Computational advances in large n & large p'' sparse Bayesian regression for binary & survival outcomes

 $\begin{array}{l} \mbox{Aki Nishimura}^1 \\ \mbox{joint work with Marc Suchard} \end{array}$

Department of Biomathematics, University of California - Los Angeles ¹

November 27, 2018

Observational Health Data Sciences and Informatics





COLUMBIA UNIVERSITY IN THE CITY OF NEW YORK





- 84 organizations
- 64 databases
- 1.26 billion patient records in total
- Open source analytic suite:
 - github.com/ohdsi
- More info at ohdsi.org

"Large n and large p" regression for observational data

Example study based on observational data from healthcare databases:

- Goal compare two different drugs (e.g. *dabigatran* and *warfarin*)
 - Efficacy of treatment? (e.g. prevention of blood clot formation)
 - Risk of serious side effects? (e.g. bleeding inside brain)?

"Large n and large p" regression for observational data

Example study based on observational data from healthcare databases:

- Goal compare two different drugs (e.g. *dabigatran* and *warfarin*)
 - Efficacy of treatment? (e.g. prevention of blood clot formation)
 - Risk of serious side effects? (e.g. bleeding inside brain)?



"Large n and large p" regression for observational data

OHDSI's approach addresses "replication crisis" (we try, at least):



- Scale of data:
 - Sample size: $n \approx 10^5 \sim 10^6$ (e.g. n = 72,489)
 - Number of features: $p \approx 10^4 \sim 10^5$ (e.g. p = 22,175); pre-existing conditions, prior treatments / drugs taken, and etc.

- Scale of data:
 - Sample size: $n \approx 10^5 \sim 10^6$ (e.g. n = 72,489)
 - Number of features: p ≈ 10⁴ ~ 10⁵ (e.g. p = 22,175); pre-existing conditions, prior treatments / drugs taken, and etc.
- Data characteristics:
 - Design matrix X is sparse, a small fraction of non-zero entries.
 - "Positive" outcome y_i = 1 is rare when regressing on serious adverse events (e.g. 192 out of 72,489)

- Two stage estimation of treatment (dabigatran) effect:
 - Propensity score estimation
 - regress treatment indicator on predictors
 - Doubly-robust treatment effect estimation
 - regress outcome on treatment indicator
 - $+\ propensity\ score\ strata\ indicators\ +\ predictors$

Table of Contents

1 Case for Bayesian approach in search of weak signals

- 2 Global-local shrinkage prior for sparse Bayesian regression
- Conjugate gradient (CG) sampler for multivariate Gaussians
- Preconditioning for rapid CG convergence in sparse regression
- 5 Numerical results / Applications
- 6 Sparse Bayesian regression for survival outcomes

7 Summary

calibrates hyper-parameters via their marginal likelihoods.
 ▶ vs. resampling / cross-validation methods

 — suspicious quality for rare outcomes
 (e.g. 192 cases out of n = 72,489 & p = 22,175)

- calibrates hyper-parameters via their marginal likelihoods.
 - vs. resampling / cross-validation methods

- suspicious quality for rare outcomes

(e.g. 192 cases out of n = 72,489 & p = 22,175)

- tends to have better non-asymptotic properties
- provides model selection uncertainty
 - vs. post-selection inference
- can incorporate structures behind given data
 - (e.g. hierarchical modeling across different hospitals)

① Case for Bayesian approach in search of weak signals

- 2 Global-local shrinkage prior for sparse Bayesian regression
- 3 Conjugate gradient (CG) sampler for multivariate Gaussians
- 4 Preconditioning for rapid CG convergence in sparse regression
- 5 Numerical results / Applications
- 6 Sparse Bayesian regression for survival outcomes

7 Summary

Global-local shrinkage prior for sparse Bayesian regression

- Global-local shrinkage of regression coefficients β :
 - Prior of the form $\beta_j | \tau, \lambda_j \sim \mathcal{N}\left(0, \tau^2 \lambda_j^2\right)$.
 - τ and λ_j are called global and local shrinkage parameters.
 - ▶ $\hat{\tau} \ll 1$, but can have $\hat{\tau}\hat{\lambda}_j = O(1)$ for a small number of j's by virtue of heavy tailed priors on λ_j

