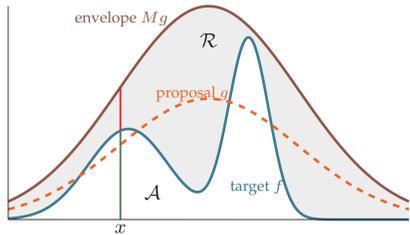


# PLIABLE REJECTION SAMPLING

## SIMPLE REJECTION SAMPLING

**Goal:** Sample from a target density  $f$  (not easy to sample from)  
**Tool:** Use a proposal density  $g$  (from which sampling is quite easy)  
**Property:** Smaller  $\mathcal{R} \implies$  fewer rejections (good!)



$M$  verifies  $f \leq Mg$ . The sampling algorithm:

- Sample  $x$  from  $g$
- Accept  $x$  as a sample from  $f$  with probability  $\frac{f(x)}{Mg(x)}$

## SETTING

Let  $d \geq 1$  and let  $f$  be a density on  $\mathbb{R}^d$ .

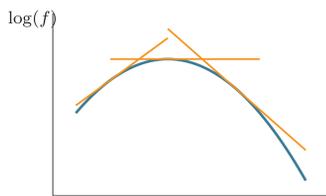
**Question:** Given a number  $n$  of requests to  $f$ , what is the number  $T$  of samples  $Y_1, \dots, Y_T$  that one can generate such that they are i.i.d. according to  $f$ ?

$$\text{acceptance rate} = \frac{T}{n}$$

Can we increase the acceptance rate?

**Adaptive Rejection Sampling (ARS)** [Gilks and Wild 1992]

- The target  $f$  is assumed to be **log-concave** (unimodal)
- The envelope is made of tangents at a set of points  $\mathcal{S}$
- At each rejection, the sample is added to  $\mathcal{S}$



**Adaptive Rejection Metropolis Sampling (ARMS)**

[Gilks, Best and Tan 1995]

- Performs a Metropolis-Hastings step for each accepted sample (which correlates the samples)

**Convex-Concave Adaptive Rejection Sampling**

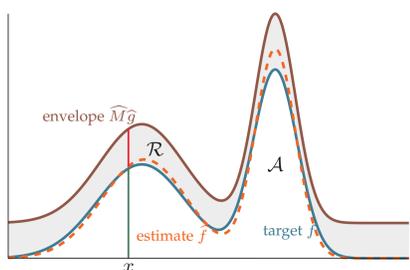
[Gorur and Tuh 2011]

- Decomposes the target as convex + concave
- Builds piecewise linear upper bounds (tangents, secant lines)

## PLIABLE REJECTION SAMPLING

**Better proposal** means smaller rejection area  $\mathcal{R}$   
**Smaller  $\mathcal{R}$**  means  $g$  should have a similar "shape" to  $f$   
**For this purpose:**

- Build an estimate  $\hat{f}$
- Translate it *uniformly*



⚠ It should be easy to sample from  $\hat{g}$  ... and  $\hat{f}$ !

## CHOICE OF THE ESTIMATE

**Assumption 1** (on the density).

- $f$  defined on  $[0, A]^d$  and bounded.
- It admits a Taylor expansion in any point up to some degree  $0 < s \leq 2$ .

**Assumption 2** (on the kernel).

- Let  $K = \prod_{i=1}^d K_0$
- $K_0$  is a density kernel: defined on  $\mathbb{R}$ , uniformly bounded, normalized and non-negative
- $K_0$  is  $\varepsilon$ -Hölder for some  $\varepsilon > 0$
- $K_0$  is of degree 2, i.e.:

$$\int_{\mathbb{R}} x K_0(x) dx = 0 \quad \text{and} \quad \int_{\mathbb{R}} x^2 K_0(x) dx < \infty$$

Let  $X_1, \dots, X_N \sim \mathcal{U}_{[0, A]^d}$ . The (modified) kernel regression estimate is

$$\hat{f}(x) = \frac{A^d}{N h^d} \sum_{k=1}^N f(X_k) K\left(\frac{X_k - x}{h}\right) \quad (1)$$

$K_0$  Gaussian kernel  $\implies \hat{f}$  is a Gaussian mixture!

