

## Sampling II

Side note: if you scroll through the past lectures, you will see I often created fake data for the lecture. [Fake stars, fake telescope noise, fake random fields.] It is these sampling techniques which do the maths behind those pictures.

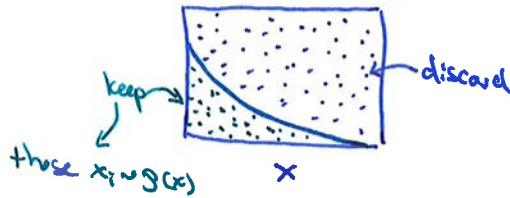
• We had the basics of sampling:

1)  $\mathcal{M}(\vec{\theta})$  &  $\mathcal{P}(\vec{\theta})$

2) if  $x_i \sim \mathcal{P}(x|M)$  then  $\int dx \mathcal{P}(x|M) g(x) \approx \frac{1}{N} \sum_{n=1}^N g(x_n)$   
 Samples

3) Markov chain:  $x_i$  depends only on  $x_{i-1}$ , not  $x_{i-n}$  with  $n > 1$ .

4) Sampling via rejection sampling or by transformation of probabilities



$$\mathcal{P}(x) dx = \mathcal{P}(y) dy$$

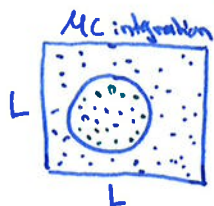
Mostly, these algorithms occur simultaneously / in conjunction / and protocol other interesting variables along the way.

Example: the chains you analyzed in the tutorials protocolled weights,  $z^2$  and log-likelihoods as auxiliary variables.

Q: what is the difference between Monte Carlo Integration and rejection sampling?

A: 1) They solve different problems: • MC-Integration happens to use random numbers, but would like the resulting integral to be a "classical" (noise-free) value.  $\approx 10^5$   
 • Rejection sampling ~~also~~ generates random numbers which follow a distribution.

2) Also the algorithms differ slightly:

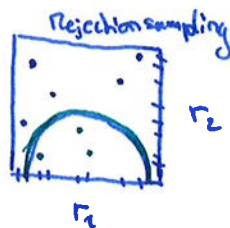


$\rightarrow$  2D Integral

$\rightarrow$  computes fraction inside/outside

$$A_{\text{square}} = L^2$$

$$A_{\text{circle}} = ?$$



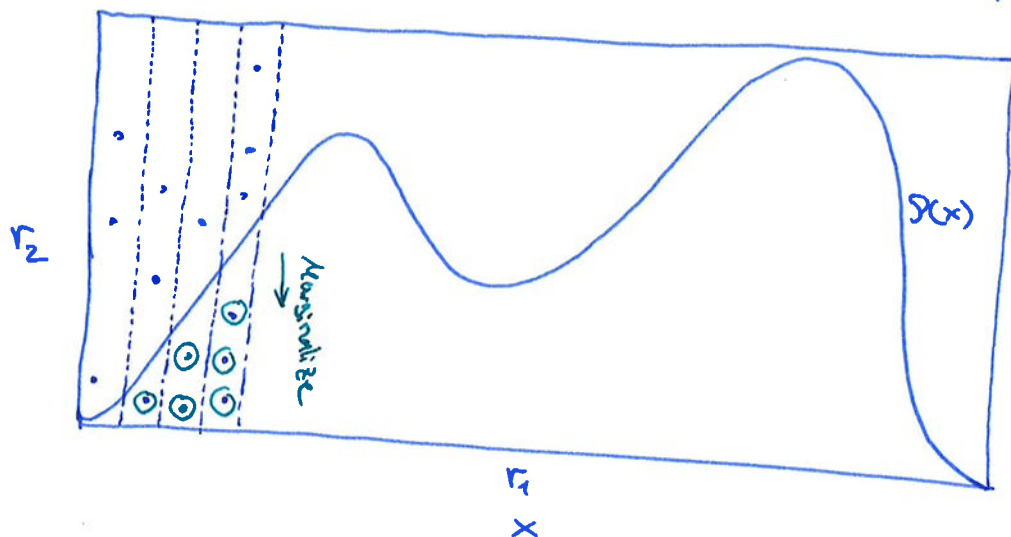
$r_2 \leftarrow$  this is an auxiliary samples?

