

HOW TO BECOME THE MOST HATED BAND IN THE WORLD:

RECORD AN ALBUM THAT'S NOTHING BUT BRILLIANT, CATCHY INSTANT CLASSICS GUARANTEED POPULARITY AND AIRTIME,



WITH A SAMPLE OF A CAR HORN, CELL PHONE, OR ALARM CLOCK INSERTED RANDOMLY IN EACH SONG.

<https://xkcd.com/780/>

Sampling can be hard to get right ...

Class update:

Exam info posted. Friday “week 13”. 3 June ,

Week 12 Monday - no lecture, finish your video assignment + good luck with all end-of-S1 deadlines!

Week 12 Wednesday - Q&A, Talk about selected questions in 2021 Exam.

“Week 13” – drop-in sessions will be announced.

Sampling methods

Motivation, why? Sampling and ML

Basic sampling algorithms

- Sampling standard distributions from $U(0, 1)$
- Rejection sampling
- Importance sampling

Markov chain Monte Carlo (MCMC)

- Markov chains
- Metropolis-Hastings

Gibbs sampling 11.3

Bishop Chap 11

Intro, 11.1.1, 11.1.2, 11.1.4, 11.1.6
11.2, 11.2.1, 11.2.2, 11.2.3
11.3

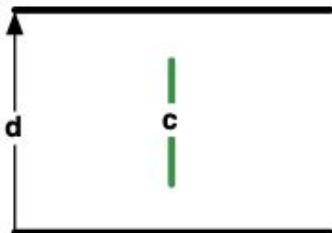
Conceptual ,

Estimating π – Buffon's needle

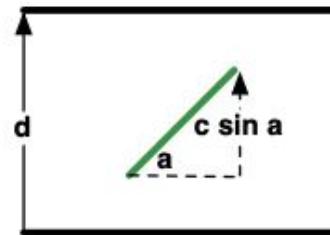
- Drop length c needle on parallel lines distance d apart
- Needle falls perpendicular:
Probability of crossing the line is c/d .
- Needle falls at an arbitrary angle a :
Probability of crossing the line $c \sin(a)/d$.
- Every angle is equally probable. Calculate the mean:

$$p(\text{crossing}) = \frac{c}{d} \int_0^\pi \sin(a) \, dp(a) = \frac{1}{\pi} \frac{c}{d} \int_0^\pi \sin(a) \, da = \frac{2}{\pi} \frac{c}{d}$$

(iv) n crossings in N experiments results in $\frac{n}{N} \approx \frac{2}{\pi} \frac{c}{d}$



(i) Needle falls
perpendicular ($a = \pi/2$).



(ii) Needle falls at
arbitrary angle a .

- $N \rightarrow \infty$ estimate "correct"
- how fast?

<https://mathworld.wolfram.com/BuffonsNeedleProblem.html>

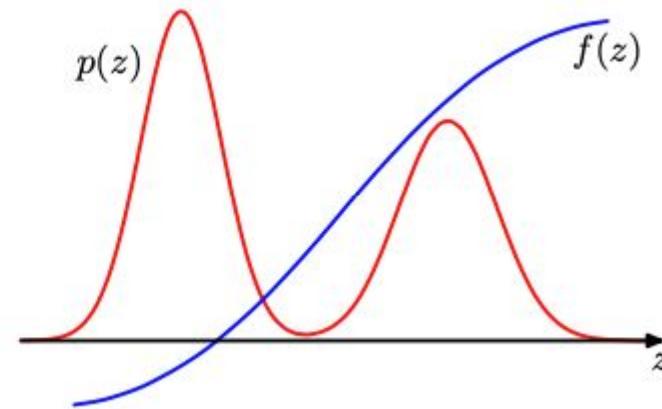
Why do we need sampling?

Distributions can be quite complex, e.g. posterior distributions, mixture distributions, graphical models (coming next).

Many ML tasks can be stated as estimating the expectation of functions under a distribution, e.g. posterior mean and variance, moments.

$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z}) d\mathbf{z} \quad (11.1)$$

Figure 11.1 Schematic illustration of a function $f(z)$ whose expectation is to be evaluated with respect to a distribution $p(z)$.



