

Controlling for confounders through approximate sufficiency

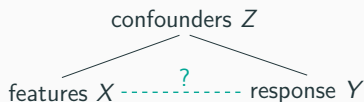
Rina Foygel Barber (joint with Lucas Janson)

<http://www.stat.uchicago.edu/~rina/>



Lucas Janson (Harvard U.)

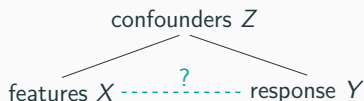
Intro: testing conditional independence



Classical (parametric) approach:

- Assume a parametric model such as $Y \mid X, Z \sim f(\cdot; \alpha^\top X + \beta^\top Z)$
- Parametric inference to test $H_0 : \alpha = 0$

Intro: testing conditional independence



Classical (parametric) approach:

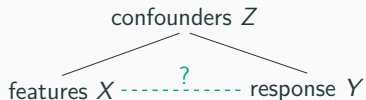
- Assume a parametric model such as $Y \mid X, Z \sim f(\cdot; \alpha^\top X + \beta^\top Z)$
- Parametric inference to test $H_0 : \alpha = 0$

Model-X approach a.k.a. Conditional Randomization Test (Candès et al 2018)

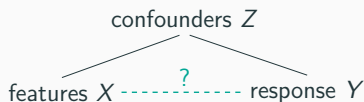
- Known distribution of $X \mid Z$ (distrib. of Y unknown)
- Choose function $T(X; Y, Z)$ that measures association
- Resample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)} \stackrel{\text{iid}}{\sim}$ (distrib. of $X \mid Z$)

$$\rightsquigarrow \text{pval} = \frac{1 + \sum_m \mathbb{1}\{T(\tilde{X}^{(m)}; Y, Z) \geq T(X; Y, Z)\}}{1 + M}$$

Intro: testing conditional independence



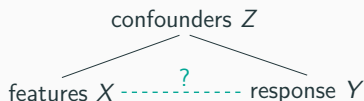
Intro: testing conditional independence



Model-X approach via sufficient statistics (Huang & Janson 2019)

- Distribution of $X \mid Z$ is only partially known
- By conditioning on sufficient statistic $S(X, Z)$,
can resample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)} \stackrel{\text{iid}}{\sim}$ (distrib. of $X \mid S(X, Z)$)
& compute p-value for test statistic T as before

Intro: testing conditional independence



Model-X approach via sufficient statistics (Huang & Janson 2019)

- Distribution of $X \mid Z$ is only partially known
- By conditioning on sufficient statistic $S(X, Z)$,
can resample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)} \stackrel{\text{iid}}{\sim} (\text{distrib. of } X \mid S(X, Z))$
& compute p-value for test statistic T as before
- Example: canonical GLMs
 - $X_i \sim \exp \{X_i \cdot Z_i^\top \theta - a(Z_i^\top \theta)\}$, $i = 1, \dots, n$, with θ unknown
 - $S(X, Z) = \sum_i X_i Z_i$ is suff. stat. for $X = (X_1, \dots, X_n)$

Intro: testing goodness-of-fit (GoF)

More generally...

Goodness-of-fit test

Testing $H_0: X \sim P_\theta$ for some $\theta \in \Theta$,
where $\{P_\theta : \theta \in \Theta\}$ is a parametric family

Intro: testing goodness-of-fit (GoF)

More generally...

Goodness-of-fit test

Testing $H_0: X \sim P_\theta$ for some $\theta \in \Theta$,
where $\{P_\theta : \theta \in \Theta\}$ is a parametric family

Conditional independence testing can be a special case:

- Assume $X \mid Z \sim P_\theta(\cdot|Z)$ for some $\theta \in \Theta$
- Null hypothesis $H_0 : X \perp\!\!\!\perp Y \mid Z$
- Equivalently... $H_0: X \mid Y, Z \sim P_\theta(\cdot|Z)$ for some $\theta \in \Theta$
- Note: we condition on Y and Z (i.e., treat as fixed)

