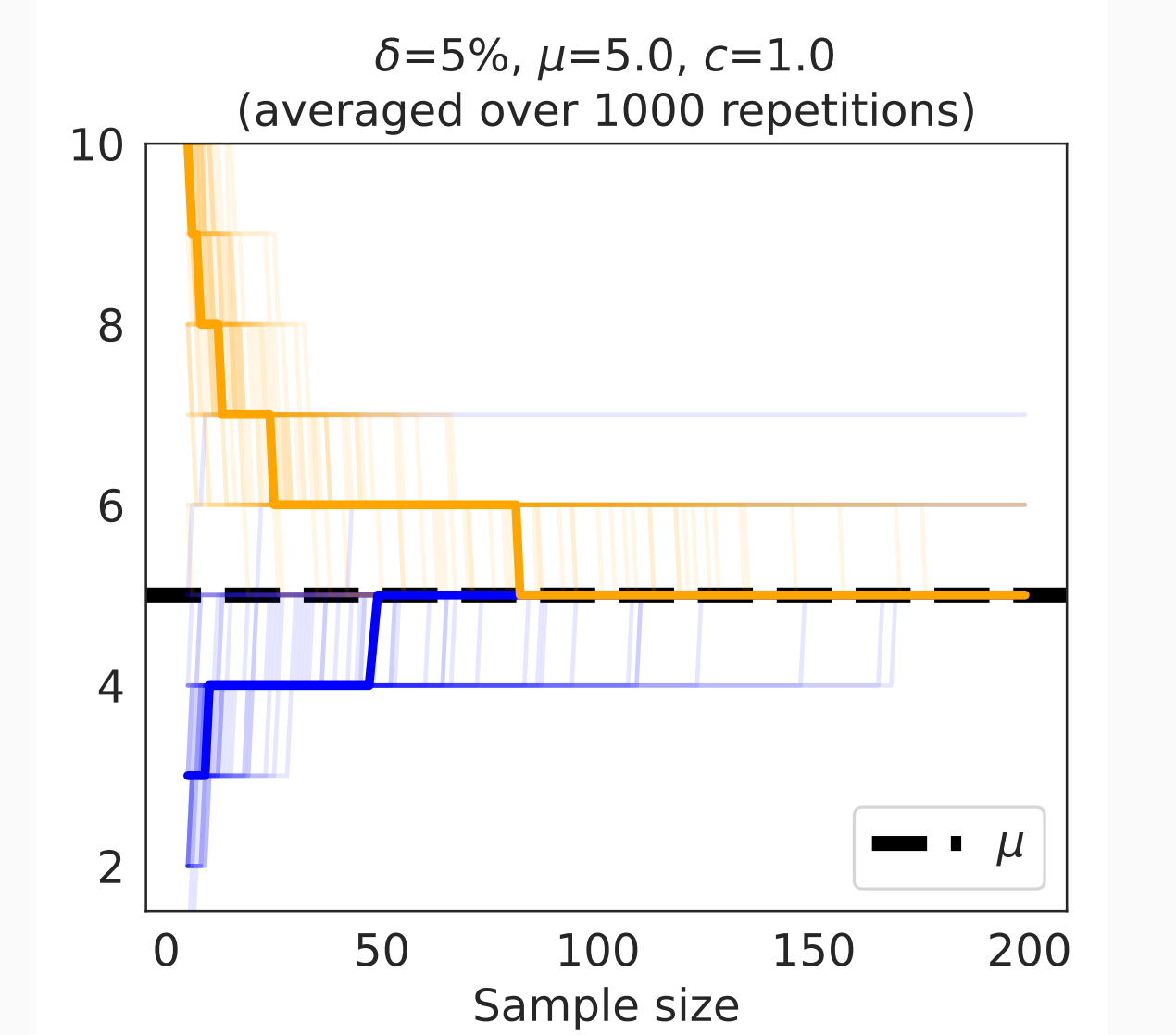
Conclusion

Ville's inequality (supermartingale + Doob's optional stopping). \square

