

Approximate methods

Lecturer: Riccardo Corradin

Being approximate: why

- MCMC methods could be expensive to compute, especially for large sample sizes n.
- The computational cost increases mainly by two components: evaluating the (log)likelihood function, usually when we need its normalization constant and is hard to compute, and proposing a new candidate.
- Moreover, many mcmc algorithms require a rough estimate of some key posterior quantities, such as the posterior variance. Recall, e.g., the MALA.
- These issues motivate the development of deterministic approximations of the posterior distribution or the definition of approximate sampling schemes, that avoid the usual bottlenecks.
- Compared to MCMC methods, the accuracy of this class of approximations can not be reduced by running the algorithm longer.
- On the other hand, deterministic approximations are typically very fast to compute and sufficiently reliable in several applied contexts, while approximate sampling schemes can be used in challenging and intractable scenarios.

Laplace Approximation

The Laplace approximation

- Let $\pi(\theta \mid X)$ be a continuous and differentiable posterior density in $\Theta \subseteq \mathbb{R}^{p}$.
- The Laplace approximation is one of the first approximation methods that has been proposed. It was known even before the advent of MCMC.
- The key idea is approximating the log-posterior density log π(θ | X) using a Taylor expansion around the mode θ̂_{MAP}, yielding

$$\log \pi(\boldsymbol{\theta} \mid \boldsymbol{X}) \approx \log \pi(\hat{\boldsymbol{\theta}}_{MAP} \mid \boldsymbol{X}) - \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_{MAP})^{\mathsf{T}} \hat{M}^{-1} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_{MAP}) + \text{const},$$

where \hat{M} is the negative Hessian of log $\pi(\theta \mid X)$ evaluated at $\hat{\theta_{MAP}}$, that is

$$\hat{M} = -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\mathsf{T}}} \log \pi(\boldsymbol{\theta} \mid \boldsymbol{X}) \Big|_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_{MAP}}$$

 Hence, the above quadratic expansion leads to the following multivariate Gaussian approximate posterior

$$\pi(\boldsymbol{\theta} \mid \boldsymbol{X}) \approx N_{\rho}(\hat{\boldsymbol{\theta}}_{MAP}, \hat{M}^{-1}).$$

Comments on the Laplace approximation

- A fairly strong asymptotic justification of the Laplace approximation is based on the Bernstein-von Mises theorem.
- Suppose the data X_1, \ldots, X_n are iid from a "true" model P_{θ_0} .
- Very roughly speaking, under suitable regularity and sampling conditions

$$||\pi(\boldsymbol{\theta} \mid \boldsymbol{X}) - N_{\mathcal{P}}(\hat{\boldsymbol{\theta}}_{MAP}, \hat{M}^{-1})||_{TV} \stackrel{P_{\boldsymbol{\theta}_{0}}}{\longrightarrow} 0, \qquad n \to \infty,$$

meaning that the total variation distance between the posterior and the Laplace approximation weakly converges to 0 w.r.t. to the law of the sampling process P_{θ_0} .

- Here we are also assuming that $\hat{\theta}_{MAP}$ and $n\hat{M}^{-1}$ are consistent estimators for the "true" parameter value θ_0 and for the inverse Fisher information matrix, respectively.
- Hence, in several cases and for *n* large enough, the law $\pi(\theta \mid \mathbf{X})$ is roughly a Gaussian centered at the mode and with variance depending on the Fisher information.

Comments on the Laplace approximation

- The Laplace approximation is an old and simple method with appealing asymptotic guarantees. Moreover, it only requires the computation of the Hessian and the MAP.
- Refined higher order improvements of expected posterior functionals can be obtained as in Tierney and Kadane (1987).
- On the other hand, especially when the sample size *n* is relatively small, the quadratic approximation of $\log \pi(\theta \mid \mathbf{X})$ may perform poorly.
- For example, if the posterior is not symmetric and unimodal, the map is not a good estimate for the posterior mean, thus leading to inaccurate Gaussian approximations.
- Furthermore, if the parameter space Θ is bounded, a Gaussian approximation could be quite problematic ⇒ a reparametrization should be considered.
- Finally, it is unclear how to handle discrete parameter spaces.

More general approximations

- Let π(θ | X) ∈ P_θ be a posterior distribution (intractable), and let q(θ) ∈ Q be a density, where Q_θ ⊆ P_θ denotes the space tractable densities with support Θ.
- In a general fashion, an optimal approximation $\hat{q}(\theta) \in \mathbb{Q}$ of the posterior distribution is defined as

$$\hat{q}(oldsymbol{ heta}) = rgmin_{q \in \mathbb{Q}} \mathcal{D}(q(oldsymbol{ heta}), \pi(oldsymbol{ heta} \mid oldsymbol{X}))$$

where $\mathcal{D}(\cdot, \cdot)$ is some divergence or metric over the space of probability distribution.

- Among other, one choice of \mathcal{D} commonly used is the Kullback-Leibler divergence $\mathcal{D}(\cdot, \cdot) = \mathrm{KL}(\cdot \mid\mid \cdot)$, resulting in well-known approximation approaches.
- Depending on the choice of the divergence D(·, ·) and of the space of approximating densities Q, the problem could be computationally feasible or not.
- Clearly, this approach is relevant whenever the posterior π(θ | X) is not included in the subspace of tractable density Q.

