# New sampler for cosmology

Anže Slosar, Brookhaven National Laboratory

BNL group meeting, 12/16

# *Introduction*

- This is a half-baked project that I've been sitting on for over a year
- Now it is handed to Jose, who is fumbling around trying to incorporate it into SimpleMC
- It is a method to sample likelihood surface, whose main advantage is *scalability*: you don't win by having fewer likelihood evaluations but instead by having a algorithm that scales better to thousands of cores.
- Allows one to get answers in minutes on a couple of thousand of cores

# *The problem*

- Question: How to get marginalised constraints for cosmological parameters when you have likelihood has $N > 10$ dimension and each evaluation is computationally expensive ($>$1s)
- Answer: use `CosmoMC` which runs Markov Chain Monte Carlo (MCMC)
- MCMC is an algorithm that "walks" around the likelihood and produces *samples*
- Integrals over likelihood can be converted to sums over samples

MCMC

# Problems with MCMC

Markov Chain Monte Carlo does not scale very well:

- Scales perfectly for small number of chains, but not on modern architectures with 1000s of cores
- To run a `CosmoMC` chain you still run on 64 cores and wait for two days, instead of running on 10000 cores and wait 20 mins.
- But can't you run 1000 chains?
- Yes, but burn-in is a constant time process: one always needs to throw away some $\sim 1000$ steps, because either:
  - You start chains randomly – they need to burn in
  - You start chains at high-likelihood region - they need to decorrelate
- Both are inefficient: You take 1000 samples to burn-in, but then 100 samples on each chain is enough to get 100,00 samples – quite inefficient
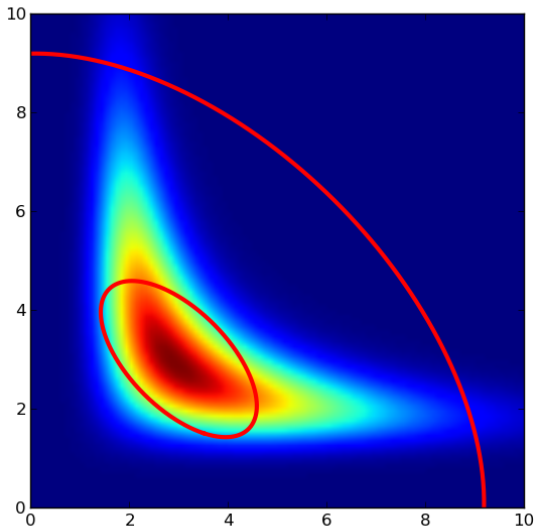
# *Importance Sampling*

Assuming that one can sample from a known distribution, then one can weight samples to recover the effective sampling from a target distribution (whose properties one would like to study)
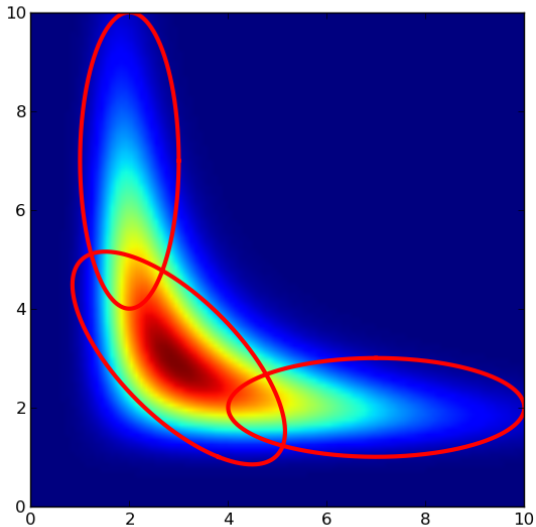
$$w_i = A \frac{L_t(\mathbf{x}_i)}{L_s(\mathbf{x}_i)}, \tag{1}$$

- ▶ Used to add a dataset to existing chains
- ▶ People tried to use it to sample cosmological likelihood directly using a Gaussian, but it fails miserably with bananas:
  - ▶ Either your Gaussian does not encompass the banana: weights blow up at the edges
  - ▶ Your Gaussian covers the banana, but also empty volume around it: most weights zero.