$\rightarrow$  2D sampling, but 1D distribution

$\rightarrow$  Aim: Samples that follow the turquoise (semi-circle) distribution

## • Rejection sampling, zoomed in: why does it work?

Had: draw  $x$ , compute  $\mathcal{P}(x)$ . Draw  $q \sim \text{Unif}[0,1]$ . If  $\mathcal{P}(x) \geq q$ , keep  $x$ .  
 [also occurs in the accept/reject step of Metropolis algorithm for MCMC]



$$\mathcal{P}(r_1, r_2) = \mathcal{P}(r_1) \mathcal{P}(r_2)$$

$$\Rightarrow \mathcal{P}(r_1) = \int_0^1 \mathcal{P}(r_1, r_2) dr_2 = \mathcal{P}(r_1) \cdot 1 \quad \text{because } \int \mathcal{P}(r_2) dr_2 = 1$$

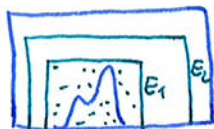
But now we keep only those  $r_2$ s below the curve  $\mathcal{P}(x)$

$$\int_0^{\mathcal{P}(x)} \mathcal{P}(r_1, r_2) dr_2 = \int_0^{\mathcal{P}(x)} \mathcal{P}(r_1) \mathcal{P}(r_2) dr_2 = \mathcal{P}(x)$$

• Gibbs-Sampling: a sampler without tunable parameters (good!)  
 (also a sampler which lets you get around evaluating  $\mathcal{P}(x)$ , which can be under- and over-flow unstable)  $\downarrow$  a exp(...)

### • Tunable parameters:

- 1) Rejection sampling:
  - maximum of  $\mathcal{P}(x)$  needs to be estimated
  - all samples need to be discarded if  $\mathcal{P} > \mathcal{P}^{\text{thought}}$  is found [i.e. maximum was not correct]
  - more general: sampling envelope needs to be tuned for efficiency



$\Rightarrow$  envelope  $E_1$  more efficient than  $E_2$

### 2) Metropolis algorithm:

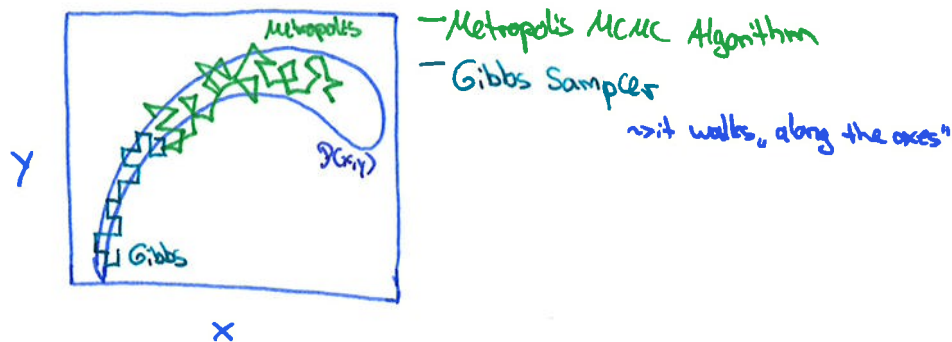
$\vec{\Delta}$ , the proposal step width and direction



$\Rightarrow$  for "wrong"  $\vec{\Delta}$ , the acceptance ratio will drastically decrease.

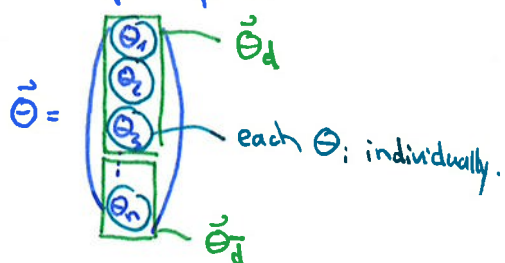
→ the Gibbs sampler does not have free parameters which affect its performance.

†: for many problems, Gibbs sampling is naturally highly effective (few correlations)



- Preconditions for Gibbs-Sampling: If  $P(x, y, z)$  is to be sampled, then all conditionals need to be known.
- Let  $\vec{\Theta}$  be the collection of all random variables in the game.

→ Ways to split up  $\vec{\Theta}$ :



+ Need all  $P(\vec{\Theta}_d | \vec{\Theta}_{-d})$

Gibbs sampling in Pseudo code:

0) initialize  $\vec{\Theta}$

1) For  $i=0, \dots, N_{\text{samples}}$ :