More general approximations

- As for the choice of D(·, ·), it would be theoretically appealing to consider metrics such as the Hellinger distance, the total variation distance, or the Wasserstein distance.
- Unfortunately, even when we let Q be the space of multivariate Gaussians, finding the optimal approximate density $\hat{q}(\theta)$ could be problematic.
- A basic requirement is that the optimization procedure should **not depend** on the intractable **normalizing constant** of the posterior.
- We will consider two different though quite related divergences.
- The $KL(q(\theta) || \pi(\theta | X))$ divergence, leading to the variational Bayes method.
- The $KL(\pi(\theta \mid X) \mid\mid q(\theta))$ divergence, leading to the expectation propagation method.

The evidence lower bound (ELBO)

In the first place, we recall that

$$extsf{KL}(q(oldsymbol{ heta}) \mid \pi(oldsymbol{ heta} \mid oldsymbol{X})) = -\int_{\Theta} q(oldsymbol{ heta}) \log rac{\pi(oldsymbol{ heta} \mid oldsymbol{X})}{q(oldsymbol{ heta})} \mathrm{d}oldsymbol{ heta}.$$

• Hence, by multiplying and dividing the fractional term by $\pi(\mathbf{X})$, we obtain

$$egin{aligned} & extsf{KL}(q(heta) \mid\mid \pi(heta \mid extsf{X})) = -\int_{\Theta} q(heta) \log rac{\pi(extsf{X})\pi(heta \mid extsf{X})}{\pi(extsf{X})q(heta)} \mathrm{d} heta \ &= -\int_{\Theta} q(heta) \log rac{\pi(heta, extsf{X})}{q(heta)} \mathrm{d} heta + \log \pi(extsf{X}) \end{aligned}$$

where the first term in the last expression is called evidence lower bound $ELBO(q(\theta))$, and the last term does not depend on θ .

• Hence, if we want the best approximating density function $q(\theta) \in \mathbb{Q}_{\theta}$, we have

$$\hat{q}(\boldsymbol{\theta}) = \operatorname*{arg\,min}_{q \in \mathbb{Q}} \mathsf{KL}(q(\boldsymbol{\theta}) \mid\mid \pi(\boldsymbol{\theta} \mid \boldsymbol{X})) = \operatorname*{arg\,max}_{q \in \mathbb{Q}} \mathsf{ELBO}(q(\boldsymbol{\theta})),$$

and the optimization does not depend on the intractable normalizing constant.

The evidence lower bound (ELBO)

• The ELBO is indeed a lower bound of the marginal likelihood, because the divergence $KL(q(\theta) || \pi(\theta | \mathbf{X})) \ge 0$, implying that

 $ELBO(q(\theta)) \leq \log \pi(X).$

- This property of the ELBO has led to using the variational bound as a **model selection criterion**, assuming that the ELBO is a good approximation of the marginal.
- Even when the optimal distribution $\hat{q}(\theta)$ can be found, there is no guarantee that the minimized KL

 $KL(q(\boldsymbol{\theta}) \mid\mid \pi(\boldsymbol{\theta} \mid \boldsymbol{X})) \geq 0$

will be small in absolute terms.

- Moreover, quantifying the value of $KL(\hat{q}(\theta) || \pi(\theta | \mathbf{X})) = \log \pi(\mathbf{X}) ELBOq(\theta)$ would require the knowledge of the normalizing constant, which is intractable.
- Essentially, it is currently hard to assess the quality of the obtained approximation without comparing it with some "gold standard" such as MCMC.

The evidence lower bound (ELBO)

- The VB optimization problem is ill-posed if we do not specify a tractable class \mathbb{Q} .
- For reasons that will become clear later on, a convenient assumption is restricting the focus on the class $\mathbb Q$ of mean-field approximations, in which we assume

$$q(oldsymbol{ heta}) = \prod_{b=1}^B q(oldsymbol{ heta}_b),$$

implying that we are forcing independence among B groups of parameters.

- It is important to notice that dependence is preserved within each block of parameters.
- Moreover, note that we are not forcing $q(\theta)$ to belong to any known parametric family of distributions. The only assumption we are making is independence.

Derivation of the CAVI algorithm

- Under the mean-field approximation, we can derive the so-called coordinate ascent variational inference (CAVI) algorithm.
- The optimization of the ELBO can be written as optimization of

$$\mathsf{ELBO}(q(oldsymbol{ heta})) = \int_{\Theta} \prod_{b=1}^B \left[q(oldsymbol{ heta}_b) \pi(oldsymbol{ heta},oldsymbol{X})
ight] \mathrm{d}oldsymbol{ heta} - \int_{\Theta} \prod_{b=1}^B \left[q(oldsymbol{ heta}_b) \log q(oldsymbol{ heta}_b)
ight] \mathrm{d}oldsymbol{ heta}$$

• We aim at maximizing the bth component $q(\theta_b)$, keeping the others fixed. Thus, we express the elbo isolating the term $q(\theta_b)$, obtaining

$$\int q(\boldsymbol{\theta}_b) \left[\int \log \pi(\boldsymbol{\theta} \mid \boldsymbol{X}) \prod_{j \neq b} q(\boldsymbol{\theta}_j) \mathrm{d}\boldsymbol{\theta}_{-b} \right] \mathrm{d}\boldsymbol{\theta}_b - \int q(\boldsymbol{\theta}_b) \log q(\boldsymbol{\theta}_b) \mathrm{d}\boldsymbol{\theta}_b + c_b,$$

where c_b denotes a term not depending on θ_b .

