# Leave Pima Indians alone

Nicolas Chopin
(joint work with James Ridgway)

ENSAE-CREST

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

# Outline

1. Introduction

2. Fast approximations

3. Sampling-based methods

4. Numerical study

5. Variable selection

6. Conclusions

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Binary regression models

Models wih data $y_i \in \{-1, 1\}$, predictors $\mathbf{x}_i \in \mathbb{R}^p$, and likelihood

$$p(\mathcal{D}|\boldsymbol{\beta}) = \prod_{i=1}^{n_{\mathcal{D}}} F(y_i \boldsymbol{\beta}^T \mathbf{x}_i)$$

where $F : \mathbb{R} \to [0, 1]$ is a CDF.
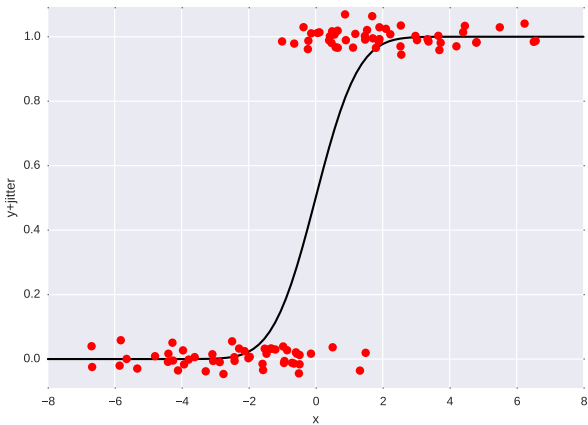
Common examples:

- $F = \Phi$ (probit),
- $F = L$ (logit), where $L(z) = 1/(1 + e^{-z})$.

Introduction
Fast approximations
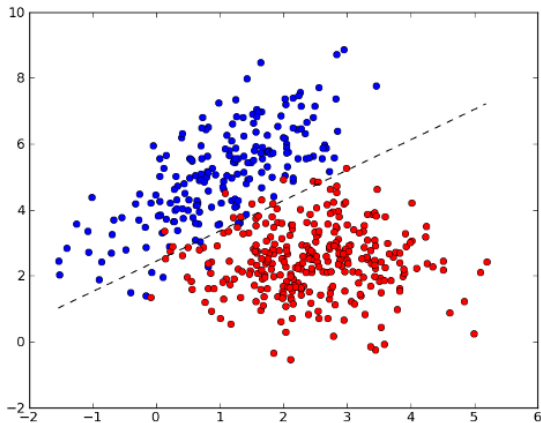Sampling-based methods
Numerical study
Variable selection
Conclusions

# When $p = 1$

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

# Connection with classification

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Properties

- Unless there is **complete separation** in the data, the log-likelihood is concave: MLE is uniquely defined.
- One nice way to deal with complete seperation is to add a proper prior, e.g. Gaussian or Cauchy. (Under Gaussian prior, log-post is concave.)
- Good practice is to standardise the predictors before eliciting the prior (Gelman et al, 2008).

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Binary regression in Bayesian Computation papers

- a long chain of papers on Gibbs sampling for different variants of binary regression models (Albert & Chib, 1993; Holmes & Held, 2006; Fruwirth-Schnatter (2009); Gramacy and Polson, 2012; Polson et al, 2013)

- nearly any paper introducing any new **generic** way to compute a posterior includes a binary regression example:
  - SMC: C (2002), Del Moral et al (2006)
  - HMC and variants: Neal (2010), Shahbaba & Neal (2011), Girolami & Calderhead (2011)
  - NUTS: Hoffman and Gelman (2013)
  - nested sampling: C & Robert (2007)

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Questions

1. Does it make sense to promote binary regression as a **benchmark** for Bayesian computation? (see similar practice in optimisation)

2. In practice, which method one should use???

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Plan

1. review of fast approximation schemes:
   - Laplace (and variants)
   - EP
   - Variational Bayes? (see Consonni & Marin, 2007)

2. review of sampling-based approaches:
   - importance sampling
   - MCMC (Gibbs, RWHM)
   - HMC (and variants)
   - SMC

3. Discussion and comparison

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Considered scenarios

- Model: probit and logit.
- prior: Gaussian and Cauchy (predictors are standardised).