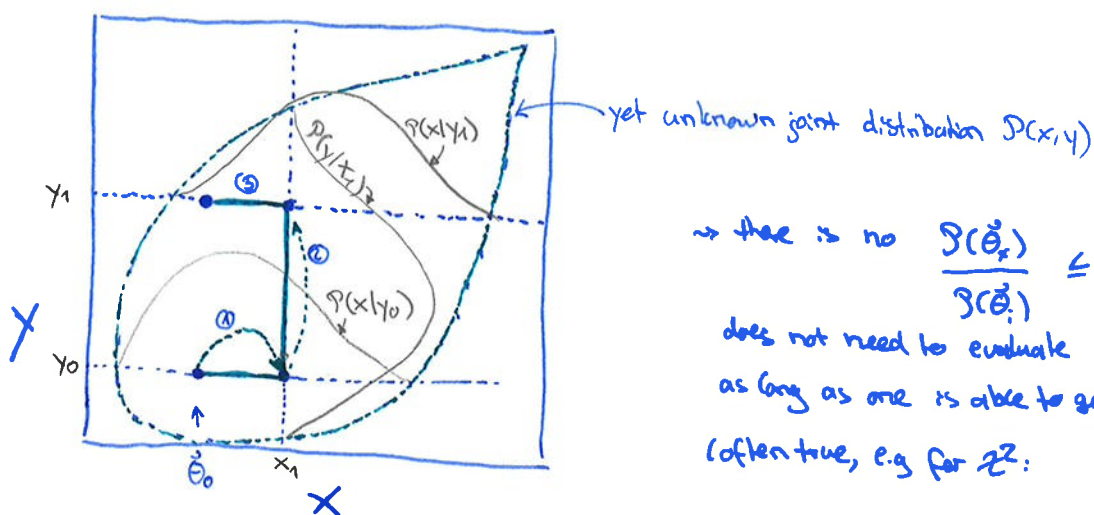
a)  $\forall d$  draw  $\vec{\Theta}_d \sim P(\vec{\Theta}_d | \vec{\Theta}_{-d})$

b) replace elements of  $\vec{\Theta}$  with  $\vec{\Theta}_d$

c) [randomize in which order components  $d$  are drawn]

↑  
for performance reasons

Example:



→ there is no  $\frac{P(\vec{\Theta}_d)}{P(\vec{\Theta})} \leq \alpha$ , hence one

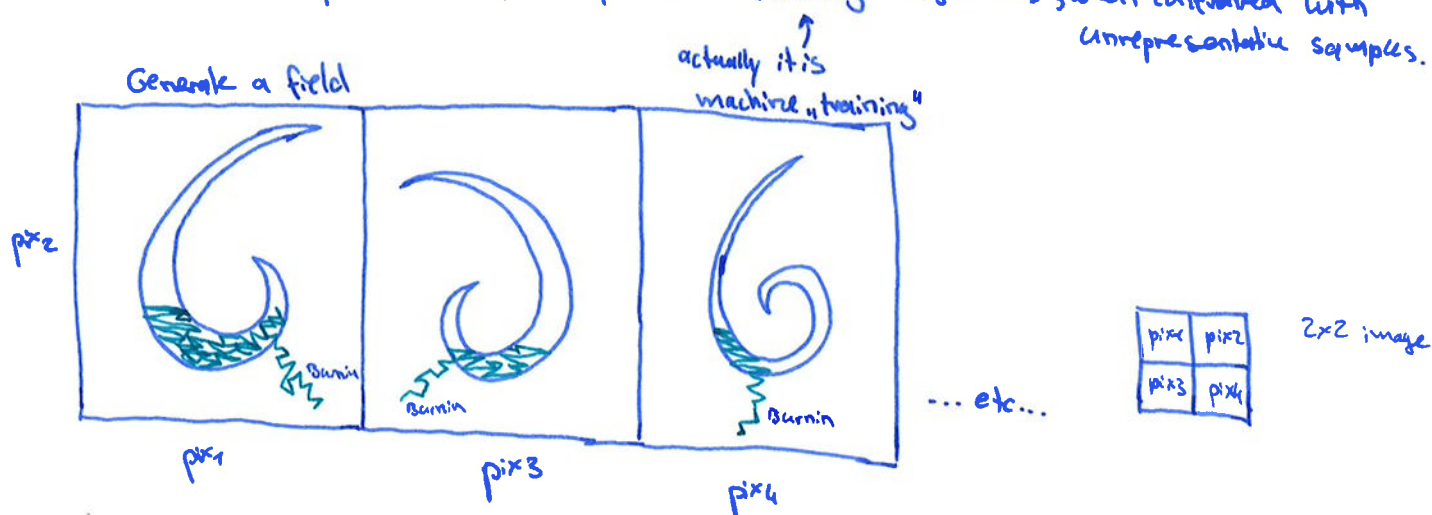
does not need to evaluate the  $P$  function, as long as one is able to generate samples from it

(often true, e.g. for  $\mathbb{R}^2$ :  $x \sim \mathcal{N}^2_p : x = \sum_{i=1}^p x_i^2$ ,  $x_i \sim G(0, 1)$ )



• The No 1 problem with sampling: when samples are not representative of the distribution

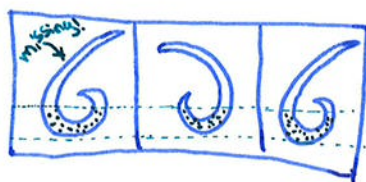
→ it is the same problem as faced by machine "learning" algorithms, when confronted with unrepresentative samples.