Bayesian logistic regression

(GP2)

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) \quad (4.87)$$

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\} \quad (4.90)$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) \quad (4.140) \quad \rightarrow \text{prior}$$

$$p(\mathbf{w}|\mathbf{t}) \propto p(\mathbf{w})p(\mathbf{t}|\mathbf{w}) \quad (4.141)$$

$$\begin{aligned} \ln p(\mathbf{w}|\mathbf{t}) &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &\quad + \sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln(1-y_n)\} + \text{const} \end{aligned} \quad (4.142)$$

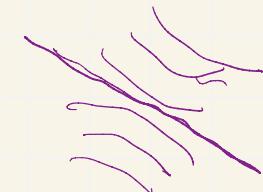
$$\mathbf{S}_N = -\nabla \nabla \ln p(\mathbf{w}|\mathbf{t}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1-y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T. \quad (4.143)$$

Laplace approximation $q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{S}_N).$

$$\underline{p(\mathcal{C}_1|\phi, \mathbf{t})} = \int p(\mathcal{C}_1|\phi, \mathbf{w})p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} \quad (4.145)$$

$$= \int \underbrace{\sigma(a)}_{\text{sample}} \underbrace{\mathcal{N}(a|\mu_a, \sigma_a^2)}_{\text{sample}} da.$$

$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z}) d\mathbf{z} \quad (11.1)$$



Sampling 101

Goal: estimate

$$\mathbb{E}[f] = \int f(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \quad (11.1)$$

Simply

sample from $p(\mathbf{z})$ i.i.d.

$$\hat{f} = \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}^{(l)}). \quad (11.2)$$

$$\underline{\text{var}}[\hat{f}] = \frac{1}{L} \mathbb{E} [(f - \mathbb{E}[f])^2] \quad (11.3)$$

$$\mathbb{E}[\hat{f}] = \mathbb{E}[f]$$

The good:

- the accuracy of the estimator does not depend on the dimensionality of \mathbf{z}
- in principle, high accuracy may be achievable with a relatively small number of samples

"effective # of samples"

The bad:

- samples $\{\mathbf{z}^{(l)}\}$ might not be independent, and so the effective sample size might be much smaller than the apparent sample size.

The good:

- the accuracy of the estimator does not depend on the dimensionality of z
- in principle, high accuracy may be achievable with a relatively small number of samples

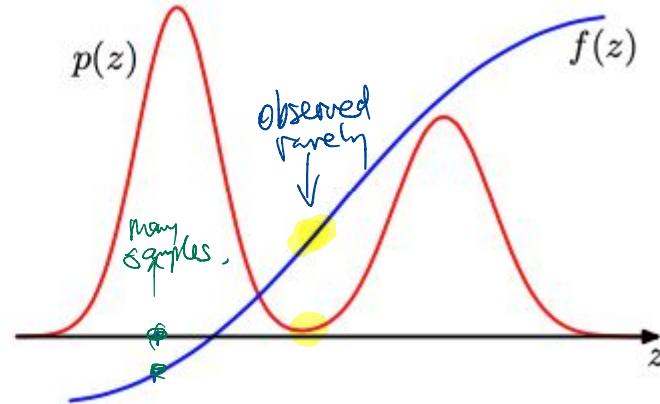
The bad:

- samples $\{z^{(l)}\}$ might not be independent, and so the effective sample size might be much smaller than the apparent sample size.

The ugly:

$f(z)$ is small in regions where $p(z)$ is large, and vice versa, then the expectation may be dominated by regions of small probability, implying that **relatively large sample sizes will be required to achieve sufficient accuracy**.

Figure 11.1 Schematic illustration of a function $f(z)$ whose expectation is to be evaluated with respect to a distribution $p(z)$.



Sampling from uniform distributions

- In a computer usually via pseudorandom number generator : an algorithm generating a sequence of numbers that approximates the properties of random numbers.
- Use a mathematically well crafted pseudorandom number generator.
- From now on we will assume that we have a good pseudorandom number generator for uniformly distributed data available.
- If you don't trust any algorithm :
Three carefully adjusted radio receivers picking up atmospheric noise to provide real random numbers at
<http://www.random.org/>

For this class: assume we can generate $z \sim U(0, 1)$

Q: given $z \sim U(0, 1)$
How do you obtain $z \sim U(a, b)$?

Sampling a distribution $p(y)$

$$p(y) = p(z) \left| \frac{dz}{dy} \right| \quad (11.5)$$

$$z = h(y) \equiv \int_{-\infty}^y p(\hat{y}) d\hat{y} \quad (11.6)$$

Multiple variables

$$p(y_1, \dots, y_M) = p(z_1, \dots, z_M) \left| \frac{\partial(z_1, \dots, z_M)}{\partial(y_1, \dots, y_M)} \right|. \quad (11.9)$$