Defining the log-density π̃(θ_b, X) = E_{-b}[π(θ, X)] + const and re-arranging the terms, we get

$$extsf{ELBO}(q(m{ heta})) = \int q(m{ heta}_b) \log rac{ ilde{\pi}(m{ heta}) m{ extsf{X}}}{q(m{ heta}_b)} \mathrm{d} m{ heta}_b + ilde{c}_b = - extsf{KL}(q(m{ heta}_b) \mid\mid ilde{\pi}(m{ heta}_b, m{ extsf{X}})) + ilde{c}_b.$$

Properties and convergence

• The above previous chain of equations implies that the local maximization of the $ELBO(q(\theta))$ with respect to the bth term of $q(\theta_b)$ is obtained by setting

$$\hat{q}(\boldsymbol{\theta}_b) \propto \exp\Big\{ \mathbb{E}_{-b}[\log \pi(\boldsymbol{\theta}, \boldsymbol{X})] \Big\},$$

for any $b = 1, \ldots, B$.

- In practice, the above expectation is often straightforward to compute, and some known kernel can usually be recognized (as in the Gibbs sampling).
- In the CAVI algorithm, we iteratively update the factors $q(\theta_b)$ by using the locally maximized terms given the others.
- By construction, this produces a monotonic sequence that convergences to a local optimum of the ELBO.
- The CAVI is an appealing algorithm for maximizing the ELBO under the mean-field assumption, but in principle, one could use any other optimizer.
- The necessary computations and expectations are usually doable if the full conditional distributions belong to some exponential family.
- The algorithm stops whenever the ELBO sequence has converged.

Underestimation of the variability

- The combination of mean-field assumption + VB approach typically leads to a sensible underestimation of the variability.
- In the first place, this is a consequence of the insufficient flexibility of the mean-field class of approximating densities.
- Indeed, if the densities in Q were arbitrarily close to the posterior, this phenomenon would be negligible.
- In second place, this is a consequence of the chosen divergence. Indeed, the quantity

$$extsf{KL}(q(oldsymbol{ heta}) \parallel \pi(oldsymbol{ heta} \mid oldsymbol{X})) = -\int_{\Theta} q(oldsymbol{ heta}) \log rac{\pi(oldsymbol{ heta} \mid oldsymbol{X})}{q(oldsymbol{ heta})} \mathrm{d}oldsymbol{ heta}$$

favors the choice of densities $q(\theta)$ which are included in the support of $\pi(\theta \mid \mathbf{X})$.

• Further, there is a large positive contribution to the above KL for those values of θ such that $\pi(\theta \mid \mathbf{X}) \approx 0$, unless $q(\theta) \approx 0$ as well.

The cavi for a Gaussian example

• Let us assume the observations (x_1, \ldots, x_n) are draws from

$$X_i \mid \mu, \tau \sim N(\mu, \tau^{-1}), \qquad i = 1, \dots, n,$$

with independent priors $\mu \sim N(m_0, s_0^2)$ and $\tau \sim \textit{Ga}(a_0, b_0)$.

- Assuming a mean-field approximation q(μ, τ) = q(μ)q(τ), the cavi algorithm iterates between the following steps simple steps.
- Update $q(\mu)$. The locally optimal variational distribution for $q(\mu)$ is

$$q(\mu) \stackrel{d}{=} N(m_n, s_n^2), \qquad m_n = s_n^2 \left(\frac{m_0}{s_0^2} + \operatorname{E}_q[\tau] \sum_{i=1}^n x_i \right), \qquad s_n^2 = \left(n \operatorname{E}_q[\tau] + \frac{1}{s_0^2} \right)^{-1}.$$

• Update $q(\tau)$. The locally optimal variational distribution for $q(\tau)$ is

$$q(\tau) \stackrel{d}{=} Ga(a_n, b_n), \qquad a_n = a_0 + \frac{n}{2}, \qquad b_n = b_0 + \frac{1}{2} \sum_{i=1}^n \mathrm{E}_q \Big[(x_i - \mu)^2 \Big].$$

The Pima dataset

- The logistic regression case has often been presented as an example in which mean-field variational Bayes can not be applied. See, for example, Section 10.5 of Bishop (2006).
- A priori, we set $\beta \sim N_{\rho}(\boldsymbol{b}_0, B_0)$.
- Recently, a mean-field approximation for an augmented version of the model has been proposed (Durante, D. and Rigon, T., 2019).
- We consider a sequence of augmented variables $Z_1, \ldots, Z_n \sim PG(1,0)$, where $PG(\alpha, \gamma)$ denotes the Pólya-gamma random variable, i.e.

$$Z \stackrel{d}{=} rac{1}{2\pi^2} \sum_{j \geq 1} rac{{{{{\it G}_j}}}}{{(j - 1/2)^2} + {\gamma^2}/{(4\pi^2)}},$$

with $G_j \sim Ga(\alpha, 1)$, $\alpha > 0$ and $\gamma \in \mathbb{R}$.

• The density of a Pólya-gamma random variable is expressed in terms of an infinite summation, but it can be easily simulated.