Global-local shrinkage prior for sparse Bayesian regression

- Global-local shrinkage of regression coefficients β :
 - Prior of the form $\beta_j | \tau, \lambda_j \sim \mathcal{N}\left(0, \tau^2 \lambda_j^2\right)$.
 - τ and λ_j are called global and local shrinkage parameters.
 - $\hat{\tau} \ll 1$, but can have $\hat{\tau}\hat{\lambda}_j = O(1)$ for a small number of j's by virtue of heavy tailed priors on λ_j
- Example Bayesian bridge (Polson et al., 2014):



 $\lambda_j \sim \text{alpha-stable distribution}$ with index of stability $\alpha/2$ $\label{eq:alpha-stable}$ $\pi(\beta_j \mid \tau) \propto \tau^{-1} \exp ig(- |\beta_j / \tau|^{lpha} ig)$

Gibbs sampler for sparse Bayesian logistic regression

- Computation for Bayesian logistic regression is commonly based on the Polya-Gamma data augmentation scheme of Polson et al. (2013).
- Through an auxiliary parameter ω , the conditional likelihood of a binary outcome y becomes

$$y'_i | \boldsymbol{X}, \boldsymbol{\beta}, \boldsymbol{\omega} \sim \mathcal{N}(\boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}, \omega_i^{-1}) \text{ for } y'_i := \omega_i^{-1} (y_i - 1/2)$$
 (1)

Gibbs sampler for sparse Bayesian logistic regression

- Computation for Bayesian logistic regression is commonly based on the Polya-Gamma data augmentation scheme of Polson et al. (2013).
- Through an auxiliary parameter ω , the conditional likelihood of a binary outcome y becomes

$$y'_i | \boldsymbol{X}, \boldsymbol{\beta}, \boldsymbol{\omega} \sim \mathcal{N}(\boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}, \omega_i^{-1}) \text{ for } y'_i := \omega_i^{-1} (y_i - 1/2)$$
 (1)

• Correspondingly, the full conditional distribution of β is given by $\beta \mid \boldsymbol{\omega}, \boldsymbol{\lambda}, \tau, \boldsymbol{y}, \boldsymbol{X} \sim \mathcal{N}(\boldsymbol{\Phi}^{-1}\boldsymbol{X}^{\mathsf{T}}\boldsymbol{\Omega}\boldsymbol{y}', \boldsymbol{\Phi}^{-1})$ for $\boldsymbol{\Phi} = \boldsymbol{X}^{\mathsf{T}}\boldsymbol{\Omega}\boldsymbol{X} + \tau^{-2}\boldsymbol{\Lambda}^{-2}$ (2)

Gibbs sampler for sparse Bayesian logistic regression

- Computation for Bayesian logistic regression is commonly based on the Polya-Gamma data augmentation scheme of Polson et al. (2013).
- Through an auxiliary parameter ω , the conditional likelihood of a binary outcome y becomes

$$y'_i | \boldsymbol{X}, \boldsymbol{\beta}, \boldsymbol{\omega} \sim \mathcal{N}(\boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}, \omega_i^{-1}) \text{ for } y'_i := \omega_i^{-1} (y_i - 1/2)$$
 (1)

- Correspondingly, the full conditional distribution of β is given by $\beta \mid \boldsymbol{\omega}, \boldsymbol{\lambda}, \tau, \boldsymbol{y}, \boldsymbol{X} \sim \mathcal{N}(\boldsymbol{\Phi}^{-1}\boldsymbol{X}^{\mathsf{T}}\boldsymbol{\Omega}\boldsymbol{y}', \boldsymbol{\Phi}^{-1})$ for $\boldsymbol{\Phi} = \boldsymbol{X}^{\mathsf{T}}\boldsymbol{\Omega}\boldsymbol{X} + \tau^{-2}\boldsymbol{\Lambda}^{-2}$ (2)
- Generating samples from (2) is the main computational bottleneck:
 - O(np²) operations for computing X^TΩX
 O(p³) operations for the Cholesky factorization Φ = LL^T
 L^{-T}δ with δ ~ N(0, I_p) ⇒ Var(L^{-T}δ) = L^{-T}L⁻¹ = Φ⁻¹