Intro: testing goodness-of-fit (GoF)

A general framework:

- Choose any test statistic $T : \mathcal{X} \rightarrow \mathbb{R}$
- Draw copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$
- Compute rank-based p-value

$$\text{pval} = \frac{1 + \sum_m \mathbb{1}\{T(\tilde{X}^{(m)}) \geq T(X)\}}{1 + M}$$

- If $X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ are exchangeable under $H_0 \rightsquigarrow$ p-value is valid

Co-sufficient sampling (CSS)

Co-sufficient sampling

Sample copies $\tilde{X}^{(m)} \sim (\text{distrib. of } X \mid S(X))$,
where $S(X)$ is a sufficient statistic for the family $\{P_\theta : \theta \in \Theta\}$

Can be applied to:

1. Test goodness-of-fit (GoF)
(Engen & Lillegård 1997, Lockhart et al 2007, Stephens 2012, Hazra 2013)
2. Test conditional independence (special case of GoF)
(Rosenbaum 1984, Kolassa 2003, Huang & Janson 2019)
3. Construct conf. intervals for a parameter of interest
(by inverting GoF tests)

Co-sufficient sampling (CSS)

Co-sufficient sampling

Sample copies $\tilde{X}^{(m)} \sim (\text{distrib. of } X \mid S(X))$,

where $S(X)$ is a sufficient statistic for the family $\{P_\theta : \theta \in \Theta\}$

Co-sufficient sampling (CSS)

Co-sufficient sampling

Sample copies $\tilde{X}^{(m)} \sim (\text{distrib. of } X \mid S(X))$,
where $S(X)$ is a sufficient statistic for the family $\{P_\theta : \theta \in \Theta\}$

Permutation tests are an example of CSS

- $H_0: X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{D}$ for $\mathcal{D} \in (\text{some set})$
- The order statistics $X_{(1)} \leq \dots \leq X_{(n)}$ are sufficient under the null
- Permutation test \Leftrightarrow resampling X conditional on order statistics
- Application: testing $X \perp\!\!\!\perp Y$
 H_0 : conditional on Y_1, \dots, Y_n , it holds that X_1, \dots, X_n are i.i.d.

Co-sufficient sampling (CSS)

Limitation of co-sufficient sampling... no power in many settings!

Example—logistic model:

- $X = (X_1, \dots, X_n) \in \{0, 1\}^n$, $Z = (Z_1, \dots, Z_n) \in (\mathbb{R}^k)^n$
- If the Z_i 's are in general position,
then $\sum_i X_i Z_i \in \mathbb{R}^k$ uniquely determines X

(so if we resample, will have $\tilde{X}^{(1)} = \dots = \tilde{X}^{(M)} = X \rightsquigarrow$ zero power)

Co-sufficient sampling (CSS)

Limitation of co-sufficient sampling... no power in many settings!

Co-sufficient sampling (CSS)

Limitation of co-sufficient sampling... no power in many settings!

For many other models, the minimal sufficient statistic $S(X)$ is essentially the data itself, e.g.,

- Mixture of Gaussians or mixture of GLMs
- Non-canonical GLMs
- Heavy tailed distributions (e.g., multivariate t)
- Models with missing or corrupted data

Approximate sufficiency

For a family $\{P_\theta : \theta \in \Theta\}$, a function $S(X)$ is a *sufficient statistic* if
(distrib. of $X \mid S(X)$, $X \sim P_\theta$) = (distrib. of $X \mid S(X)$, $X \sim P_{\theta'}$) $\forall \theta, \theta'$.

Asymptotic sufficiency: (Le Cam, Wald, ...)

Informally...

(distrib. of $X \mid S(X)$, $X \sim P_\theta$) \approx (distrib. of $X \mid S(X)$, $X \sim P_{\theta'}$) $\forall \theta, \theta'$.

- Under regularity conditions, $S(X) = \hat{\theta}_{\text{MLE}}(X)$ is asymp. suff.