The Pima dataset

- The algorithm iterates between two simple steps.
- Update $q(\beta)$. The locally optimal variational distribution for $q(\beta)$ is

$$egin{split} q(oldsymbol{eta}) \propto \exp\left[\mathrm{E}_q \{ \log \pi(oldsymbol{y},oldsymbol{z} \mid eta) + \log \pi(oldsymbol{eta}) \}
ight] \ &\propto \pi(oldsymbol{eta}) \left\{ \sum_{i=1}^n (y_i - 1/2) oldsymbol{x}_i^{\mathsf{T}} oldsymbol{eta} - rac{1}{2} \mathrm{E}_q[z_i] (oldsymbol{x}_i^{\mathsf{T}} oldsymbol{eta})^2
ight\} \end{split}$$

Re-arranging the above equation, we obtain that $q(\beta) \stackrel{d}{=} N_p(\mu, \Sigma)$ with

$$\boldsymbol{\mu} = \boldsymbol{\Sigma}(\boldsymbol{X}^{\mathsf{T}}(\boldsymbol{y} - 1/2) + B_0^{-1}\boldsymbol{b}_0), \qquad \boldsymbol{\Sigma} = (\boldsymbol{X}^{\mathsf{T}} \mathbf{E}_q[\boldsymbol{Z}] \boldsymbol{X} + B^{-1})^{-1},$$

where $Z = \text{diag}(z_1, \ldots, z_n)$ and its expectation is taken w.r.t. q(z).

 Hence, the optimal variational distribution for β is Gaussian. This is an implication of the mean-field structure and not an assumption.

The Pima dataset

- The second step involves the variational distribution q(z).
- Update q(z). The locally optimal variational distribution for q(z) is

$$\begin{aligned} q(\boldsymbol{z}) &\propto \exp \mathrm{E}_q[\log p(\boldsymbol{y}, \boldsymbol{z} \mid \boldsymbol{\beta})] \\ &\propto \prod_{i=1}^n p(z_i \mid 1, 0) \exp\left(-\frac{z_i}{2} \mathrm{E}_q[\eta_i^2]\right) \end{aligned}$$

Re-arranging the above equation, we have

$$q(\mathbf{z}) = \prod_{i=1}^{n} PG(1, \mathbb{E}_{q}[\eta_{i}^{2}])$$

 Hence, the optimal variational distribution for z are independent Pólya-gamma distributions. As before, this is an implication and not an assumption.

- We consider again the Pima dataset, using the CAVI algorithm.
- We use a $N_p(\mathbf{0}, 100I_p)$ as prior distribution.

```
# prior parameters
b <- rep(0, 8)
B <- diag(100, 8)
# Running the CAVI
fit_CAVI <- logit_CAVI(y, X, B, b)
est_coef <- fit_CAVI$mu
# X Xnpreg Xglu Xbp Xskin Xbmi Xped Xage
# -0.990 0.406 1.097 -0.095 0.072 0.569 0.452 0.284</pre>
```

Approximate Bayesian computation

Even more approximate

- In the previous case, we introduced approximate methods to avoid the intractability of $\pi(\theta \mid \mathbf{X})$, with $X_i \in \mathbb{X}$ and $\theta \in \Theta \subseteq \mathbb{R}^p$, but we assumed that we can evaluate/deal with $\pi(\theta, \mathbf{X})$.
- Recall that $\pi(\theta, \mathbf{X}) = L(\mathbf{X} \mid \theta)\pi(\theta)$, where the likelihood term is given by

$$L(\boldsymbol{X} \mid \boldsymbol{\theta}) = \prod_{i=1}^{n} f(x_i \mid \boldsymbol{\theta}).$$

• In many scenarios, even evaluating the likelihood is unfeasible, since we have

$$f(x_i \mid \boldsymbol{\theta}) = \frac{1}{Z_{\boldsymbol{\theta}}}g(x_i \mid \boldsymbol{\theta}),$$

where $Z_{\theta} < +\infty$ is an intractable normalizing constant, which depends on the value of θ .

 At the same time, even if we cannot evaluate the density function f(x_i | θ) of X_i, in many scenarios we can easily sample X_i.

Some examples

 g-and-k distribution, extends the Gaussian distribution with added skewness and heavier/lighter tails, defined by the quantile function

$$F_{gk}^{-1}(u) = a + b \Big[1 + c \tanh(gu/2) \Big] \Phi^{-1}(u) \Big[1 + (\Phi^{-1}(u))^2 \Big]^k,$$

with $\Phi^{-1}(\cdot)$ being the quantile function of a standard Gaussian distribution, $a \in \mathbb{R}, b > 0, g \ge 0, k \ge 0$ are location, scale, shape (affecting the skewness) and shape (affecting the kurtosis) parameters respectively.

Hidden Potts model, let i ∈ {1,..., n} denotes the notes of a lattice (e.g. pixels in an image), with y_i ∈ {1,..., k} denotes the node's ith node's state. The Potts model is a Markov random field of the form

$$P(y_i \mid y_{i \sim \ell}) = \frac{\exp\left(\theta \sum_{i \sim \ell} \mathbb{1}_{[y_i = y_\ell]}\right)}{\sum_{j=1}^k \exp\left(\theta \sum_{i \sim \ell} \mathbb{1}_{[y_i = y_\ell]}\right)}$$

where $i \sim \ell$ are neighboring nodes of i, and $\theta \geq 0$.