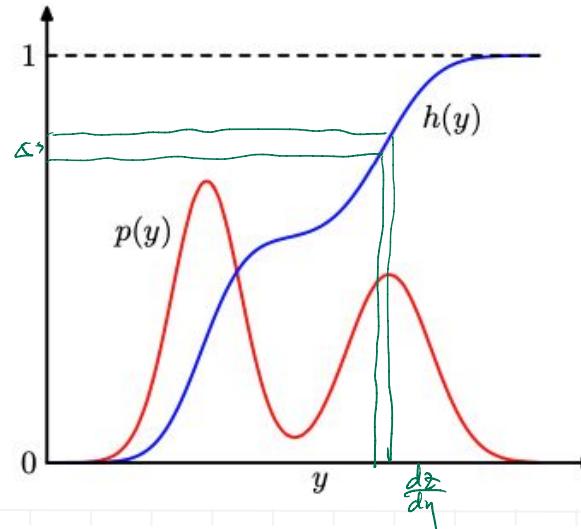
Steps:

- Generate $z \sim U(0, 1)$
- Use inverse CDF $y=h^{-1}(z)$ to obtain y

determinant of Jacobian

$$M \times M \quad ij \quad \frac{\partial z_i}{\partial y_j}$$

Figure 11.2 Geometrical interpretation of the transformation method for generating nonuniformly distributed random numbers. $h(y)$ is the indefinite integral of the desired distribution $p(y)$. If a uniformly distributed random variable z is transformed using $y = h^{-1}(z)$, then y will be distributed according to $p(y)$.



Example: exponential distribution

$$p(y) = \lambda \exp(-\lambda y)$$

$$(11.7) \quad 0 \leq y < \infty.$$

cdf \underline{z}

$$h(y) = 1 - \exp(-\lambda y)$$

$$y = -\lambda^{-1} \ln(1 - z)$$

$$1 - z = \exp^{-\lambda y}$$

$$\log(1 - z) = -\lambda y.$$

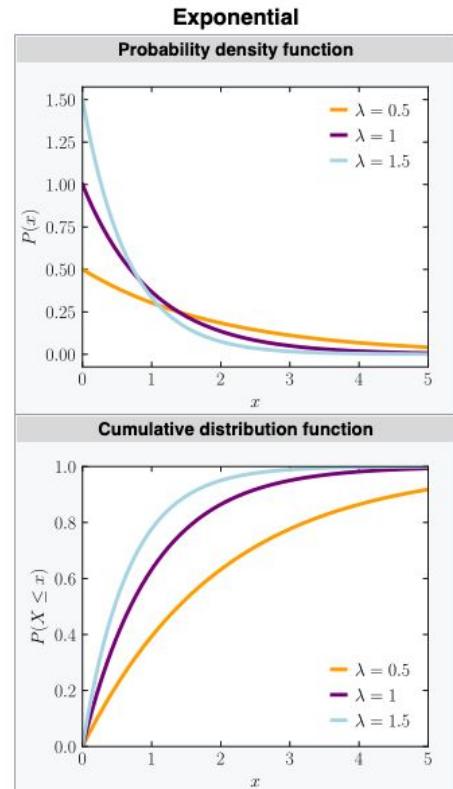
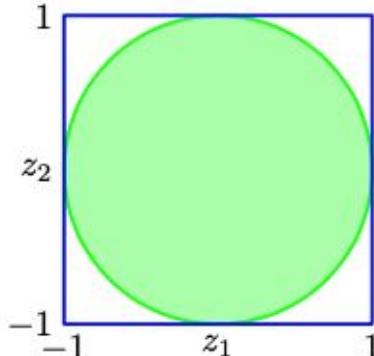


Figure 11.3 The Box-Muller method for generating Gaussian distributed random numbers starts by generating samples from a uniform distribution inside the unit circle.



$$r^2 \geq z_1^2 + z_2^2$$

Typo in the book (online pdf version, fixed in the latest print version)