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Laplace

Based on a second order Taylor expansion of the log posterior:

$$\log p(\boldsymbol{\beta}|\mathcal{D}) \approx \log p(\boldsymbol{\beta}_{\mathrm{MAP}}|\mathcal{D}) - \frac{1}{2}\left(\boldsymbol{\beta} - \boldsymbol{\beta}_{\mathrm{MAP}}\right)^{T}\mathbf{Q}\left(\boldsymbol{\beta} - \boldsymbol{\beta}_{\mathrm{MAP}}\right)$$

where $\mathbf{Q}$ is minus the Hessian of $\log p(\boldsymbol{\beta}|\mathcal{D})$ at $\boldsymbol{\beta} = \boldsymbol{\beta}_{\mathrm{MAP}}$.

Exponentiate to get a Gaussian approximation of the posterior. In practice, use Newton-Raphson to obtain $\boldsymbol{\beta}_{\mathrm{MAP}}$ and $\mathbf{Q}$.

Very fast. May not converge if $p$ is very large.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Impoved Laplace

For each marginal:

$$p(\beta_j|\mathcal{D}) \propto \frac{p(\boldsymbol{\beta})p(\mathcal{D}|\boldsymbol{\beta})}{p(\boldsymbol{\beta}_{-j}|\beta_j, \mathcal{D})}$$

Choose a fine grid of $\beta_j$ values; for each $\beta_j$ value, compute a Laplace approximation of $p(\boldsymbol{\beta}_{-j}|\beta_j, \mathcal{D})$.

Note: more expensive, connection with INLA.

Introduction
**Fast approximations**
Sampling-based methods
Numerical study
Variable selection
Conclusions

## EM-Laplace

For a Student prior, Gelman et al (2008) derive an approximate EM scheme based on

$$\beta_j | \sigma_j^2 \sim \mathrm{N}_1(0, \sigma_j^2), \quad \sigma_j^2 \sim \mathrm{Inv} - \mathrm{Gamma}(\nu/2, s_j \nu/2)$$

However, we will observe in our simulations that Laplace still works well for such a prior.

Introduction
**Fast approximations**
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Expectation Propagation

From the following decomposition:

$$p(\boldsymbol{\beta}|\mathcal{D}) = \frac{1}{p(\mathcal{D})} \prod_{i=0}^{n_{\mathcal{D}}} l_i(\boldsymbol{\beta}), \quad l_i(\boldsymbol{\beta}) = F(y_i \boldsymbol{\beta}^T \mathbf{x}_i) \text{ for } i \geq 1,$$

and $l_0$ is prior, EP computes iteratively a parametric approximation of the posterior with the same structure:

$$q_{\mathrm{EP}}(\boldsymbol{\beta}) = \prod_{i=0}^{n_{\mathcal{D}}} \frac{1}{Z_i} q_i(\boldsymbol{\beta}).$$

Taking $q_i$ to be an unnormalised Gaussian density

$$q_i(\boldsymbol{\beta}) = \exp \left\{ -\frac{1}{2} \boldsymbol{\beta}^T \mathbf{Q}_i \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{r}_i \right\},$$

$q_{\mathrm{EP}}$ is a Gaussian with parameters $\mathbf{Q} = \sum_{i=0}^{n} \mathbf{Q}_i$, $\mathbf{r} = \sum_{i=0}^{n} \mathbf{r}_i$

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## EP site update

Update each 'site' in turn: update $q_i$, while keeping $q_j$, $j \neq i$ fixed, by minimising the Kullback-Leibler divergence between

$$h(\boldsymbol{\beta}) \propto l_i(\boldsymbol{\beta}) \prod_{j \neq i} q_j(\boldsymbol{\beta})$$

and $q(\boldsymbol{\beta}) \propto \prod_j q_j$.

Thanks to nice properties of exponential families, this boils to match the moments of $h$ and $q$.

In binary regression, these site updates lead to explicit expressions (probit) or one-dimensional integrals that are easy to approximate accurately (logit).