Approximate co-sufficient sampling (aCSS)

Main idea:

- Let $\hat{\theta} \in \Theta$ be an approximate MLE given the data X
- Let $p_{\theta}(\cdot|\hat{\theta}) = \text{distrib. of } X \mid \hat{\theta}$, if marginally $X \sim P_{\theta}$
 \rightsquigarrow under the null, $X \mid \hat{\theta} \sim p_{\theta_0}(\cdot|\hat{\theta})$ for the unknown true θ_0
- Sample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ from $\underbrace{p_{\hat{\theta}}(\cdot|\hat{\theta}) \approx p_{\theta_0}(\cdot|\hat{\theta})}_{\text{by approx. sufficiency}}$

$X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)} \approx \text{exchangeable under } H_0 \rightsquigarrow \text{p-value is } \approx \text{valid}$

Approximate co-sufficient sampling (aCSS)

Distance to exchangeability

$$d_{\text{exch}}(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}) = \inf_{\substack{\text{Exch. distrib.} \\ \mathcal{D} \text{ on } \mathcal{X}^{M+1}}} \left\{ d_{\text{TV}}\left((X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}), \mathcal{D}\right) \right\}$$

For any test statistic $T(X)$, the p-value

$$\text{pval} = \frac{1 + \sum_m \mathbb{1}\{T(\tilde{X}^{(m)}) \geq T(X)\}}{1 + M}$$

satisfies

$$\mathbb{P}\{\text{pval} \leq \alpha\} \leq \alpha + d_{\text{exch}}(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}).$$

- Step 1: choose a test statistic $T : \mathcal{X} \rightarrow \mathbb{R}$
- Step 2: observe data X , and compute an approximate MLE $\hat{\theta}$
- Step 3: sample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ from \approx distribution of $X \mid \hat{\theta}$
- Step 4: compute a rank-based p-value to test H_0 :

$$\text{pval} = \frac{1 + \sum_m \mathbb{1}\{T(\tilde{X}^{(m)}) \geq T(X)\}}{1 + M}$$

aCSS algorithm

- Step 1: choose a test statistic $T : \mathcal{X} \rightarrow \mathbb{R}$
- Step 2: observe data X , and compute an approximate MLE $\hat{\theta}$
- Step 3: sample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ from \approx distribution of $X \mid \hat{\theta}$
- Step 4: compute a rank-based p-value to test H_0 :

$$\text{pval} = \frac{1 + \sum_m \mathbb{1}\{T(\tilde{X}^{(m)}) \geq T(X)\}}{1 + M}$$

aCSS algorithm

- Step 2: observe data X , and compute an approximate MLE $\hat{\theta}$

Ideally would like to minimize

$$\mathcal{L}(\theta; X, W) = \underbrace{\mathcal{L}(\theta; X)}_{\substack{\text{penalized neg. log-likelihood} \\ -\log f(X; \theta) + \mathcal{R}(\theta)}} + \underbrace{\sigma \cdot W^\top \theta}_{\substack{\text{perturb with } W \sim \mathcal{N}(0, \frac{1}{d} \mathbf{I}_d) \\ \text{(choose } \sigma \ll n^{1/2})}}$$

(see also Tian & Taylor 2018—random perturbation for selective inference)

- Step 2: observe data X , and compute an approximate MLE $\hat{\theta}$

Ideally would like to minimize

$$\mathcal{L}(\theta; X, W) = \underbrace{\mathcal{L}(\theta; X)}_{\substack{\text{penalized neg. log-likelihood} \\ -\log f(X; \theta) + \mathcal{R}(\theta)}} + \underbrace{\sigma \cdot W^\top \theta}_{\substack{\text{perturb with } W \sim \mathcal{N}(0, \frac{1}{d} \mathbf{I}_d) \\ \text{(choose } \sigma \ll n^{1/2})}}$$

(see also Tian & Taylor 2018—random perturbation for selective inference)

But... what if nonconvex? what if no global minimum?