$$y_1 = z_1 \left(\frac{-2 \ln z_1}{r^2} \right)^{1/2} \quad (11.10)$$

$$y_2 = z_2 \left(\frac{-2 \ln z_2}{r^2} \right)^{1/2} \quad (11.11)$$



$$y_1 = z_1 \left(\frac{-2 \ln r^2}{r^2} \right)^{1/2}$$

$$y_2 = z_2 \left(\frac{-2 \ln r^2}{r^2} \right)^{1/2}$$

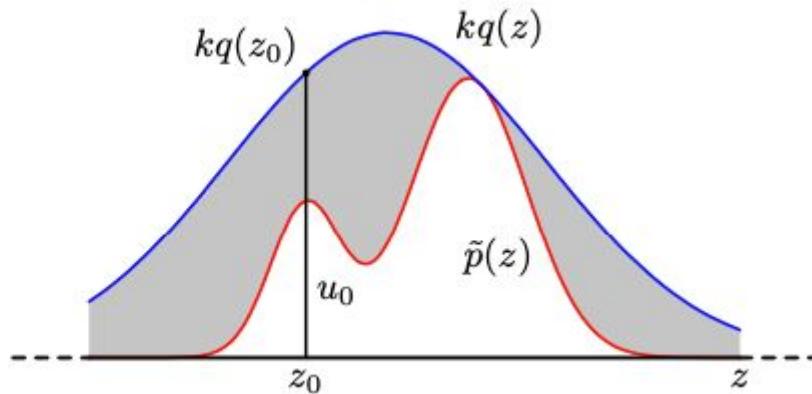
$$\begin{aligned} p(y_1, y_2) &= p(z_1, z_2) \left| \frac{\partial(z_1, z_2)}{\partial(y_1, y_2)} \right| \\ &= \left[\frac{1}{\sqrt{2\pi}} \exp(-y_1^2/2) \right] \left[\frac{1}{\sqrt{2\pi}} \exp(-y_2^2/2) \right] \end{aligned} \quad (11.12)$$

Full derivation: through change of variable to the polar coordinates (r, θ) , somewhat involved

$$r \sim U(0, 1), \quad \theta \sim U(0, 2\pi)$$

Rejection sampling

Figure 11.4 In the rejection sampling method, samples are drawn from a simple distribution $q(z)$ and rejected if they fall in the grey area between the unnormalized distribution $\tilde{p}(z)$ and the scaled distribution $kq(z)$. The resulting samples are distributed according to $p(z)$, which is the normalized version of $\tilde{p}(z)$.



Assumption 1 : Sampling directly from $p(z)$ is difficult, but we can evaluate $p(z)$ up to some unknown normalisation constant Z_p

$$p(z) = \frac{1}{Z_p} \tilde{p}(z)$$

Assumption 2 : We can draw samples from a simpler distribution $q(z)$ and for some constant k and all z holds

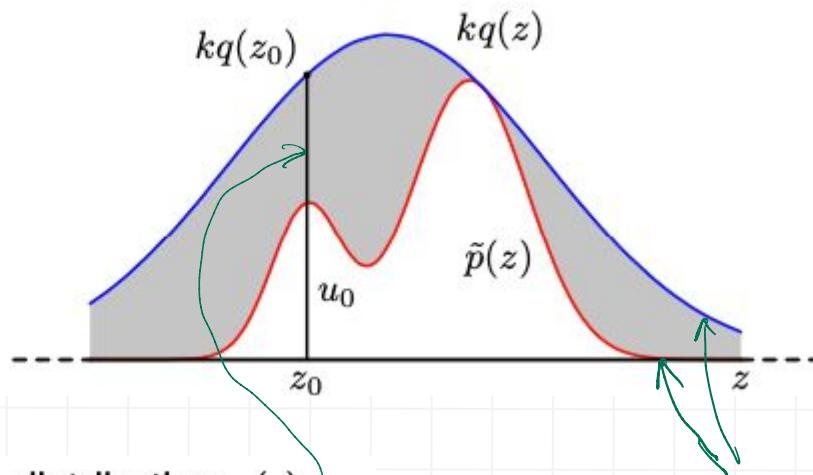
$$\underline{kq(z)} \geq \tilde{p}(z)$$

Impossible:

$$q(z) \geq p(z) \quad \forall z$$

Rejection sampling

Figure 11.4 In the rejection sampling method, samples are drawn from a simple distribution $q(z)$ and rejected if they fall in the grey area between the unnormalized distribution $\tilde{p}(z)$ and the scaled distribution $kq(z)$. The resulting samples are distributed according to $p(z)$, which is the normalized version of $\tilde{p}(z)$.



- 1 Generate a random number z_0 from the distribution $q(z)$.
- 2 Generate a number from the u_0 from the uniform distribution over $[0, k q(z_0)]$.
- 3 If $u_0 > \tilde{p}(z_0)$ then reject the pair (z_0, u_0) .
- 4 The remaining pairs have uniform distribution under the curve $\tilde{p}(z)$.
- 5 The z values are distributed according to $p(z)$.

*q(z) should have
"heavier" tails
than p(z)*

Rejection sampling example

Figure 11.5 Plot showing the gamma distribution given by (11.15) as the green curve, with a scaled Cauchy proposal distribution shown by the red curve. Samples from the gamma distribution can be obtained by sampling from the Cauchy and then applying the rejection sampling criterion.