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## General remarks

- Since the approximation methods covered in the previous section are faster by orders of magnitude than sampling-based methods, we will assume that a Gaussian approximation $q(\beta)$ (from Laplace or EP) has been computed in a preliminary step.
- Complexity: Laplace is $O(n_{\mathcal{D}} + p^3)$, EP is $O(n_{\mathcal{D}}p^3)$.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Importance sampling

Proposal $q$ set to some Gaussian approx of the posterior. Then to approximate $p(\mathcal{D})$, generate $\beta_1, \ldots, \beta_N \sim q$, compute

$$Z_N = \frac{1}{N} \sum_{n=1}^{N} w(\beta_n), \quad w(\beta) := \frac{p(\beta)p(\mathcal{D}|\beta)}{q(\beta)}$$

and to approximate the posterior expectation of $\varphi$, compute

$$\varphi_N = \frac{\sum_{n=1}^{N} w(\beta_n)\varphi(\beta_n)}{\sum_{n=1}^{N} w(\beta_n)}.$$

Introduction
Fast approximations
**Sampling-based methods**
Numerical study
Variable selection
Conclusions

## IS pros and cons

Pros:

- simple, generic
- embarassingly parallel
- approximates the marginal likelihood at no extra cost
- IID sampling: MC error is easy to assess
- can plug in QMC points

Cons:

- ESS may collapse when $p$ is large.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## MCMC general remarks

The following points

- choice of starting point
- MCMC convergence assessment

are not big issues for binary regression models.

More important issues for us are:

- chain autocorrelations
- difficulty to parallelise

Introduction
Fast approximations
**Sampling-based methods**
Numerical study
Variable selection
Conclusions

## Gibbs

Well-known, based on data augmentation:

$$z_i = \boldsymbol{\beta}^T \mathbf{x}_i + \epsilon_i$$
$$y_i = \text{sgn}(z_i)$$

then sample iteratively (probit/Gaussian case):

1. $\boldsymbol{\beta} | \mathbf{z}$ (regression posterior, tractable)
2. $\mathbf{z} | \boldsymbol{\beta}, \mathbf{y}$ (product of truncated Gaussians)

Gibbs is particularly **not generic**: any change in the prior of $F$ requires deriving a new algorithm. This can also change the complexity (e.g. from $\mathcal{O}(p^2)$ to $\mathcal{O}(p^3)$ when using a Student prior).

Introduction
Fast approximations
**Sampling-based methods**
Numerical study
Variable selection
Conclusions

# Random walk Metropolis-Hastings

### One iteration of RWMH

Input: $\boldsymbol{\beta}$
Output: $\boldsymbol{\beta}'$
1. Sample $\boldsymbol{\beta}^\star \sim \mathrm{N}_p(\boldsymbol{\beta}, \boldsymbol{\Sigma})$
2. With probability $1 \wedge r$,

$$r = \frac{p(\boldsymbol{\beta}^\star)p(\mathcal{D}|\boldsymbol{\beta}^\star)}{p(\boldsymbol{\beta})p(\mathcal{D}|\boldsymbol{\beta})},$$

set $\boldsymbol{\beta}' = \boldsymbol{\beta}^\star$; otherwise set $\boldsymbol{\beta}' = \boldsymbol{\beta}$

In practice, choose $\boldsymbol{\Sigma}$ as some fraction of $\boldsymbol{\Sigma}_q$.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## HMC

Consider $(\beta, \alpha)$, $\beta \sim p(\beta|\mathcal{D})$, $\alpha \sim N_p(0, M^{-1})$, with joint un-normalised density $\exp\{-H(\beta, \alpha)\}$,

$$H(\beta, \alpha) = E(\beta) + \frac{1}{2}\alpha^T \mathbf{M}\alpha, \quad E(\beta) = -\log\{p(\beta)p(\mathcal{D}|\beta)\}.$$

The physical interpretation of HMC is that of a particle at position $\beta$, with velocity $\alpha$, potential energy $E(\beta)$, kinetic energy $\frac{1}{2}\alpha^T M\alpha$, and thus total energy given by $H(\beta, \alpha)$. The particle is expected to follow a trajectory such that $H(\beta, \alpha)$ remains constant over time.