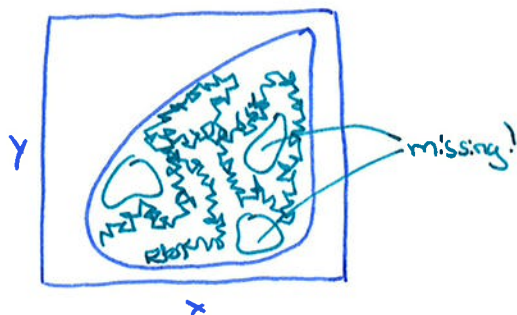
1) Burnin not representative

2) Sampler gets stuck: it doesn't succeed to get around the curve.

Result: it generates samples which capture only part of the variability.

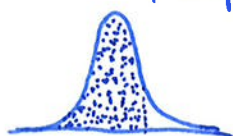


Another problem are strong correlations in chains.



→ if a artificial neural network were now trained on those 2x2 images, then it would not "learn" (= "fit to") the true distribution. Because unrepresentative samples cause biases.

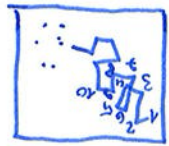
Extremely simple case:



→ obviously one Gaussian tail is missing

The #1 user-error for using samples which are representative of the true distribution.

→ using neighbouring samples → don't!



Example: MCMC:

French: together! All samples count!  
 100 does not depend on 2  
 10 does not depend on 1  
 ...  
 4 depends on 3  
 3 depends on 2  
 2 depends on 1

user selects samples 1, 2, 3.  
 ↓  
 user selection induces errors, although samples may fine.

→ Samples form an ensemble. If you use samples from an MCMC chain for post processing, Pick randomly from the chain!

The #2 user-error of MCMC sampling: concatenating unconverged chains

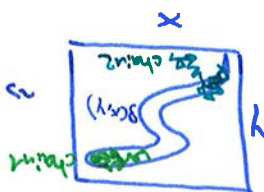
→ converged chain:  $n \times \mathcal{D}$   
 → 2nd converged chain:  $n_2 \times \mathcal{D}$

Since both converged → can concatenate (= "paste both into same .txt file")

→ 1st unconverged chain:  $n_1 \times \mathcal{D}$   
 → 2nd unconverged chain:  $n_2 \times \mathcal{D}$

concatenation does not improve the quality of the chain! chain + chain still not  $\propto \mathcal{D}$ !

Example:



→ chain + chain still misses middle part of distribution

# Mini-crash course: Machine "Learning" <sup>Training</sup>

~~Distribution~~  $\Rightarrow$  Samples  $\Rightarrow$  "Learning from Examples"

$\rightarrow$  then called "training data"

traditional inference

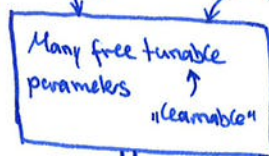
$\rightarrow$  since we gave up on the distribution [in general] replace it by "a loss function"

if  $D(x)$  is the sampling distribution of data  $x$ , and I want to learn about parameters  $\theta$ , then

$D(x|\theta)$  is my likelihood.

$\partial_{\theta} D(x|\theta) \stackrel{!}{=} 0$  then determines my best fit

Training data + lossfunction



$\leftarrow$  "the Machine" e.g. neural net, SVM, etc

output: values of the tunable parameters, for which the loss function is small.

$\Downarrow$   
save to disk.

Simplest example:

$$\text{the "machine"} \rightarrow A \vec{x}_i \stackrel{!}{=} \vec{b}_i \quad \forall i$$

$\uparrow$   
data

$\uparrow$   
the physics which generated it

$A_{11}$	$A_{12}$	$A_{13}$	$A_{14}$
$A_{21}$	$A_{22}$	$A_{23}$	$A_{24}$
$A_{31}$	$A_{32}$	$A_{33}$	$A_{34}$
$A_{41}$	$A_{42}$	$A_{43}$	$A_{44}$

$\uparrow$   
the "learnable" (= "fittable") parameters

$x_0$
$x_1$
$x_2$
$x_3$

$\uparrow$   
many trainings data

$\hat{=}$

$b_0$
$b_1$
$b_2$
$b_3$

$\uparrow$   
many provided wanted outcomes

$\oplus$  Loss function: e.g.  $|\vec{b}_M - \vec{b}_{true}|$

Aim: find  $A$  such that it maps  $\vec{x}_i$  correctly to  $\vec{b}_i$ , by fitting the  $A_{ij}$  elements.

$\rightarrow$  obviously this is mathematically solvable (linear algebra), hence "learnable" from unbiased examples.

After  $N \rightarrow$  many training runs with  $N$  examples, feed in one more data vector, this time the true one. Since  $\frac{N+1}{N}$  small (because  $N \gg 1$ ), loss will again be small (if all went well).

$\rightarrow$  it "learned sth".