Sample Gamma distribution with $a > 1$

$$\text{Gam}(z|a, b) = \frac{b^a z^{a-1} \exp(-bz)}{\Gamma(a)}$$

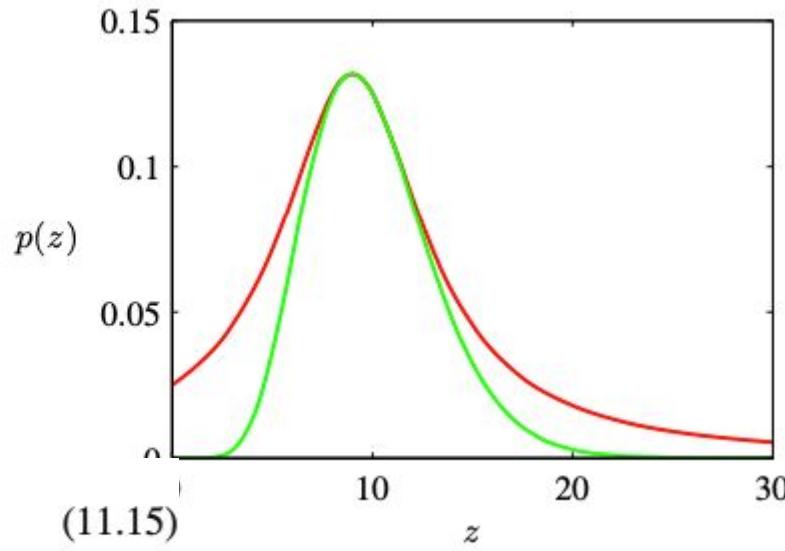
Proposal distribution: Cauchy

$$p(y) = \frac{1}{\pi} \frac{1}{1 + y^2}. \quad (11.8)$$

Or,

$$q(z) = \frac{k}{1 + (z - c)^2/b^2}$$

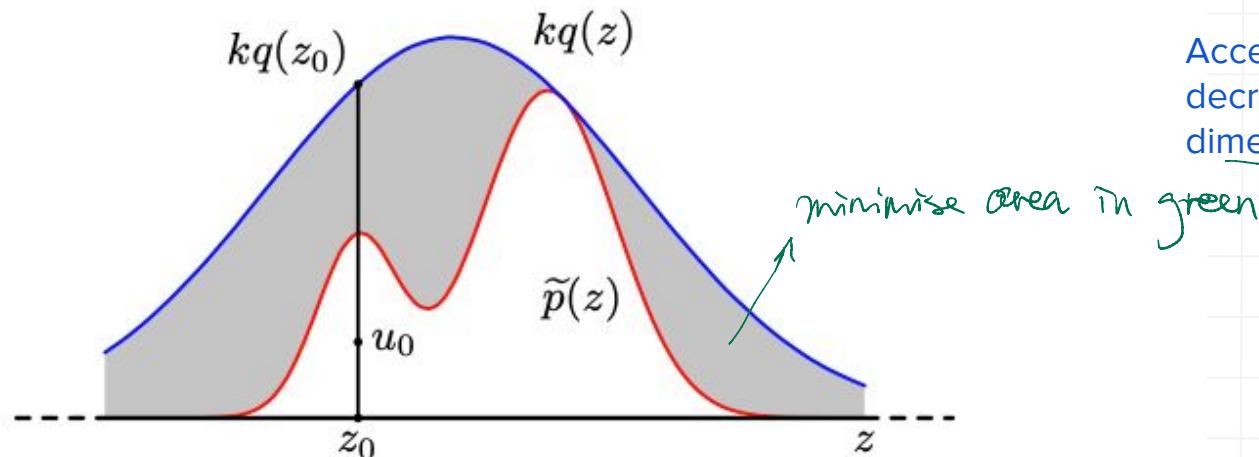
$$\begin{aligned} u &\sim U[0, 1] \\ z &= b \tan u + c \end{aligned}$$



Rejection sampling - challenges

- Need to find a proposal distribution $q(z)$ which is a **close upper bound** to $p(z)$; otherwise many samples are rejected.
- Curse of dimensionality for multivariate distributions.

Proposal distribution need to have heavier tails than the target distribution.



Acceptance rate can decrease exponentially as dimensionality increases.