- Function $\hat{\theta} : \mathcal{X} \times \mathbb{R}^d \rightarrow \Theta$, returns $\hat{\theta}(X, W)$.
- If $\hat{\theta}(X, W)$ is a strict SOSP of $\mathcal{L}(\theta; X, W)$, proceed to next step.
- Otherwise return $\tilde{X}^{(1)} = \dots = \tilde{X}^{(M)} = X \rightsquigarrow \text{pval} = 1$.


- Step 3: sample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ from \approx distribution of $X \mid \hat{\theta}$

aCSS algorithm

- Step 3: sample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ from \approx distribution of $X \mid \hat{\theta}$

Density of $X \mid \hat{\theta}$, conditional on the event that $\hat{\theta}(X, W)$ is strict SOSp:


$$\propto f(x; \theta_0) \cdot \exp \left\{ -\frac{\|\nabla_{\theta} \mathcal{L}(\hat{\theta}; x)\|}{2\sigma^2/d} \right\} \cdot \det \left(\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}; x) \right) \cdot \mathbb{1}_{x \in \mathcal{X}_{\hat{\theta}}}$$

 support of $X \mid \hat{\theta}$

- Step 3: sample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ from \approx distribution of $X \mid \hat{\theta}$

Density of $X \mid \hat{\theta}$, conditional on the event that $\hat{\theta}(X, W)$ is strict SOSP:

$$\propto f(x; \theta_0) \cdot \exp \left\{ -\frac{\|\nabla_{\theta} \mathcal{L}(\hat{\theta}; x)\|}{2\sigma^2/d} \right\} \cdot \det \left(\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}; x) \right) \cdot \mathbb{1}_{x \in \mathcal{X}_{\hat{\theta}}}$$

 support of $X \mid \hat{\theta}$


θ_0 unknown \rightsquigarrow use $\hat{\theta}$ as plug-in estimate:

$$\propto f(x; \hat{\theta}) \cdot \exp \left\{ -\frac{\|\nabla_{\theta} \mathcal{L}(\hat{\theta}; x)\|}{2\sigma^2/d} \right\} \cdot \det \left(\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}; x) \right) \cdot \mathbb{1}_{x \in \mathcal{X}_{\hat{\theta}}}$$

- Step 3: sample copies $\tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}$ from \approx distribution of $X \mid \hat{\theta}$

Density of $X \mid \hat{\theta}$, conditional on the event that $\hat{\theta}(X, W)$ is strict SOSP:

$$\propto f(x; \theta_0) \cdot \exp \left\{ -\frac{\|\nabla_{\theta} \mathcal{L}(\hat{\theta}; x)\|}{2\sigma^2/d} \right\} \cdot \det \left(\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}; x) \right) \cdot \mathbb{1}_{x \in \mathcal{X}_{\hat{\theta}}}$$


support of $X \mid \hat{\theta}$

θ_0 unknown \rightsquigarrow use $\hat{\theta}$ as plug-in estimate:

$$\propto f(x; \hat{\theta}) \cdot \exp \left\{ -\frac{\|\nabla_{\theta} \mathcal{L}(\hat{\theta}; x)\|}{2\sigma^2/d} \right\} \cdot \det \left(\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}; x) \right) \cdot \mathbb{1}_{x \in \mathcal{X}_{\hat{\theta}}}$$

If sampling directly is impossible,

can use an exchangeable form of MCMC (Besag & Clifford 1989)

Assumption 1: regularity conditions

- $\Theta \subseteq \mathbb{R}^d$ convex & open
- P_θ has positive density $f(\cdot; \theta)$ w.r.t. base measure $\nu_{\mathcal{X}}$ for all $\theta \in \Theta$
- Log-likelihood $\log f(x; \theta)$ & penalty $\mathcal{R}(\theta)$ are continuously twice diff.

Type I error guarantee

Assumption 2: approximate MLE

For $X \sim P_{\theta_0}$ and $W \sim \mathcal{N}(0, \frac{1}{d}I_d)$, with prob. at least $1 - \delta$,
 $\|\hat{\theta}(X, W) - \theta_0\| \leq r$ and $\hat{\theta}(X, W)$ is a strict SOSP of $\mathcal{L}(\theta; X, W)$.