## BOUNDING THE GAP

**Theorem 1.** The estimate  $\hat{f}$  is such that with probability larger than  $1 - \delta$ , for any point  $x \in [0, A]^d$ ,

$$|\hat{f}(x) - f(x)| \leq H_0 \left( \left( \frac{\log(NAd/\delta)}{N} \right)^{\frac{\alpha}{2\alpha+d}} \right)$$

where  $H_0$  is a constant that depends on the problem parameters.

## THE PLIABLE PROPOSAL

- Remaining requests to  $f$ :  $n - N$
- Let  $r_N = A^d H_C \left( \frac{\log(NAd/\delta)}{N} \right)^{\frac{\alpha}{2\alpha+d}}$
- Construct the *pliable* proposal  $\hat{g}$  out of  $\hat{f}$ :

$$\hat{g} = \frac{\hat{f} + r_N \mathcal{U}_{[0, A]^d}}{\frac{1}{N} \sum_{i=1}^N f(X_i) + r_N}$$

- Perform rejection sampling using  $\hat{g}$  and the empirical rejection sampling constant

$$\hat{M} = \frac{\frac{1}{N} \sum_i f(X_i) + r_N}{\frac{1}{N} \sum_i f(X_i) - 5r_N}$$

## ALGORITHM: PRS

**Parameters:**  $s, n, \delta, H_C$

**Initial sampling**

Draw uniformly at random  $N$  samples on  $[0, A]^d$  and evaluate  $f$  on them

**Estimation of  $f$**

Estimate  $f$  by  $\hat{f}$  on these  $N$  samples (Equation 1)

**Generating the samples**

Sample  $n - N$  samples from the compact pliable proposal  $\hat{g}^*$

Perform rejection sampling on these samples

using  $\hat{M}$  as a rejection constant to get  $\hat{n}$  samples

**Output:** Return the  $\hat{n}$  samples

## NUMBER OF ACCEPTED SAMPLES

**Theorem 2.** Under Theorem 1's assumptions and if:

$$H_0 < H_C \quad \bullet \quad 8r_N \leq \int_{[0, A]^d} f(x) dx$$

For  $n$  large enough, we have with probability larger than  $1 - \delta$  that

$$\hat{n} \geq n \left[ 1 - \mathcal{O} \left( \frac{\log(nAd/\delta)}{n} \right)^{\frac{\alpha}{3\alpha+d}} \right].$$

Convergence Rate  $\uparrow$  with smoothness

Convergence Rate  $\downarrow$  with dimensionality

## PRS PROPERTIES

- PRS deals with a wider class of functions and not necessarily normalized
- PRS has guarantees: asymptotically we accept everything (whp).
- PRS is a **perfect sampler** (whp) the samples are iid (unlike MCMC)
- PRS empirical performance is comparable to state of the art
  - PRS deals better with peakiness than A\* sampling
  - in general PRS does not scale to high dimensions
- An extension to densities with unbounded support is provided

**Some notes on (very) high dimensionality**

Let the  $\gamma$ -support of  $f$  be

$$\text{Supp}_{f, \gamma} = \bar{\Lambda}_{f, \gamma} \quad \text{where} \quad \Lambda_{f, \gamma} \stackrel{\text{def}}{=} \{x \in \mathcal{D} : f(x) > \gamma\}$$