- Significant progress in $n \ll p$ case:
 - Samples from β | ω, λ, τ, y, X can be generated with only O(n²p) operations (Bhattacharya et al., 2016).
 - Further speed-up via approximation + improved mixing of the global shrinkage parameter through partial marginalization (Johndrow et al., 2018).

Experimental set-up:

- Data: \boldsymbol{X} is of size $72,489 \times 22,175$ and is 95% sparse \boldsymbol{y} indicates the treatment by dabigatran over warfarin $y_i = 1$ account for 27.3% of the 72,489 cases
- Computing environment: iMac 2015 with Intel Core i7 CPUs
- Implementation: reasonably optimized Python code (vectorized & optimal matrix formats used)

Experimental set-up:

- Data: \boldsymbol{X} is of size $72,489 \times 22,175$ and is 95% sparse \boldsymbol{y} indicates the treatment by dabigatran over warfarin $y_i = 1$ account for 27.3% of the 72,489 cases
- Computing environment: iMac 2015 with Intel Core i7 CPUs
- Implementation: reasonably optimized Python code (vectorized & optimal matrix formats used)

Computing time:

- 2 min 30 sec for the matrix-matrix multiplication $X^{\intercal}\Omega X$
- 30 sec for the Cholesky factorization of ${f \Phi} = {m X}^{\intercal} {m \Omega} {m X} + au^{-2} {m \Lambda}^{-2}$

Experimental set-up:

- Data: \boldsymbol{X} is of size $72,489 \times 22,175$ and is 95% sparse \boldsymbol{y} indicates the treatment by dabigatran over warfarin $y_i = 1$ account for 27.3% of the 72,489 cases
- Computing environment: iMac 2015 with Intel Core i7 CPUs
- Implementation: reasonably optimized Python code (vectorized & optimal matrix formats used)

Computing time:

- 2 min 30 sec for the matrix-matrix multiplication $X^{\mathsf{T}}\Omega X$
- 30 sec for the Cholesky factorization of ${f \Phi} = {m X}^{\intercal} {m \Omega} {m X} + au^{-2} {m \Lambda}^{-2}$
- In total, 15 hours for 300 iterations of the Gibbs sampler.

Compare with arguably the most popular sparse regression tool: glmnet

Loaded glmnet 2.0-13 Warning message: package 'glmnet' was built under R version 3.4.2 > runtime < system.time(+ propensity_score_fit <- (cv.glmnet(X, factor(y)) family='binomial', standardize=FALSE) +) > print(runtime) user system elapsed 4938.046 22.916 4947.262

How bad is the standard computational approach?

With the new algorithm, Bayesian method is competitive with glmnet.



① Case for Bayesian approach in search of weak signals

- 2 Global-local shrinkage prior for sparse Bayesian regression
- 3 Conjugate gradient (CG) sampler for multivariate Gaussians
- 4 Preconditioning for rapid CG convergence in sparse regression
- 5 Numerical results / Applications
- 6 Sparse Bayesian regression for survival outcomes
 - 7 Summary

Conjugate gradient sampler for multivariate Gaussians

The problem of generating a sample from (3) can be recast as that of solving a linear system via the algorithm below:

$$\boldsymbol{eta} \sim \mathcal{N}(\boldsymbol{\Phi}^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{y}, \boldsymbol{\Phi}^{-1}) \;\; ext{for} \;\; \boldsymbol{\Phi} = \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{X} + \tau^{-2} \boldsymbol{\Lambda}^{-2} \; (3)$$

Algorithm (Nishimura & Suchard, 2018)

The following procedure generates a sample β from the distribution (3):