Assumption 3: Hessian of the log-likelihood

$$\mathbb{E} \left[\exp \left\{ \sup_{\theta \in \mathbb{B}(\theta_0, r) \cap \Theta} r^2 \|\nabla^2 \log f(X; \theta) - \mathbb{E} [\nabla^2 \log f(X; \theta)]\| \right\} \right] \leq e^\varepsilon$$

Type I error guarantee

Assumption 2: approximate MLE

For $X \sim P_{\theta_0}$ and $W \sim \mathcal{N}(0, \frac{1}{d} \mathbf{I}_d)$, with prob. at least $1 - \delta$,
 $\|\hat{\theta}(X, W) - \theta_0\| \leq r$ and $\hat{\theta}(X, W)$ is a strict SOSP of $\mathcal{L}(\theta; X, W)$.

Assumption 3: Hessian of the log-likelihood

$$\mathbb{E} \left[\exp \left\{ \sup_{\theta \in \mathbb{B}(\theta_0, r) \cap \Theta} r^2 \|\nabla^2 \log f(X; \theta) - \mathbb{E} [\nabla^2 \log f(X; \theta)]\| \right\} \right] \leq e^\varepsilon$$

In standard settings with n independent observations...

$$r, \varepsilon, \delta = \tilde{\mathcal{O}}(n^{-1/2})$$

Type I error guarantee

Theorem

Under Assumptions 1, 2, & 3, the copies produced by aCSS satisfy

$$d_{\text{exch}}(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}) \leq 3\sigma r + \delta + \varepsilon$$

under H_0 .

Therefore, for any test statistic T , Type I error for testing H_0 satisfies

$$\mathbb{P}\{\text{pval} \leq \alpha\} \leq \alpha + 3\sigma r + \delta + \varepsilon$$

Type I error guarantee

Theorem

Under Assumptions 1, 2, & 3, the copies produced by aCSS satisfy

$$d_{\text{exch}}(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}) \leq 3\sigma r + \delta + \varepsilon$$

under H_0 .

Therefore, for any test statistic T , Type I error for testing H_0 satisfies

$$\mathbb{P}\{\text{pval} \leq \alpha\} \leq \alpha + 3\sigma r + \delta + \varepsilon$$



Excess Type I error should be $o(1)$...

- $r, \delta, \varepsilon \asymp n^{-1/2}$ from the assumptions
- σ = noise level, chosen by analyst
→ choose $\sigma \asymp n^c$ for some $c \in [0, \frac{1}{2})$

Examples

Examples where CSS has no power, but aCSS assumptions hold:

- Canonical GLMs such as logistic regression (low-dim.):

$$X_i \stackrel{\parallel}{\sim} \text{Bernoulli} \left(\frac{e^{Z_i^\top \beta}}{1 + e^{Z_i^\top \beta}} \right) \text{ for unknown } \beta$$

- Two-sample difference-of-means (the Behrens–Fisher problem):

$$X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_X, \sigma_X^2), \quad Y_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_Y, \sigma_Y^2), \quad \text{test } H_0 : \mu_X = \mu_Y$$

(An aCSS-like approach for this problem was considered by Lillegård 2001)

Examples

Examples where CSS has no power, but aCSS assumptions hold:

- Spatial process on integer lattice: for unknown ρ ,

$$X \sim \mathcal{N}(0, \Sigma) \text{ where } \Sigma_{ij} = \rho^{D_{ij}} \text{ for known pairwise distances } D_{ij}$$

- Multivariate t distribution (low-dim.):

$$X_i \stackrel{\text{iid}}{\sim} t_\gamma(0, \Sigma) \text{ for known } \gamma \text{ \& unknown } \Sigma$$

- And maybe missing data, latent variables, and more ...