Introduction
Fast approximations
**Sampling-based methods**
Numerical study
Variable selection
Conclusions

## HMC iteration

### One iteration of HMC

Input: $\boldsymbol{\beta}$
Output: $\boldsymbol{\beta}'$
1. Sample momentum $\boldsymbol{\alpha} \sim \mathrm{N}_p(0, \mathbf{M})$.
2. Perform $L$ leap-frog steps, starting from $(\boldsymbol{\beta}, \boldsymbol{\alpha})$; call $(\boldsymbol{\beta}^\star, \boldsymbol{\alpha}^\star)$ the final position.
3. With probability $1 \wedge r$, $r = \exp\{H(\boldsymbol{\beta}, \boldsymbol{\alpha}) - H(\boldsymbol{\beta}^\star, \boldsymbol{\alpha}^\star)\}$ set $\boldsymbol{\beta}' = \boldsymbol{\beta}^\star$; otherwise set $\boldsymbol{\beta}' = \boldsymbol{\beta}$.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

# Leapfrog step

### Leapfrog step

Input: $(\boldsymbol{\beta}, \boldsymbol{\alpha})$
Output: $(\boldsymbol{\beta}_1, \boldsymbol{\alpha}_1)$
1. $\boldsymbol{\alpha}_{1/2} \leftarrow \boldsymbol{\alpha} - \frac{\epsilon}{2} \nabla_\beta E(\boldsymbol{\beta})$
2. $\boldsymbol{\beta}_1 \leftarrow \boldsymbol{\beta} + \epsilon \boldsymbol{\alpha}_{1/2}$
3. $\boldsymbol{\alpha}_1 \leftarrow \boldsymbol{\alpha}_{1/2} - \frac{\epsilon}{2} \nabla_\beta E(\boldsymbol{\beta}_1)$

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## HMC variants

- Riemanian HMC (Girolami and Calderhead, 2011): simply too expensive
- NUTS (No U-Turn Sampler, Hoffman & Gelman, 2013): HMC with on-the-fly calibration of $L$ and $\epsilon$. Included in our comparisons.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## SMC

We consider tempering SMC, i.e. SMC for sequence

$$\pi_t(\beta) \propto q(\beta)^{1-\delta_t} \left\{ p(\beta)p(\mathcal{D}|\beta) \right\}^{\delta_t}$$

**Principle**: sequence of importance sampling steps, from $\pi_{t-1}$ to $\pi_t$. When weight degeneracy becomes too high, resample, and move particles through MCMC (e.g. random walk Metropolis).

The algorithm can choose the $\delta_j$ on the fly (Jasra et al, 2011).

Introduction
Fast approximations
**Sampling-based methods**
Numerical study
Variable selection
Conclusions

## SMC algorithm

Operations involving index $n$ must be performed for all $n \in 1:N$.

**0** Sample $\beta_n \sim q(\beta)$ and set $\underline{\delta} \leftarrow 0$.

**1** Let, for $\delta \in [\underline{\delta}, 1]$,

$$\text{EF}(\delta) = \frac{1}{N} \frac{\left\{ \sum_{n=1}^{N} w_\gamma(\beta_n) \right\}^2}{\left\{ \sum_{n=1}^{N} w_\gamma(\beta_n)^2 \right\}}, \quad u_\delta(\beta) = \left\{ \frac{p(\beta)p(\mathcal{D}|\beta)}{q(\beta)} \right\}^\delta.$$

If $\text{EF}(1) \geq \tau$, stop and return $(\beta_n, w_n)_{n=1:N}$ with $w_n = u_1(\beta_n)$; otherwise, use the bisection method to solve numerically in $\delta$ the equation $\text{EF}(\gamma) = \tau$.

**2** Resample according to normalised weights
$W_n = w_n / \sum_{m=1}^{N} w_m$; with $w_n = u_\delta(\beta_n)$.

**3** Update the $\beta_n$'s through $m$ MCMC steps that leaves invariant

Introduction
Fast approximations
**Sampling-based methods**
Numerical study
Variable selection
Conclusions

## Remarks on SMC

- Completely automatic: we can use the current set of particles to adjust the random walk proposal, the number of MCMC steps, and so on.
- Will often collapse to a **single** IS step (when ESS from $q$ to posterior is not too low)

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## First set of datasets

| Dataset | $n_{\mathcal{D}}$ | $p$ |
|---|---|---|
| Pima (Indian diabetes) | 532 | 8 |
| German (credit) | 999 | 25 |
| Heart (Statlog) | 270 | 14 |
| Breast (cancer) | 683 | 10 |
| Liver (Indian Liver patient) | 579 | 11 |
| Plasma (blood screening data) | 32 | 3 |
| Australian (credit) | 690 | 15 |
| Elections | 2015 | 52 |

This is a superset of datasets considered in most papers.