The problem of generating a sample from (3) can be recast as that of solving a linear system via the algorithm below:

$$\boldsymbol{\beta} \sim \mathcal{N}(\boldsymbol{\Phi}^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{y}, \boldsymbol{\Phi}^{-1}) \text{ for } \boldsymbol{\Phi} = \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{X} + \tau^{-2} \boldsymbol{\Lambda}^{-2}$$
 (3)

Algorithm (Nishimura & Suchard, 2018)

The following procedure generates a sample β from the distribution (3):

• Generate $\boldsymbol{b} \sim \mathcal{N}(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{y}, \boldsymbol{\Phi})$ by sampling independent Gaussian vectors $\boldsymbol{\eta} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_n)$ and $\boldsymbol{\delta} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_p)$ and then setting

$$\boldsymbol{b} = \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{y} + \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega}^{1/2} \boldsymbol{\eta} + \tau^{-1} \boldsymbol{\Lambda}^{-1} \boldsymbol{\delta}$$
(4)

The problem of generating a sample from (3) can be recast as that of solving a linear system via the algorithm below:

$$\boldsymbol{eta} \sim \mathcal{N}(\boldsymbol{\Phi}^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{y}, \boldsymbol{\Phi}^{-1}) \ \ \text{for} \ \ \boldsymbol{\Phi} = \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{X} + \tau^{-2} \boldsymbol{\Lambda}^{-2}$$
(3)

Algorithm (Nishimura & Suchard, 2018)

The following procedure generates a sample β from the distribution (3):

• Generate $\boldsymbol{b} \sim \mathcal{N}(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{y}, \boldsymbol{\Phi})$ by sampling independent Gaussian vectors $\boldsymbol{\eta} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_n)$ and $\boldsymbol{\delta} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_p)$ and then setting

$$\boldsymbol{b} = \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega} \boldsymbol{y} + \boldsymbol{X}^{\mathsf{T}} \boldsymbol{\Omega}^{1/2} \boldsymbol{\eta} + \tau^{-1} \boldsymbol{\Lambda}^{-1} \boldsymbol{\delta}$$
(4)

2 Solve the following linear system for $oldsymbol{eta}$

$$\Phi eta = b$$
 (5)

so that $\operatorname{Var}(\boldsymbol{\beta}) = \operatorname{Var}(\boldsymbol{\Phi}^{-1}\boldsymbol{b}) = \boldsymbol{\Phi}^{-1}\operatorname{Var}(\boldsymbol{b})\boldsymbol{\Phi}^{-1}.$

- CG is an *iterative method* for solving a positive definite system.
- Given an initial guess β₀, CG generates a sequence {β_k}_{k=1,2,...} of increasingly accurate approximations to the solution.

- CG is an *iterative method* for solving a positive definite system.
- Given an initial guess β₀, CG generates a sequence {β_k}_{k=1,2,...} of increasingly accurate approximations to the solution.
- Cost of update $eta_k o eta_{k+1}$ is dominated by the operation $v o \Phi v.$

- CG is an *iterative method* for solving a positive definite system.
- Given an initial guess β₀, CG generates a sequence {β_k}_{k=1,2,...} of increasingly accurate approximations to the solution.
- Cost of update $eta_k o eta_{k+1}$ is dominated by the operation $m{v} o m{\Phi}m{v}.$
- Note: the operation $v o \Phi v$ requires no explicit formation of Φ
 - Multiplication by $\Phi = X^{\mathsf{T}} \Omega X + \tau^{-2} \Lambda^{-2}$ can be carried out via $v \to Xv$, $w \to X^{\mathsf{T}}w$, and some element-wise multiplications.
 - Memory advantage: X^TΩX are often much denser than X.
 (e.g. 74.5 GB of memory to allocate p × p matrix when p = 10⁵)