Importance sampling

Desired: $\mathbb{E}[f] \simeq \sum_{l=1}^L p(\mathbf{z}^{(l)}) f(\mathbf{z}^{(l)}).$ (11.18)

- Directly calculate $\mathbb{E}_p [f(z)]$ w.r.t. some $p(z).$
- Does not sample $p(z)$ as an intermediate step.
- Again use samples $z^{(l)}$ from some proposal $q(z).$

All samples retained.

$$\begin{aligned}\mathbb{E}[f] &= \int f(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\ &= \int f(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})} q(\mathbf{z}) d\mathbf{z} \\ &\simeq \frac{1}{L} \sum_{l=1}^L \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})} f(\mathbf{z}^{(l)}).\end{aligned}\quad (11.19)$$

Importance weights

$$r_l = p(\mathbf{z}^{(l)})/q(\mathbf{z}^{(l)})$$

r_l can be > 1

Estimating the (ratio of) normalization constants

Easy to evaluate distribution up to a normalising constant, but hard to know what the constant is. $p(\mathbf{z}) = \tilde{p}(\mathbf{z}) / Z_p$ $q(\mathbf{z}) = \tilde{q}(\mathbf{z}) / Z_q$

$$\begin{aligned}\mathbb{E}[f] &= \int f(\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\ &= \frac{Z_q}{Z_p} \int f(\mathbf{z}) \frac{\tilde{p}(\mathbf{z})}{\tilde{q}(\mathbf{z})} q(\mathbf{z}) d\mathbf{z} \\ &\approx \frac{Z_q}{Z_p} \frac{1}{L} \sum_{l=1}^L \tilde{r}_l f(\mathbf{z}^{(l)}). \quad (11.20)\end{aligned}$$

$$\begin{aligned}\tilde{r}_l &\neq r_l = \frac{p(\mathbf{z})}{q(\mathbf{z})} \\ &\rightarrow \frac{Z_p}{Z_q} r_l.\end{aligned}$$

$$\begin{aligned}\frac{Z_p}{Z_q} &= \frac{1}{Z_q} \int \tilde{p}(\mathbf{z}) d\mathbf{z} = \int \frac{\tilde{p}(\mathbf{z})}{\tilde{q}(\mathbf{z})} q(\mathbf{z}) d\mathbf{z} \\ &\approx \frac{1}{L} \sum_{l=1}^L \tilde{r}_l \quad \xrightarrow{\text{average of unnormalized importance weights!}} \quad (11.21)\end{aligned}$$

Using the same samples!

Importance sampling - comments

- Importance weights r_l correct the bias introduced by sampling from the proposal distribution $q(z)$ instead of the wanted distribution $p(z)$.
- Success depends on how well $q(z)$ approximates $p(z)$.
- If $p(z) > 0$ in same region, then $q(z) > 0$ necessary.

effective sample size can be much smaller than the apparent sample size L . The problem is even more severe if none of the samples falls in the regions where $p(z)f(z)$ is large. In that case, the apparent variances of r_l and $r_l f(z^{(l)})$ may be small even though the estimate of the expectation may be severely wrong. Hence a major drawback of the importance sampling method is the potential to produce results that are arbitrarily in error and with no diagnostic indication. This also highlights a key requirement for the sampling distribution $q(z)$, namely that it should not be small or zero in regions where $p(z)$ may be significant.

! Support of
 $q(z)$
is a superset of
support of $p(z)$

Sampling and the EM algorithm

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \int p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{Z}, \mathbf{X}|\boldsymbol{\theta}) d\mathbf{Z}. \quad (11.28)$$

$\overbrace{p(\mathbf{Z})}$ $\overbrace{f(\mathbf{z})}$

Sample this $\sum^{(l)}$

Monte Carlo EM

$$\underbrace{Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})}_{\text{---}} \simeq \frac{1}{L} \sum_{l=1}^L \ln p(\mathbf{Z}^{(l)}, \mathbf{X}|\boldsymbol{\theta}). \quad (11.29)$$

$\overbrace{\sum_{l=1}^L}$

M-step: as usual,

Being Bayesian in EM

IP Algorithm

I-step. We wish to sample from $p(\mathbf{Z}|\mathbf{X})$ but we cannot do this directly. We therefore note the relation

$$p(\mathbf{Z}|\mathbf{X}) = \int p(\mathbf{Z}|\boldsymbol{\theta}, \mathbf{X})p(\boldsymbol{\theta}|\mathbf{X}) d\boldsymbol{\theta} \quad (11.30)$$

and hence for $l = 1, \dots, L$ we first draw a sample $\boldsymbol{\theta}^{(l)}$ from the current estimate for $p(\boldsymbol{\theta}|\mathbf{X})$, and then use this to draw a sample $\mathbf{Z}^{(l)}$ from $p(\mathbf{Z}|\boldsymbol{\theta}^{(l)}, \mathbf{X})$.