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Fast approximations

Logit/Cauchy scenario. We compare: Laplace, Improved Laplace, EM-Laplace, and EP, in term of

- marginal accuracies (one minus half the $L_1$ distance between approximate and true marginals)
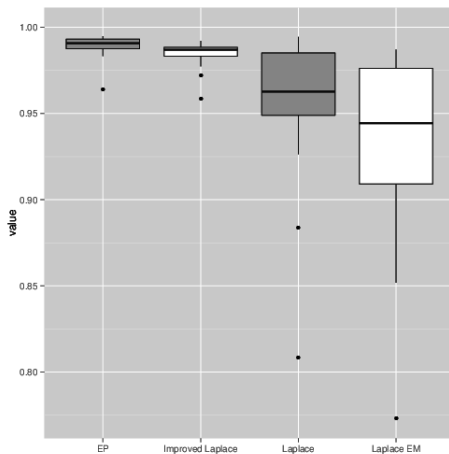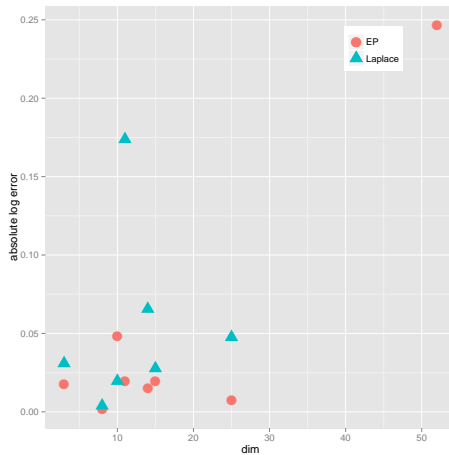- approximation error for marginal likelihood

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Pima

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

# Heart

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Breast

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## German credit

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

# Marginal likelihoods

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Sampling-based methods: importance sampling

| | IS | | | IS-QMC | |
|---|---|---|---|---|---|
| Dataset | EF | CPU | MT | MSE x | MSE x |
| | $= \mathrm{ESS}/N$ | time | speed-up | (expect) | (evid) |
| Pima | 99.5% | 37.54 s | 4.39 | 28.9 | 42.7 |
| German | 97.9% | 79.65 s | 4.51 | 13.2 | 8.2 |
| Breast | 82.9% | 50.91 s | 4.45 | 2.6 | 6.2 |
| Heart | 95.2% | 22.34 s | 4.53 | 8.8 | 9.3 |
| Liver | 74.2 % | 35.93 s | 4.76 | 7.6 | 11.3 |
| Plasma | 90.0% | 2.32 s | 4.28 | 2.2 | 4.4 |
| Australian | 95.6% | 53.32 s | 4.57 | 12 | 20.3 |
| Elections | 21.39% | 139.48 s | 3.87 | 617.9 | 3.53 |

(Probit/Gaussian scenario, to make like easier for Gibbs)

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## comparison with MCMC



IRIS = Inefficiency relative to IS

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Bigger datasets

| Dataset | $n_{\mathcal{D}}$ | $p$ |
|---------|------|-----|
| Musk | 476 | 95 |
| Sonar | 208 | 61 |
| DNA | 400 | 180 |

Bigger datasets, but also with higher correlations between predictors. We will look at the probit/Gaussian case.

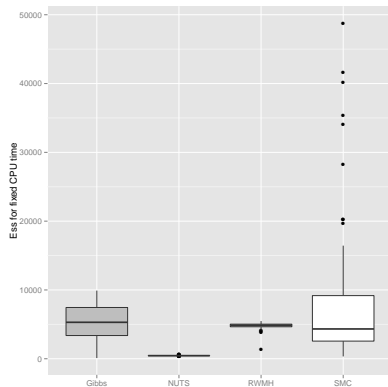IS no longer an option.