- CG is an *iterative method* for solving a positive definite system.
- Given an initial guess β₀, CG generates a sequence {β_k}_{k=1,2,...} of increasingly accurate approximations to the solution.
- Cost of update $eta_k o eta_{k+1}$ is dominated by the operation $v o \Phi v.$
- Note: the operation $v o \Phi v$ requires no explicit formation of Φ
 - Multiplication by $\Phi = X^{\mathsf{T}} \Omega X + \tau^{-2} \Lambda^{-2}$ can be carried out via $v \to Xv, w \to X^{\mathsf{T}}w$, and some element-wise multiplications.
 - Memory advantage: X^TΩX are often much denser than X.
 (e.g. 74.5 GB of memory to allocate p × p matrix when p = 10⁵)
- CG yields an exact solution in p iterations
 - In the worst case, the number of required arithmetic operations is comparable to that a direct method.
 - But it is possible to achieve $\beta_k \approx \beta = \Phi^{-1} b$ for $k \ll p$.

Illustration of the CG Gaussian sampler

• Dabigatran vs. warfarin comparison data (n = 72,489, p = 22,175)

• Compare CG iterates β_k 's to $\beta = \Phi^{-1} b$ for $b \sim \mathcal{N}(X^{\mathsf{T}} \Omega y, \Phi)$.

Illustration of the CG Gaussian sampler

- Dabigatran vs. warfarin comparison data (n = 72,489, p = 22,175)
- Compare CG iterates $oldsymbol{eta}_k$'s to $oldsymbol{eta} = oldsymbol{\Phi}^{-1}oldsymbol{b}$ for $oldsymbol{b} \sim \mathcal{N}ig(oldsymbol{X}^{\mathsf{T}} oldsymbol{\Omega} oldsymbol{y}, oldsymbol{\Phi}ig).$



Figure: Error is quantified as the average of $|(\beta_k)_j - \beta_j|/|\beta_j|$'s.

Case for Bayesian approach in search of weak signals

- 2 Global-local shrinkage prior for sparse Bayesian regression
- Conjugate gradient (CG) sampler for multivariate Gaussians

Preconditioning for rapid CG convergence in sparse regression

- 5 Numerical results / Applications
- 6 Sparse Bayesian regression for survival outcomes

7 Summary

Preconditioning CG to accelerate its convergence

• Direct application of CG rarely leads to a rapid convergence.

Preconditioning CG to accelerate its convergence

- Direct application of CG rarely leads to a rapid convergence.
- CG is typically applied to a preconditioned system

$$ilde{m \Phi} ilde{m eta} = ilde{m b}$$
 for $ilde{m \Phi} = M^{-1/2} m \Phi M^{-1/2}$ and $ilde{m b} = M^{-1/2} m b$ (6)

with a positive definite *preconditioner* matrix M.

- Direct application of CG rarely leads to a rapid convergence.
- CG is typically applied to a preconditioned system

$$ilde{\Phi} ilde{eta}= ilde{m{b}}$$
 for $ilde{m{\Phi}}=M^{-1/2}m{\Phi}M^{-1/2}$ and $ilde{m{b}}=M^{-1/2}m{b}$ (6)

with a positive definite *preconditioner* matrix M.

- Key considerations when designing a preconditioner:
 - Convergence rate of CG applied to the preconditioned system.
 - One-time cost of computing M; better be less than inverting Φ .

- Direct application of CG rarely leads to a rapid convergence.
- CG is typically applied to a preconditioned system

$$ilde{\Phi} ilde{eta}= ilde{m{b}}$$
 for $ilde{m{\Phi}}=M^{-1/2}m{\Phi}M^{-1/2}$ and $ilde{m{b}}=M^{-1/2}m{b}$ (6)

with a positive definite *preconditioner* matrix M.

- Key considerations when designing a preconditioner:
 - Convergence rate of CG applied to the preconditioned system.
 - One-time cost of computing M; better be less than inverting Φ .

"Finding a good preconditioner to solve a given sparse linear system is often viewed as a combination of art and science." Saad (2003)

- Two most famous results on CG convergence rates *not* the most useful ones (Nishimura & Suchard, 2018).
- Following rule of thumb is more useful:
 - \blacktriangleright CG converges quickly if the eigenvalues of Φ are "clustered."
 - if not too many, large eigenvalues cause little delay in convergence.
 - small eigenvalues tend to delay convergence longer.