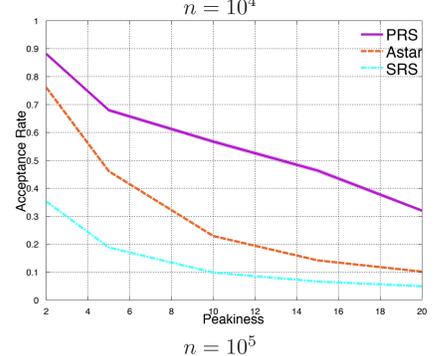
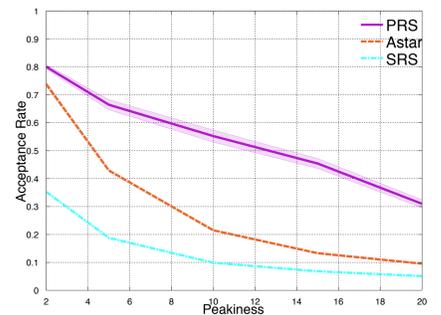
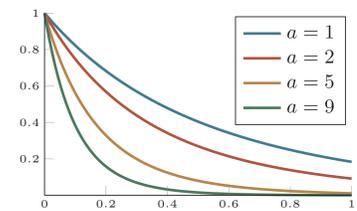
- In general, localizing the 0-support of  $f$  may require an exponential cost in  $d$
- $\text{Supp}_{f, \gamma}$  is localizable with a less than exponential cost, if:
  - $f|_{\text{Supp}_{f, \gamma}^c}$  the restriction of  $f$  on the complementary of  $\text{Supp}_{f, \gamma}$  is convex
  - One can evaluate  $f$  and its gradient pointwise

$\text{Supp}_{f, \gamma-\varepsilon}$  localization cost  $\mathcal{O}(d^2/\varepsilon^2)$

**The trick:** find a point  $x_0$  in  $\text{Supp}_{f, \gamma}^c$  and use standard gradient based optimization to find a maximum on  $\partial \text{Supp}_{f, \gamma}$ .

## EXPERIMENTS - SCALING WITH PEAKINESS

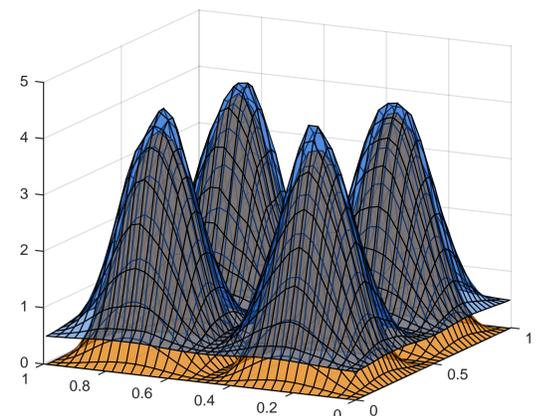
$f \propto \frac{e^{-x}}{(1+x)^a}$ , parameter  $a$  defines the **peakiness** level



## EXPERIMENTS - MULTIMODAL EXAMPLES

**A 2D example**

$$f(x, y) \propto \left(1 + \sin\left(4\pi x - \frac{\pi}{2}\right)\right) \left(1 + \sin\left(4\pi y - \frac{\pi}{2}\right)\right).$$



$n = 10^6$	acceptance rate	standard deviation
PRS	66.4%	0.45%
A* sampling	76.1%	0.80%
SRS	25.0%	0.01%

**The Clutter Problem** [Thomas P. Minka, UAI '01]

- Consider  $K$  data points  $(X_i)_{i=1}^K$  with half within  $[-5, -3]^d$  and half within  $[2, 4]^d$

- These points are assumed to be generated from

$$p(x|\theta) = (1 - \pi)\mathcal{N}(x; \theta, I) + \pi\mathcal{N}(x; 0, 10I)$$

- We put a gaussian prior on the mean  $p(\theta) = \mathcal{N}(\theta; 0, 100I)$
- The goal is to sample from  $p(\theta | (X_i)_{i=1}^K) \propto p(\theta) \prod_{i=1}^K p(X_i | \theta)$

$n = 10^5$ , 1D	acceptance rate	standard deviation
PRS	79.5%	0.2%
A* sampling	89.4%	0.8%
SRS	17.6%	0.1%

$n = 10^5$ , 2D	acceptance rate	standard deviation
PRS	51.0%	0.4%
A* sampling	56.1%	0.5%
SRS	$2.1 \cdot 10^{-3}\%$	$10^{-5}\%$

## EXTENDING THIS WORK

**Iterative version:**

- PRS is a 2 step algorithm: estimation + RS
- Possibly improve the estimate on several steps
- Optimize the number of samples gathered between these "estimation update steps"