• Suppose we observe X_1, \ldots, X_n with

$$X_i \mid \boldsymbol{\theta} \stackrel{iid}{\sim} f(x_i \mid \boldsymbol{\theta}), \quad \text{and} \quad \boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}).$$

- The intuition of ABC methods is that, once we have a parameter value θ, we generate a set of synthetic data S₁,..., S_n.
- If the synthetic data S_1, \ldots, S_n are close enough to the observed one X_1, \ldots, X_n , then the data generating process of the synthetic data is similar to the one of the observed data.
- But the two data generating process share the same structure and are indexed by θ. Hence, if the two set of data are close enough, θ is a reasonable parameter also for X₁,..., X_n.
- Such approach can be combined with different sampling strategies and refined to obtain efficient sampler for posterior arising from intractable likelihood.
- This motivates also labeling these approaches as "likelihood free methods",

Rejection sampler, an old friend

- We recall one of the basic Monte Carlo sampler, the rejection sampler.
- Suppose we want to sample from $\pi(\theta \mid X)$ using an auxiliary distribution $h(\theta)$. With rejection sampling we can produce a sample as follows.

At the *r*th sampling step

- i) We generate $\theta^{(r)} \stackrel{iid}{\sim} h(\theta)$ (independent of the previous state).
- ii) Accept $\theta^{(r)}$ with probability

$$\frac{\pi(\boldsymbol{\theta}^{(r)} \mid \boldsymbol{X})}{Kh(\boldsymbol{\theta}^{(r)})}, \quad \text{ with } \quad K \geq \max_{\boldsymbol{\theta}} \frac{\pi(\boldsymbol{\theta}^{(r)} \mid \boldsymbol{X})}{h(\boldsymbol{\theta}^{(r)})}$$

otherwise go back to i).

- Here we still have to evaluate the likelihood function, we can extend this algorithm accommodating for synthetic data.
- The support of $h(\cdot)$ should cover the support of $\pi(\theta^{(r)} | \mathbf{X})$.

ABC rejection sampler

• A first intuitive ABC sampler can be obtained as follows.

At the rth sampling step

- i) We generate $\theta^{(r)} \stackrel{iid}{\sim} h(\theta)$ (independent of the previous state).
- ii) We generate $\boldsymbol{S} \stackrel{iid}{\sim} f(\boldsymbol{s} \mid \boldsymbol{\theta}^{(r)}).$
- ii) Accept $\theta^{(r)}$ with probability

$$rac{\pi(oldsymbol{ heta}^{(r)})}{Kh(oldsymbol{ heta}^{(r)})}, \qquad ext{with} \qquad K \geq \max_{oldsymbol{ heta}} rac{\pi(oldsymbol{ heta}^{(r)})}{h(oldsymbol{ heta}^{(r)})}$$

otherwise go back to i).

- We do not need anymore to evaluate the data density function, but just to use it as sampling mechanism.
- However, having $\boldsymbol{S}, \boldsymbol{X} \in \mathbb{X}^n \subseteq \mathbb{R}^d$, then we have $P(\boldsymbol{S} = \boldsymbol{X}) = 0$.
- Further, when the number of observation is large, using the whole sample information is slowing the sampler possibly without practical benefits.

Relaxing the matching condition

- In general, we do not assume an identity function to compare the observe data and the synthetic one.
- In fact, we require that the observed data are close enough to the synthetic one, for example requiring that

$$||\boldsymbol{X} - \boldsymbol{S}|| \geq \epsilon,$$

for some threshold $\epsilon \geq 0$ and some distance measure $||\cdot||$.

- When ε = 0, we go back on the previous rejection sampler and we produce a sample from π(θ | X). However, usually we consider ε > 0, producing a sample from an approximation of such a distribution.
- The previous can be further generalized by applying a kernel function to the previous distance, hence moving from

$$\mathbb{I}_{[||\boldsymbol{X}-\boldsymbol{S}|| \leq \epsilon]}$$
 to $\mathrm{K}_{\lambda}(||\boldsymbol{X}-\boldsymbol{S}||) = \frac{1}{\lambda}\mathrm{K}\left(\frac{||\boldsymbol{X}-\boldsymbol{S}||}{\lambda}\right).$

Common choices are the triangular kernel $K(u) = (1 - |u|)1_{[|u| \le 1]}$ and the Gaussian kernel $K(u) = \phi(u)$, with ϕ being the density function of a standard Gaussian distribution.

Simplify the comparison

- Secondly, when the number of observations increases, the comparison can be based on a reduced information instead of the whole sample.
- Further, is highly unlikely that $S \approx X$ can be generated from $f(s \mid \theta)$ for any choice of θ .
- This results in the need of a large scale parameter λ in the kernel function, in order to achieve decent acceptance rates with the rejection algorithm.
- A common practice is to consider summary statistics of both observed and synthetic data, and then use those statistics in the comparison

$$\mathrm{K}_{\lambda}(||\phi(X) - \phi(S)||)$$

where $\phi(\cdot) : \mathbb{X}^n \to \mathbb{R}^q$ is a function producing a vector of summaries of **X** or **S**.

This approach actually produce a sample from π_{ABC}(θ | φ(X)). However, if the vector of summary statistics is sufficient for the model parameters, then the approximating distribution corresponds to π_{ABC}(θ | X).

At the rth sampling step

- i) We generate $\theta^{(r)} \stackrel{iid}{\sim} h(\theta)$ (independent of the previous state).
- ii) We generate $\boldsymbol{S} \stackrel{iid}{\sim} f(s \mid \boldsymbol{\theta}^{(r)})$.
- ii) Compute the synthetic data summary statistics $\psi(\mathbf{S})$.
- iv) Accept $\theta^{(r)}$ with probability

$$\frac{K_{\lambda}(||\phi(X) - \phi(S)||)\pi(\theta^{(r)})}{Ch(\theta^{(r)})}, \quad \text{with} \quad C \ge K_{\lambda}(0) \max_{\theta} \frac{\pi(\theta^{(r)})}{h(\theta^{(r)})},$$
otherwise go back to *i*).