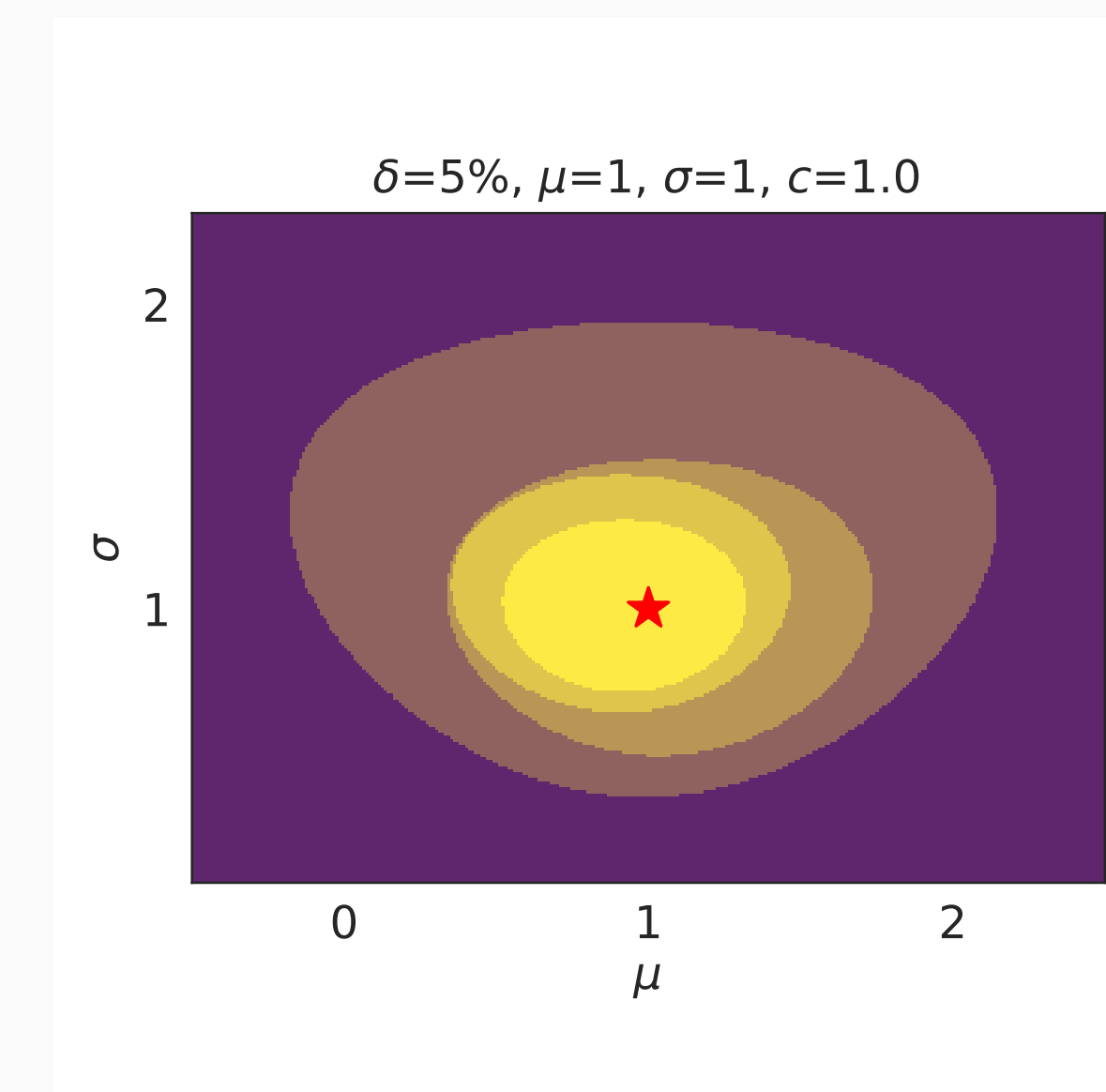
Numerical experiments



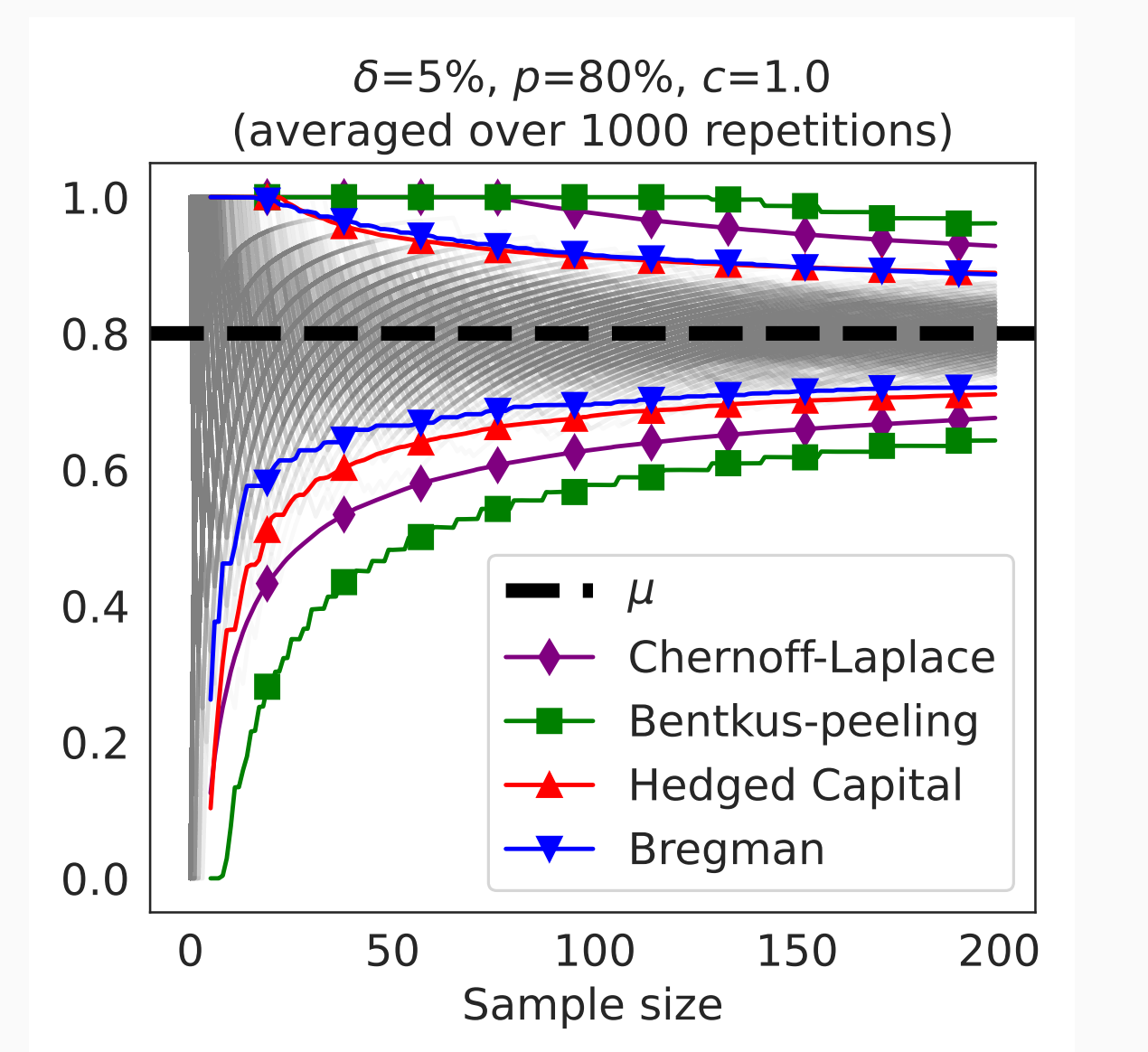
Pareto



Chi-square



Gaussian mean-variance
 $t \in \{10, 25, 50, 100\}$



Comparison of median confidence envelopes around the mean for Bernoulli $\mathcal{B}(0.8)$