# Why naive importance sampling doesn't work

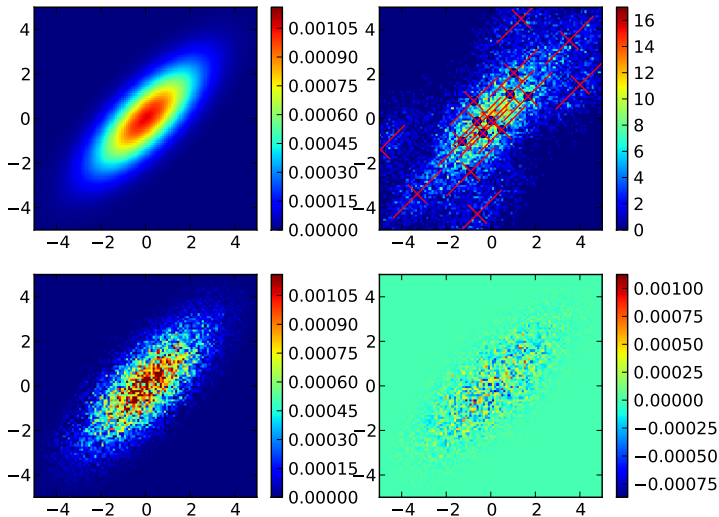# *Guassian embedding sampler*

1. Populate a list of Gaussians with a single Gaussian centered at a chosen starting point with suitable covariance.

2. Take $N$ samples from the most recently added Gaussian in the list.
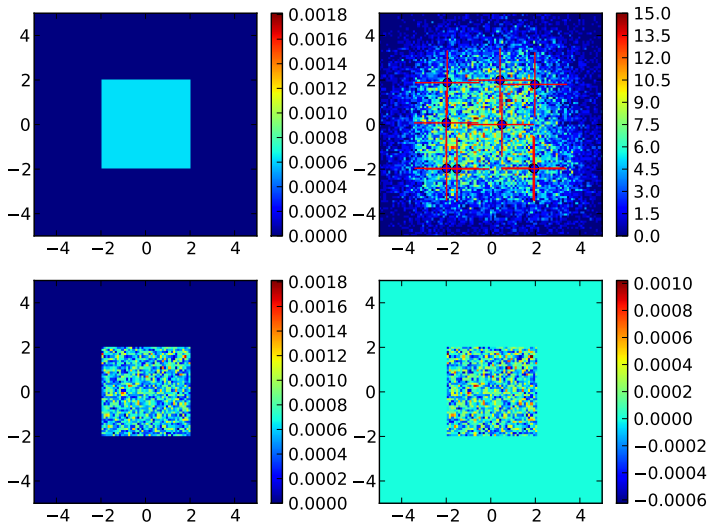
3. Calculate importance sample weights,

$$w_i = A \frac{L_t(\mathbf{x}_i)}{\sum_{j=1\ldots M} G_j(\mathbf{x}_i - \mu_j, \mathrm{C}_j)}, \qquad (2)$$

4. Add new Gaussian at the position of the largest importance weight
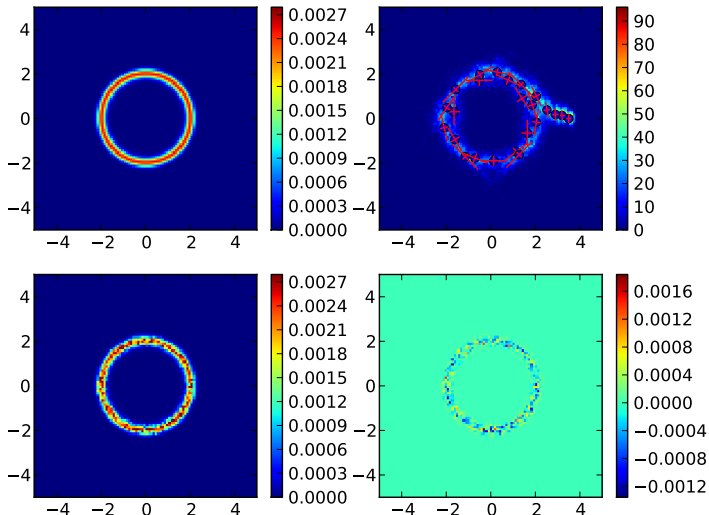
5. GOTO step 2

# Test 3: Doughnut shaped banana

# *Convergence*

Tried several convergence tests. Jose is now working on my latest idea, namely that

$$c = \frac{\Delta\theta \mathrm{C}^{-1}\Delta\theta}{N_p}, \qquad (3)$$

where $\Delta\theta$ is the shift in mean parameter values when the last Gaussian is added, is smaller than a certain number.