Illustration of the CG convergence



Figure: Comparison of two preconditioning strategies in sparse Bayesian logistic regression context (n = 72,489 & p = 22,175).

• Normalized error:
$$\left\{\frac{1}{p}\sum_{j}\hat{\xi}_{j}^{-2}(\boldsymbol{\beta}_{k}-\boldsymbol{\beta})_{j}^{2}\right\}^{1/2}$$
 with $\hat{\xi}_{j}^{2} \approx \mathbb{E}[\beta_{j}^{2} \mid \boldsymbol{y}, \boldsymbol{X}]$

• Norm of the preconditioned residual $p^{-1/2} \| ilde{m{r}}_k \|_2$ as a stopping criteria

Illustration of the CG convergence



Figure: Eigenvalue distribution under the proposed preconditioner.

Figure: Eigenvalue distribution of the Jacobi preconditioned matrix.

The magic preconditioner for sparse Bayesian regression problem for $\Phi = X^{\mathsf{T}} \Omega X + \tau^{-2} \Lambda^{-2}$ is ... (drum roll) ...

The magic preconditioner for sparse Bayesian regression problem for $\Phi = X^{\mathsf{T}} \Omega X + \tau^{-2} \Lambda^{-2}$ is ... (drum roll) ...

$$M = \tau^{-2} \Lambda^{-2} \tag{7}$$

The magic preconditioner for sparse Bayesian regression problem for $\Phi = X^{\mathsf{T}} \Omega X + \tau^{-2} \Lambda^{-2}$ is ... (drum roll) ...

$$M = \tau^{-2} \Lambda^{-2} \tag{7}$$

... isn't it ... too simple? (= boring = unpublishable?)

Prior preconditioner: made for sparse Bayesian regression

- Fancier (general-purpose) preconditioners exist: successive over-relaxation, incomplete Cholesky, sparse approximate inverse, etc.
- But they stand no chance against the proposed prior preconditioner

$$\boldsymbol{M} = \tau^{-2} \boldsymbol{\Lambda}^{-2} = \operatorname{Var}(\boldsymbol{\beta} \,|\, \tau, \boldsymbol{\lambda}, \boldsymbol{\omega})^{-1}$$

• Jacobi's $M = \operatorname{diag}(\Phi)$ is reasonable, but is significantly inferior.

- Fancier (general-purpose) preconditioners exist: successive over-relaxation, incomplete Cholesky, sparse approximate inverse, etc.
- But they stand no chance against the proposed prior preconditioner

$$M = \tau^{-2} \Lambda^{-2} = \operatorname{Var}(\beta \,|\, \tau, \lambda, \omega)^{-1}$$

ullet Jacobi's $M=\operatorname{diag}(\Phi)$ is reasonable, but is significantly inferior.

- Mathematically, prior-preconditioning is equivalent to sampling from $\tau^{-1} {\bf \Lambda}^{-1} {\bf \beta} \, | \, \tau, {\bf \lambda}, {\boldsymbol \omega}, {\boldsymbol y}, {\boldsymbol X}$
- Exploits the fundamental feature of sparse Bayesian regression: the prior dominates the likelihood — and hence the posterior looks like the prior — except for a small number of directions

- Fancier (general-purpose) preconditioners exist: successive over-relaxation, incomplete Cholesky, sparse approximate inverse, etc.
- But they stand no chance against the proposed prior preconditioner

$$M = \tau^{-2} \Lambda^{-2} = \operatorname{Var}(\beta \,|\, \tau, \lambda, \omega)^{-1}$$

• Jacobi's $oldsymbol{M}=\operatorname{diag}(oldsymbol{\Phi})$ is reasonable, but is significantly inferior.

- Mathematically, prior-preconditioning is equivalent to sampling from $\tau^{-1} {\bf \Lambda}^{-1} {\bf \beta} \, | \, \tau, {\bf \lambda}, {\boldsymbol \omega}, {\boldsymbol y}, {\boldsymbol X}$
- Exploits the fundamental feature of sparse Bayesian regression: the prior dominates the likelihood — and hence the posterior looks like the prior — except for a small number of directions
- See Nishimura & Suchard (2018) for mathematically precise theory.