Simulations

Compare to oracle method that knows θ_0 :

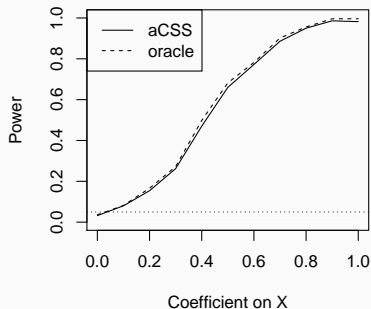
- Sample copies $\tilde{X}^{(m)} \stackrel{\text{iid}}{\sim} P_{\theta_0}$
- Compute p-value with same statistic $T(x)$

Simulations

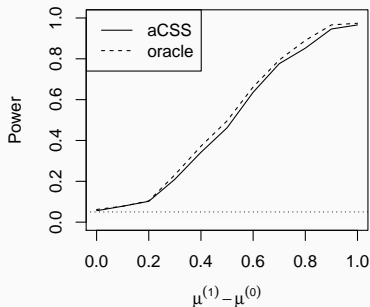
Compare to oracle method that knows θ_0 :

- Sample copies $\tilde{X}^{(m)} \stackrel{\text{iid}}{\sim} P_{\theta_0}$
- Compute p-value with same statistic $T(x)$

Logistic Regression



Behrens-Fisher

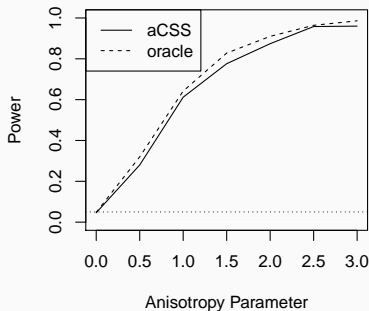


Simulations

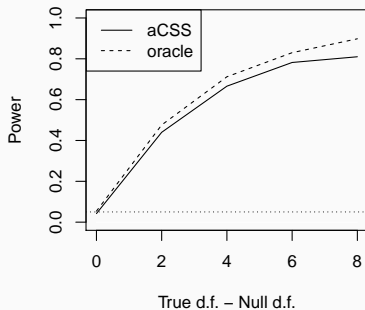
Compare to oracle method that knows θ_0 :

- Sample copies $\tilde{X}^{(m)} \stackrel{\text{iid}}{\sim} P_{\theta_0}$
- Compute p-value with same statistic $T(x)$

Gaussian Spatial



Multivariate t



Sampling

Recall: need to sample copies $\tilde{X}^{(m)}$ from

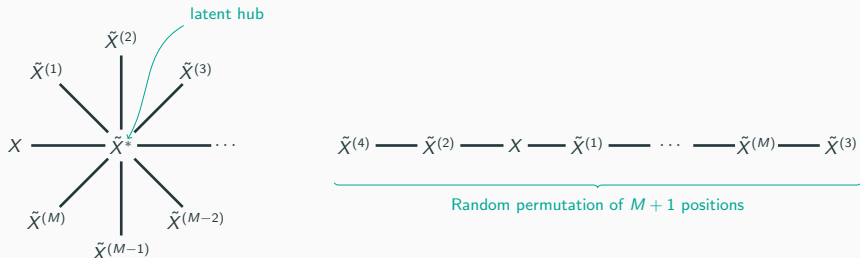
$$\propto f(x; \hat{\theta}) \cdot \exp \left\{ -\frac{\|\nabla_{\theta} \mathcal{L}(\hat{\theta}; x)\|}{2\sigma^2/d} \right\} \cdot \det \left(\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}; x) \right) \cdot \mathbb{1}_{x \in \mathcal{X}_{\hat{\theta}}}$$

Sampling

Recall: need to sample copies $\tilde{X}^{(m)}$ from

$$\propto f(x; \hat{\theta}) \cdot \exp \left\{ -\frac{\|\nabla_{\theta} \mathcal{L}(\hat{\theta}; x)\|}{2\sigma^2/d} \right\} \cdot \det \left(\nabla_{\theta}^2 \mathcal{L}(\hat{\theta}; x) \right) \cdot \mathbb{1}_{x \in \mathcal{X}_{\hat{\theta}}}$$