P-step. Given the relation

$$p(\boldsymbol{\theta}|\mathbf{X}) = \int p(\boldsymbol{\theta}|\mathbf{Z}, \mathbf{X})p(\mathbf{Z}|\mathbf{X}) d\mathbf{Z} \quad (11.31)$$

we use the samples $\{\mathbf{Z}^{(l)}\}$ obtained from the I-step to compute a revised estimate of the posterior distribution over $\boldsymbol{\theta}$ given by

$$p(\boldsymbol{\theta}|\mathbf{X}) \simeq \frac{1}{L} \sum_{l=1}^L p(\boldsymbol{\theta}|\mathbf{Z}^{(l)}, \mathbf{X}). \quad (11.32)$$

By assumption, it will be feasible to sample from this approximation in the I-step.

Sampling methods

Motivation, why? Sampling and ML

Basic sampling algorithms

- Sampling standard distributions from $U(0, 1)$
- Rejection sampling
- Importance sampling

Markov chain Monte Carlo (MCMC)

- Markov chains
- Metropolis-Hastings

Gibbs sampling

Assume: ① can evaluate $p(z)$

for any given z .

② evaluate $p(z)|_{z=z_0}$

is easier than

Sampling $z \sim p(z)$



involves inverse CDF

Markov Chain Monte Carlo (MCMC)

$$E_p[f(z)]$$

Goal: sample from $p(z)$.

Generate a sequence using a Markov chain.

- ① Generate a new sample $z^* \sim q(z | z^{(l)})$, conditional on the previous sample $z^{(l)}$.
- ② Accept or reject the new sample according to some appropriate criterion.

$$\underline{z}^{(l+1)} = \begin{cases} z^* & \text{if accepted} \\ z^{(l)} & \text{if rejected} \end{cases}$$

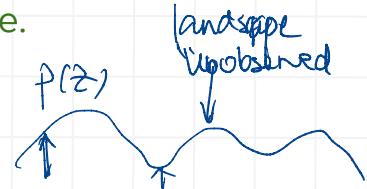
copy of the prev value.

- ③ For an appropriate proposal and corresponding acceptance criterion, as $l \rightarrow \infty$, $\underline{z}^{(l)}$ approaches an independent sample of $p(z)$.

Motivation:

Cover high-probability regions of $p(z)$ -- can we move towards it?

Scale better with the dimensionality of the sample space.



Metropolis Algorithm

- ① Choose a symmetric proposal distribution
 $q(z_A | z_B) = q(z_B | z_A)$.

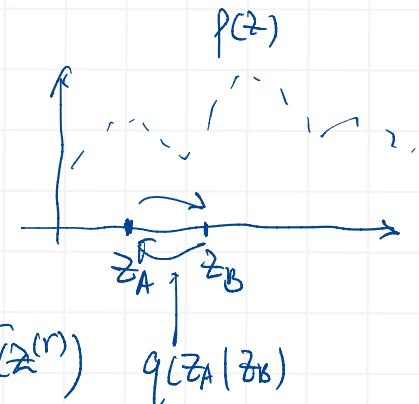
- ② Accept the new sample z^* with probability

$$A(z^*, z^{(r)}) = \min \left(1, \frac{\tilde{p}(z^*)}{\tilde{p}(z^{(r)})} \right)$$

e.g., let $u \sim \text{Uniform}(0, 1)$ and accept if $\frac{\tilde{p}(z^*)}{\tilde{p}(z^{(r)})} > u$.

- ③ Unlike rejection sampling we include the previous sample on rejection of the proposal:

$$z^{(r+1)} = \begin{cases} z^* & \text{if accepted} \\ z^{(r)} & \text{if rejected} \end{cases}$$



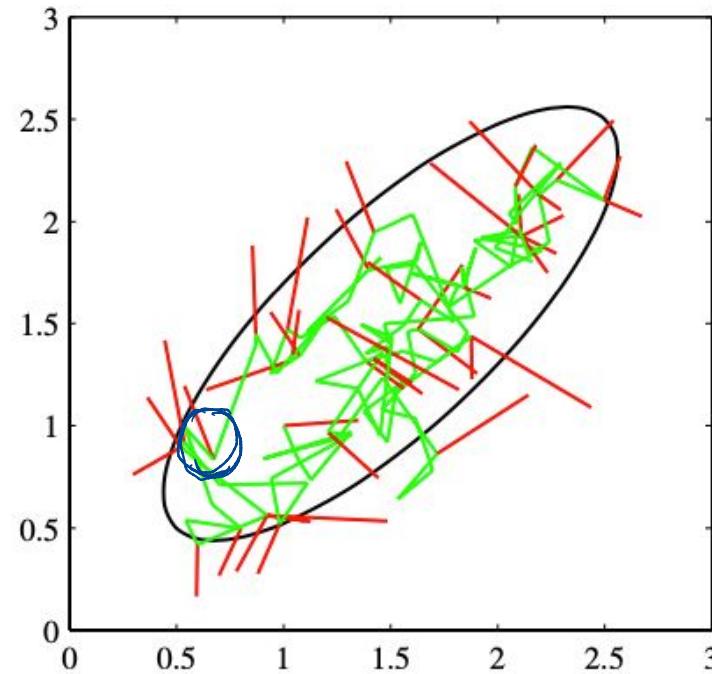
When $\tilde{p}(z^*) > \tilde{p}(z^{(r)})$
then accept