Introduction
Fast approximations
Sampling-based methods
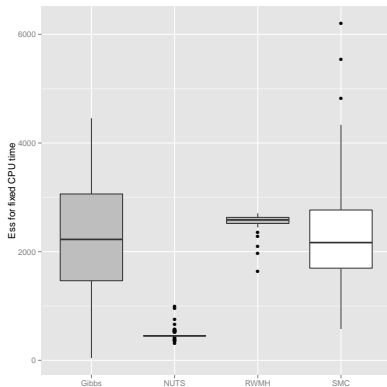**Numerical study**
Variable selection
Conclusions

## Approximations: Musk

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Approximations: Sonar

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Approximations: DNA

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Sampling-based methods: Musk



Left: posterior expectations, Right: posterior variances

Introduction
Fast approximations
Sampling-based methods
**Numerical study**
Variable selection
Conclusions

## Sampling-based methods: Sonar



Left: posterior expectations, Right: posterior variances

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Sampling-based methods: DNA



Left: posterior expectations, Right: posterior variances

Introduction
Fast approximations
Sampling-based methods
Numerical study
**Variable selection**
Conclusions

## Variable selection

Add for each predictor $\beta_j$ an indicator $\gamma_j \in \{0, 1\}$; prior for $\gamma$ is Uniform over $\{0, 1\}^p$.

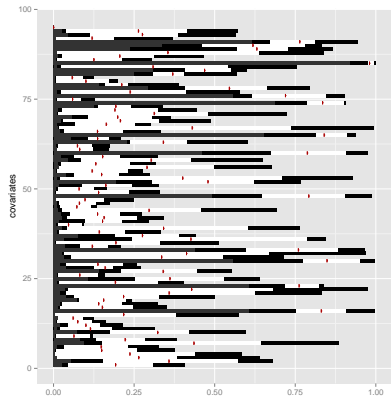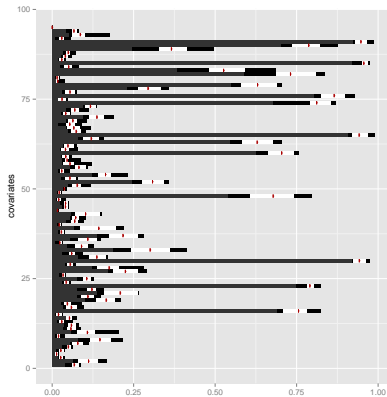The posterior mixes discrete and continuous components; $p(\gamma|\mathcal{D})$ is severely multimodal.

Introduction
Fast approximations
Sampling-based methods
Numerical study
**Variable selection**
Conclusions

## VS: proposed approach

To compute $p(\mathcal{D}|\gamma) = \int p(\mathcal{D}|\gamma, \beta)p(\beta|\gamma) \, d\beta$, use:

1. either Laplace
2. or IS based on Laplace

To simulate from $p(\gamma|\mathcal{D})$, adapt the tempering SMC sampler of Schafer and Chopin (2013), for sampling binary vectors.

Introduction
Fast approximations
Sampling-based methods
Numerical study
**Variable selection**
Conclusions

## Results

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
**Conclusions**

# Recommendations to end users (who wish to fit a binary regression model)

- EP is fast and accurate even in difficult cases.
- to improve on EP, one might run SMC; often this will collapse to IS and outperforms everything else significantly.
- That said, for large *p*, RWHM performs surprising well.
- HMC algorithms seem very difficult to calibrate.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
Conclusions

## Benchmarks for specialised algorithms

For specialised algorithms (Gibbs), benchmark=dataset.

It is not very clear that the Gibbs samplers developed for binary regression are very useful: corresponding papers tend to showcase these algorithms on datasets with $p < 50$, for which more generic methods fare much better.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
**Conclusions**

## Benchmarks for generic algorithms

For generic algorithms (e.g. RWHM), benchmark=posterior.

A binary regression posterior of dimension $< 50$ is very close to a Gaussian; i.e. it does not represent a very challenging benchmark. However, it is an useful **sanity check**.

More challenging benchmarks: $p \geq 100$, hierarchical regression, spike and slab prior, ...

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
**Conclusions**

## More general remarks

Beware ML fast approximation schemes; they are fast and getting better and better. . .

Always compare new methods to well calibrated simple algorithms, like IS and RWHM.

Introduction
Fast approximations
Sampling-based methods
Numerical study
Variable selection
**Conclusions**

## Final word

Comments most welcome!