- The observed data summary statistics can be pre-computed once, before the sampling, saving computational time.
- The degree of approximation we are introducing is tuned by specific choices of λ , K and $||\cdot||$.

The target distribution

• In practice, we are sampling from the joint distribution

 $\pi(\boldsymbol{\theta}, \phi(\boldsymbol{S}) \mid \phi(\boldsymbol{X})) \propto \mathrm{K}_{\lambda}(||\phi(\boldsymbol{X}) - \phi(\boldsymbol{S})||) \mathrm{L}(\phi(\boldsymbol{S}) \mid \boldsymbol{\theta}) \pi(\boldsymbol{\theta}),$

and the posterior ABC distribution si defined as

$$\pi_{ABC}(\boldsymbol{ heta} \mid \phi(\boldsymbol{X})) = \int_{\mathbb{X}^n} \pi(\boldsymbol{ heta}, \phi(\boldsymbol{S}) \mid \phi(\boldsymbol{X})) \mathrm{d} \boldsymbol{S}.$$

• We can see that, as the scale parameter of the kernel decreases, we have

$$\begin{split} \lim_{\lambda \to 0} \pi_{ABC}(\boldsymbol{\theta} \mid \phi(\boldsymbol{X})) &\propto \int_{\mathbb{X}^n} \lim_{\lambda \to 0} \mathrm{K}_{\lambda}(||\phi(\boldsymbol{X}) - \phi(\boldsymbol{S})||) \mathrm{L}(\phi(\boldsymbol{S}) \mid \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{S} \\ &= \int_{\mathbb{X}^n} \delta_{\phi(\boldsymbol{X})}(\phi(\boldsymbol{S})) \mathrm{L}(\phi(\boldsymbol{S}) \mid \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) \mathrm{d}\boldsymbol{S} = \pi(\boldsymbol{\theta} \mid \phi(\boldsymbol{X})) \mathrm{d}\boldsymbol{S} \end{split}$$

If φ(·) is sufficient or it is simply the whole sample, it is apparent that the previous corresponds to π(θ | X).

How much approximate

- We can also quantify the accuracy of an ABC sampler. For simplicity of illustration, let us consider the a sampler without summary statistics, and an univariate case with a single data point, x, s ∈ X = R, and ||·|| = |·|.
- We can see that the error committed in approximating the likelihood function (i.e. ignoring the prior distribution) corresponds to

$$L_{ABC}(x \mid \boldsymbol{\theta}) = \int K_{\lambda}(|x - \boldsymbol{s}|)L(\boldsymbol{s} \mid \boldsymbol{\theta})d\boldsymbol{s} = \int K(u)L(x - u\lambda \mid \boldsymbol{\theta})du.$$

Using a Taylor expansion, we have

$$L_{ABC}(x \mid \boldsymbol{\theta}) = \int K(u) \left[L(x \mid \boldsymbol{\theta}) - u\lambda \frac{d}{dx} L(x \mid \boldsymbol{\theta}) + \frac{u^2 \lambda^2}{2} \frac{d^2}{dx^2} L(x \mid \boldsymbol{\theta}) - \cdots \right] du$$

which can be truncated, with some simplification, in

$$L_{ABC}(x \mid \boldsymbol{\theta}) \simeq L(x \mid \boldsymbol{\theta}) + \frac{\lambda^2}{2} \frac{d^2}{dx^2} L(x \mid \boldsymbol{\theta}) \int u^2 \mathcal{K}(u) du$$

· Hence, the bias we are committing is given by

$$b_{\lambda}(x \mid \boldsymbol{\theta}) = rac{\lambda^2}{2} \sigma_{\mathrm{K}}^2 rac{\mathrm{d}^2}{\mathrm{d}x^2} \mathrm{L}(x \mid \boldsymbol{\theta}),$$

where we can see that is quadratic in λ and depends on the kernel choice through $\sigma_K^2=\int u^2 K(u) du.$

- There are many way to describe and interpret ABC samplers and their impact on statistical analysis.
- Commonly (e.g. Blum, 2010), in the ABC setting we are interested in the joint sample $(\theta^{(r)}, \mathbf{S}^{(r)})$, and we produce an empirical conditional density of $\pi(\theta \mid \mathbf{X})$ by weighting the $\theta^{(r)}$ s by $||\mathbf{S}^{(r)} \mathbf{X}||$
- Fearnhead and Prangle (2012) noted that the ABC approximation of the posterior is a continuous mixture

$$\pi_{ABC}(\boldsymbol{ heta} \mid \boldsymbol{X}) \propto \int w(\boldsymbol{S}) \pi(\boldsymbol{ heta} \mid \boldsymbol{S}) \mathrm{d} \boldsymbol{S},$$

where $w(S) \propto K_{\lambda}(||S - X||)m(X)$, with m(X) being the marginal distribution of X.

• Wilkinson (2013) pointed out that ABC methods can be considered as exact if ||S - X|| represents an error term (observational error or model misspecification), and K_{λ} is the error distribution.