Two exchangeable MCMC strategies (Besag & Clifford 1989)



- Run Metropolis–Hastings, where $f(x; \hat{\theta})$ stationary for proposal distrib.
- e.g., if X consists of n indep. observations (i.e., $f(x; \hat{\theta}) = \prod_{i=1}^n f_i(x_i; \hat{\theta})$), can choose proposal distrib. = resample s of n observations

Proof sketch for Theorem

Need to bound $d_{\text{exch}}(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)})$

(1) Calculate joint distribution:

$$\begin{cases} \hat{\theta} & \sim (\text{marginal distrib. of } \hat{\theta}) \\ X \mid \hat{\theta} & \sim p_{\theta_0}(\cdot \mid \hat{\theta}) \\ \tilde{X}^{(m)} \mid X, \hat{\theta} & \sim p_{\hat{\theta}}(\cdot \mid \hat{\theta}) \end{cases}$$

$$\implies d_{\text{exch}}(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}) \leq \mathbb{E}_{\hat{\theta}} \left[d_{\text{TV}}(p_{\theta_0}(\cdot \mid \hat{\theta}), p_{\hat{\theta}}(\cdot \mid \hat{\theta})) \right]$$

Proof sketch for Theorem

(2) To bound d_{TV} :

$$\frac{p_{\hat{\theta}}(X|\hat{\theta})}{p_{\theta_0}(X|\hat{\theta})} \propto \frac{f(X;\hat{\theta})}{f(X;\theta_0)} \Rightarrow \frac{p_{\hat{\theta}}(X|\hat{\theta})}{p_{\theta_0}(X|\hat{\theta})} = \frac{\frac{f(X;\hat{\theta})}{f(X;\theta_0)}}{\mathbb{E}_{p_{\theta_0}(\cdot|\hat{\theta})} \left[\frac{f(X;\hat{\theta})}{f(X;\theta_0)} \right]}$$

Proof sketch for Theorem

(2) To bound d_{TV} :

$$\begin{aligned} \frac{p_{\hat{\theta}}(X|\hat{\theta})}{p_{\theta_0}(X|\hat{\theta})} &\propto \frac{f(X;\hat{\theta})}{f(X;\theta_0)} \Rightarrow \frac{p_{\hat{\theta}}(X|\hat{\theta})}{p_{\theta_0}(X|\hat{\theta})} = \frac{\frac{f(X;\hat{\theta})}{f(X;\theta_0)}}{\mathbb{E}_{p_{\theta_0}(\cdot|\hat{\theta})} \left[\frac{f(X;\hat{\theta})}{f(X;\theta_0)} \right]} \\ \Rightarrow d_{TV}(p_{\theta_0}(\cdot|\hat{\theta}), p_{\hat{\theta}}(\cdot|\hat{\theta})) &= \mathbb{E}_{p_{\theta_0}(\cdot|\hat{\theta})} \left[\left(1 - \frac{\frac{f(X;\hat{\theta})}{f(X;\theta_0)}}{\mathbb{E}_{p_{\theta_0}(\cdot|\hat{\theta})} \left[\frac{f(X;\hat{\theta})}{f(X;\theta_0)} \right]} \right)_+ \right] \end{aligned}$$

So, we need to show that $\frac{f(X;\hat{\theta})}{f(X;\theta_0)}$ is \approx constant over distrib. $X|\hat{\theta}$.