Case for Bayesian approach in search of weak signals

- 2 Global-local shrinkage prior for sparse Bayesian regression
- Conjugate gradient (CG) sampler for multivariate Gaussians
- 4 Preconditioning for rapid CG convergence in sparse regression
- 5 Numerical results / Applications
- 6 Sparse Bayesian regression for survival outcomes

7 Summary

- Propensity score estimation
 - ▶ 1,500 iterations in 7.04 (vs 77.4) hours 11 times speed-up
- Doubly robust treatment effect estimation
 - 2,000 iteration for 4.36 (vs 107) hour 25 times speed-up

Propensity scores & difference in population



Propensity scores & difference in population



Effect of dabigatran over warfarin



Case for Bayesian approach in search of weak signals

- 2 Global-local shrinkage prior for sparse Bayesian regression
- Conjugate gradient (CG) sampler for multivariate Gaussians
- Preconditioning for rapid CG convergence in sparse regression
- 5 Numerical results / Applications
- 6 Sparse Bayesian regression for survival outcomes

7 Summary

Sparse Bayesian regression for log-concave likelihood

• Conditional distribution of $\beta | \lambda, \tau, y, X$ typically has no closed-form expression but is log-concave with predictable Hessian structures.

- Conditional distribution of $\beta | \lambda, \tau, y, X$ typically has no closed-form expression but is log-concave with predictable Hessian structures.
- *Hamiltonian Monte Carlo* (HMC) is particularly well-suited to log-concave distributions:
 - computational cost is $O\left(p^{1/4}\left(\frac{M}{m}\right)^{1/2}\right)$ where M and m are the upper and lower bound on eigenvalues of the log-density Hessian (Mangoubi and Smith, 2017).

- Conditional distribution of $\beta | \lambda, \tau, y, X$ typically has no closed-form expression but is log-concave with predictable Hessian structures.
- *Hamiltonian Monte Carlo* (HMC) is particularly well-suited to log-concave distributions:
 - computational cost is $O\left(p^{1/4}\left(\frac{M}{m}\right)^{1/2}\right)$ where M and m are the upper and lower bound on eigenvalues of the log-density Hessian (Mangoubi and Smith, 2017).
- Running HMC in practice requires knowing M and m:
 - efficient computation by Lancoz iteration through a small number of matrix-vector multiplications.

Case for Bayesian approach in search of weak signals

- 2 Global-local shrinkage prior for sparse Bayesian regression
- Conjugate gradient (CG) sampler for multivariate Gaussians
- Preconditioning for rapid CG convergence in sparse regression
- 5 Numerical results / Applications
- 6 Sparse Bayesian regression for survival outcomes



• Prior-preconditioned CG speeds up the conditional updates of β , removing the computational bottleneck of sparse Bayesian regression.

- Prior-preconditioned CG speeds up the conditional updates of β, removing the computational bottleneck of sparse Bayesian regression.
- Iterative methods are promising approaches for accelerating MCMC involving high-dimensional (conditionally) Gaussian distribution because:
 - sparsity in X is natively exploited.
 - approximations to the target covariance make for preconditioners.
 - covariance structures are related from one iteration to another.

- Prior-preconditioned CG speeds up the conditional updates of β, removing the computational bottleneck of sparse Bayesian regression.
- Iterative methods are promising approaches for accelerating MCMC involving high-dimensional (conditionally) Gaussian distribution because:
 - sparsity in X is natively exploited.
 - approximations to the target covariance make for preconditioners.
 - covariance structures are related from one iteration to another.
- Upcoming: sparse Bayesian regression with log-concave likelihoods.

Nishimura, A. and Suchard, M. (2018)

Prior-preconditioned conjugate gradient for accelerated Gibbs sampling in "large n & large p" sparse Bayesian logistic regression models. arXiv:1810.12437

Python package "bayesbridge" available from PyPI

Source code available at

https://github.com/aki-nishimura/bayes-bridge

Thank you!