ABC rejection sampler with Pima dataset

- We consider again the Pima dataset, but using a ABC rejection sampler.
- We use a $N_p(\mathbf{0}, 100I_p)$ as prior distribution, $N_p(\mu, S)$ as proposal, and as kernel, we consider $K_{\lambda}(||\mathbf{S} \mathbf{X}||) = (\frac{1}{n} \sum_{i=1}^{n} I_{[S_i = X_i]})^{\lambda}$.

```
# not informed
S <- diag(100, 8)
mu <- rep(0, 8)
lambda <- 3
# Running the ABCr (R = 3000)
ABCrejectionNI <- ABCrejection(R, lambda, y, X, S, mu)
colMeans(ABCrejectionNI$param)
# X Xnpreg Xglu Xbp Xskin Xbmi Xped Xage
# -2.975 1.978 3.760 1.060 1.640 2.047 1.793 2.402
ABCrejectionNI$acc
```

0.152

• Same but using the Laplace approximation as proposal.

```
# not informed
S <- vcov(fit_logit)
mu <- coefficients(fit_logit)
lambda <- 3
# Running the ABCr (R = 3000)
ABCrejectionNI <- ABCrejection(R, lambda, y, X, S, mu)
colMeans(ABCrejectionNI$param)
# X Xnpreg Xglu Xbp Xskin Xbmi Xped Xage
# -0.991 0.405 1.097 -0.098 0.068 0.571 0.449 0.284</pre>
```

ABCrejectionNI\$acc

0.995

Improving ABC methods

- We can improve our sampler by combining the previous strategy with other sampling strategy, for example the one presented in the previous sessions.
- In general, we target the joint distribution of parameter vector and summary statistics

 $\pi(\theta, \phi(\mathbf{S}) \mid \phi(\mathbf{X})) \propto \mathrm{K}_{\lambda}(||\phi(\mathbf{S}) - \phi(\mathbf{X})||)\mathrm{L}(\mathbf{S} \mid \theta)\pi(\theta)$

- Samples of θ can be obtained by sampling jointly (θ, S) and marginalizing by discarding S.
- Ideally, we will see ...

ABC importance sampling

- We discussed that importance sampling is a procedure that, in the spirit of rejection sampler, use an instrumental distribution to propose possible values of the quantity of interest. However, rather than calculate acceptance probabilities, produce a weighted sample.
- Once we draw a candidate θ^(r) ~ q(θ^(r)), we also compute a weight associated to this value, w(θ^(r)) = π(θ^(r) | X)/q(θ^(r)).
- From the previous, it is easy to see that

$$\mathrm{E}_q[w(oldsymbol{ heta})g(oldsymbol{ heta})] = \int g(oldsymbol{ heta})w(oldsymbol{ heta})q(oldsymbol{ heta})\mathrm{d}oldsymbol{ heta} = \int g(oldsymbol{ heta})\pi(oldsymbol{ heta}\midoldsymbol{X})\mathrm{d}oldsymbol{ heta} = \mathrm{E}_\pi[g(oldsymbol{ heta})],$$

which can be approximated via Monte Carlo methods.

 When the target distribution is not normalized, we simply used a normalized version of the weights

$$W^{(r)} = \frac{w(\theta^{(r)})}{\sum_{r=1}^{R} w(\theta^{(r)})}.$$

- In the following slides, we do not consider the summary statistics case. However, everything
 here presented can be done also using summaries of the samples, both observed and
 synthetic.
- From an ABC perspective, importance sampling works similarly to rejection sampling.
- Our target distribution of the importance sampling is $\pi_{ABC}(\theta, S \mid X)$, and the proposal (importance) distribution is of the form $q(\theta, S) = L(S \mid \theta)q(\theta)$, defined jointly on the parameter and sample spaces.
- As results, the unnormalized importance weights are

$$w(oldsymbol{ heta}) \propto rac{\pi_{ABC}(oldsymbol{ heta},oldsymbol{S}\midoldsymbol{X})}{q(oldsymbol{ heta},oldsymbol{S})} \propto rac{\mathrm{K}_{\lambda}(||oldsymbol{S}-oldsymbol{X}||)\pi(oldsymbol{ heta})}{q(oldsymbol{ heta})},$$

which is free of the intractable likelihood term.

• The choice of the importance distribution $q(\theta)$ is crucial for the algorithm performances.

At the rth sampling step

i) We generate $\theta^{(r)} \stackrel{iid}{\sim} q(\theta)$ (independent of the previous state).

ii) We generate
$$S_i \stackrel{iid}{\sim} f(s \mid \theta^{(r)})$$
, for $i = 1, ..., n$.

ii) Compute the importance weights

$$w(\boldsymbol{ heta}_r) = rac{\mathrm{K}_{\lambda}(||\boldsymbol{S}-\boldsymbol{X}||)\pi(\boldsymbol{ heta}^{(r)})}{q(\boldsymbol{ heta}^{(r)})}.$$

- The result is a weighted sample, which can be possibly resampled.
- The kernel choice impacts on the resulting weights. Specifically, if the kernel has non-compact support, the importance weights are guaranteed to be always non-zero. However, this may results in high variability and low ESS of the sampler.