Proof sketch for Theorem

$$\log \left(\frac{f(X; \hat{\theta})}{f(X; \theta_0)} \right) = -(\theta_0 - \hat{\theta})^\top \nabla_{\theta} \log f(X; \hat{\theta}) - \frac{1}{2} (\theta_0 - \hat{\theta})^\top \nabla_{\theta}^2 \log f(X; \tilde{\theta}) (\theta_0 - \hat{\theta})$$

Proof sketch for Theorem

$$\log \left(\frac{f(X; \hat{\theta})}{f(X; \theta_0)} \right) = -(\theta_0 - \hat{\theta})^\top \nabla_{\theta} \log f(X; \hat{\theta}) - \frac{1}{2} (\theta_0 - \hat{\theta})^\top \nabla_{\theta}^2 \log f(X; \tilde{\theta}) (\theta_0 - \hat{\theta})$$

$$\begin{aligned} \Rightarrow & \left| \log \left(\frac{f(X; \hat{\theta})}{f(X; \theta_0)} \right) + \frac{1}{2} (\theta_0 - \hat{\theta})^\top \mathbb{E}_{\theta_0} \left[\nabla_{\theta}^2 \log f(X; \tilde{\theta}) \right] (\theta_0 - \hat{\theta}) \right| \\ \leq & r \cdot \underbrace{\| \nabla_{\theta} \log f(X; \hat{\theta}) \|}_{= \sigma \|W\| \asymp \sigma} + \frac{1}{2} \cdot r^2 \underbrace{\left\| \nabla_{\theta}^2 \log f(X; \tilde{\theta}) - \mathbb{E}_{\theta_0} \left[\nabla_{\theta}^2 \log f(X; \tilde{\theta}) \right] \right\|}_{\asymp \varepsilon \text{ by Asm. 3}} \end{aligned}$$

$\| \theta_0 - \hat{\theta} \| \leq r$
with prob. $\geq 1 - \delta$ by Asm. 2

Proof sketch for Theorem

$$\log \left(\frac{f(X; \hat{\theta})}{f(X; \theta_0)} \right) = -(\theta_0 - \hat{\theta})^\top \nabla_{\theta} \log f(X; \hat{\theta}) - \frac{1}{2} (\theta_0 - \hat{\theta})^\top \nabla_{\theta}^2 \log f(X; \tilde{\theta}) (\theta_0 - \hat{\theta})$$

$$\begin{aligned} \Rightarrow & \left| \log \left(\frac{f(X; \hat{\theta})}{f(X; \theta_0)} \right) + \frac{1}{2} (\theta_0 - \hat{\theta})^\top \mathbb{E}_{\theta_0} \left[\nabla_{\theta}^2 \log f(X; \tilde{\theta}) \right] (\theta_0 - \hat{\theta}) \right| \\ \leq & r \cdot \underbrace{\left\| \nabla_{\theta} \log f(X; \hat{\theta}) \right\|}_{=\sigma \|W\| \asymp \sigma} + \frac{1}{2} \cdot r^2 \underbrace{\left\| \nabla_{\theta}^2 \log f(X; \tilde{\theta}) - \mathbb{E}_{\theta_0} \left[\nabla_{\theta}^2 \log f(X; \tilde{\theta}) \right] \right\|}_{\asymp \varepsilon \text{ by Asm. 3}} \end{aligned}$$

\nearrow
 $\|\theta_0 - \hat{\theta}\| \leq r$
with prob. $\geq 1 - \delta$ by Asm. 2

Rearrange \rightsquigarrow

$$d_{\text{exch}}(X, \tilde{X}^{(1)}, \dots, \tilde{X}^{(M)}) \leq \mathbb{E}_{\hat{\theta}} \left[d_{\text{TV}}(p_{\theta_0}(\cdot | \hat{\theta}), p_{\hat{\theta}}(\cdot | \hat{\theta})) \right] \leq 3\sigma r + \delta + \varepsilon$$

Summary & open questions

- Summary: aCSS can test goodness-of-fit by sampling nearly-exchangeable copies of the data, in a much broader range of settings than CSS

Summary & open questions

- Summary: aCSS can test goodness-of-fit by sampling nearly-exchangeable copies of the data, in a much broader range of settings than CSS
- How to choose σ to balance Type I error & power?
- Connections to Bayesian methods?
- Apply to high dimensional regression / covariance estimation?
- Apply to missing data / latent variables / models with singularities?
- Extend to model-X knockoffs?

Thank you!