- We consider again the Pima dataset, but using a ABC importance sampler.
- We use a N_p(0, 100I_p) as importance distribution, N_p(μ, S) as proposal, and as kernel, we consider K_λ(||S − X||) = (¹/_n ∑ⁿ_{i=1} I_[S_i=X_i])^λ.

```
# not informed
S <- diag(100, 8)
mu <- rep(0, 8)
lambda <- 3
# Running the ABCis (R = 3000)
ABCimportanceNI <- ABCimportance(R, lambda, y, X, S, mu)
colMeans(ABCrejectionNI$param)
# X Xnpreg Xglu Xbp Xskin Xbmi Xped Xage
# -3.639 1.782 3.760 1.013 1.403 2.154 1.812 2.388
ABCrejectionNI$acc
```

ESS

• Same but with Laplace approximation as importance distribution.

```
# not informed
S <- vcov(fit_logit)
mu <- coefficients(fit_logit)
lambda <- 3
# Running the ABCis (R = 3000)
ABCimportanceI <- ABCimportance(R, lambda, y, X, S, mu)
colMeans(ABCrejectionNI$param)
# X Xnpreg Xglu Xbp Xskin Xbmi Xped Xage
# -1.064 0.348 1.089 -0.092 0.086 0.539 0.428 0.326</pre>
```

ABCrejectionNI\$acc

14.263

ABC-MCMC

- As we saw earlier in the days, MCMC can be used to define sampling strategies. In
 particular, with some regularity assumptions, it is possible to produce suitable algorithm
 whose transition kernel produce a Markov Chain which is ergodic to a specific target
 distribution.
- Among others, a broad class of algorithm is given by the Metropolis-Hastings algorithm, where given a current state $\theta^{(r)}$, we propose a candidate from $q(\theta \mid \theta^{(r-1)})$ and then accept such a value with probability

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^{(r)}) = \min\left\{1, \frac{\pi(\boldsymbol{\theta} \mid \boldsymbol{X})q(\boldsymbol{\theta}^{(r-1)} \mid \boldsymbol{\theta})}{\pi(\boldsymbol{\theta}^{(r-1)} \mid \boldsymbol{X})q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(r-1)})}\right\}.$$

- it is only natural that MCMC-based ABC algorithms have been studied, as they combine both the tractability of ABC methods, and the feature of MCMC algorithms, such as Markovian dependence over time.
- ABC MCMC algorithms were originally studied by Marjoram et al. (2003), an later extended in many directions, by considering different MCMC strategies.

- In an ABC setting, the target distribution is π_{ABC}(θ, S | X), whereas possibly we marginalize S.
- Hence, our proposal distribution works on the product space Θ × Xⁿ. In practice, we want to have (Markovian) memories only for the parameter values. Hence, our proposal usually has form

$$q(\theta, \mathbf{S} \mid \theta^{(r-1)}, \mathbf{S}^{(r-1)}) = q(\theta \mid \theta^{(r-1)})L(\mathbf{S} \mid \theta),$$

so that the synthetic data are sampled independently from the past.

• The resulting acceptance rate equals

$$\begin{aligned} \alpha(\theta, \mathbf{S} \mid \theta^{(r-1)}, \mathbf{S}^{(r-1)}) &= \min\left\{ 1, \frac{\pi_{ABC}(\theta, \mathbf{S} \mid \mathbf{X})q(\theta^{(r-1)}, \mathbf{S}^{(r-1)} \mid \theta, \mathbf{S})}{\pi_{ABC}(\theta^{(r-1)}, \mathbf{S}^{(r-1)} \mid \mathbf{X})q(\theta, \mathbf{S} \mid \theta^{(r-1)}, \mathbf{S}^{(r-1)})} \right\} \\ &= \min\left\{ 1, \frac{\mathrm{K}_{\lambda}(||\mathbf{S} - \mathbf{X}||)\pi(\theta)q(\theta^{(r-1)} \mid \theta)}{\mathrm{K}_{\lambda}(||\mathbf{S}^{(r-1)} - \mathbf{X}||)\pi(\theta^{(r-1)})q(\theta \mid \theta^{(r-1)})} \right\}. \end{aligned}$$

• Note that, the previous acceptance rate does not involve the intractable likelihood term.

At the rth sampling step

i) We generate $\theta \stackrel{iid}{\sim} q(\theta \mid \theta^{(r-1)})$.

ii) We generate
$$S_i \stackrel{iid}{\sim} f(s \mid oldsymbol{ heta})$$
, for $i=1,\ldots,n$.

ii) With probability

se

$$\alpha(\boldsymbol{\theta}, \boldsymbol{S} \mid \boldsymbol{\theta}^{(r-1)}, \boldsymbol{S}^{(r-1)}) = \min\left\{1, \frac{\mathrm{K}_{\lambda}(||\boldsymbol{S} - \boldsymbol{X}||)\pi(\boldsymbol{\theta})q(\boldsymbol{\theta}^{(r-1)} \mid \boldsymbol{\theta})}{\mathrm{K}_{\lambda}(||\boldsymbol{S}^{(r-1)} - \boldsymbol{X}||)\pi(\boldsymbol{\theta}^{(r-1)})q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(r-1)})}\right\}$$

t $(\boldsymbol{\theta}^{(r)}, \boldsymbol{S}^{(r)}) = (\boldsymbol{\theta}, \boldsymbol{S})$, otherwise $(\boldsymbol{\theta}^{(r)}, \boldsymbol{S}^{(r)}) = (\boldsymbol{\theta}^{(r-1)}, \boldsymbol{S}^{(r-1)})$.

- The algorithm satisfies the detailed balanced condition on $\pi_{ABC}(\theta, S \mid X)$.
- The algorithm can be sensible to the initial condition.
- Thus, the strategy can be combined with specific proposal (e.g. MALA).