When $\tilde{p}(z^*) < \tilde{p}(z^{(r)})$
accept w. $\frac{P(z^*)}{\tilde{p}(z^{(r)})} < 1$

a copy of $z^{(r)}$

Metropolis Algorithm - illustration

Figure 11.9 A simple illustration using Metropolis algorithm to sample from a Gaussian distribution whose one standard-deviation contour is shown by the ellipse. The proposal distribution is an isotropic Gaussian distribution whose standard deviation is 0.2. Steps that are accepted are shown as green lines, and rejected steps are shown in red. A total of 150 candidate samples are generated, of which 43 are rejected.



Random walk as a Markov Chain

Consider a state space z consisting of the integers, with probabilities

$$p(z^{(\tau+1)} = z^{(\tau)}) = 0.5 \quad (11.34)$$

$$p(z^{(\tau+1)} = z^{(\tau)} + 1) = 0.25 \quad (11.35)$$

$$p(z^{(\tau+1)} = z^{(\tau)} - 1) = 0.25 \quad (11.36)$$

if $z^{(1)} = 0$,

then $\mathbb{E}[z^{(\tau)}] = 0$

$\mathbb{E}[(z^{(\tau)})^2] = \tau/2$.

after τ steps, the random walk has only travelled a distance that on average is proportional to the square root of τ
→ Random walks are very inefficient in exploring the state space.

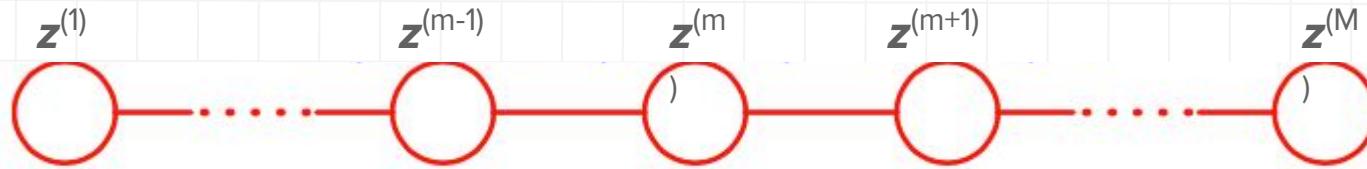
Design goal of MCMC algorithms: avoid random-walk behaviour.

First order Markov Chain

A series of random variable $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(M)}$ such that

$$p(\mathbf{z}^{(m+1)} | \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}) = p(\mathbf{z}^{(m+1)} | \mathbf{z}^{(m)}). \quad (11.37)$$

chain shaped
graphical model



Marginal probability

$$\begin{aligned} p(z^{(m+1)}) &= \sum_{z^{(m)}} p(z^{(m+1)} | z^{(m)}) p(z^{(m)}) \\ &= \sum_{z^{(m)}} T_m(z^{(m)} | z^{(m+1)}) p(z^{(m)}) \end{aligned}$$

where $T_m(z^{(m)} | z^{(m+1)})$ are the transition probabilities.

MCMC - why it works

- Marginal probability

$$p(z^{(m+1)}) = \sum_{z^{(m)}} T_m(z^{(m)} | z^{(m+1)}) p(z^{(m)})$$

drop $z^{(m+1)}$

- A Markov chain is called **homogeneous** if the transition probabilities are the same for all m , denoted by $T(z', z)$.
- A distribution is **invariant**, or **stationary**, with respect to a Markov chain if each step leaves the distribution invariant.
- For a homogeneous Markov chain, the distribution $p^*(z)$ is invariant if

$$\underbrace{p^*(z)}_{\mathbf{z}} = \sum_{\mathbf{z}'} \underbrace{T(\mathbf{z}', \mathbf{z})}_{\mathbf{z}'} \underbrace{p^*(\mathbf{z}')}_{\mathbf{z}'}. \quad (11.39)$$

(Note: There can be many. If T is the identity matrix, every distribution is invariant.)

↑ we stay put at all values of \mathbf{z} .

MCMC - why it works

- **Detailed balance**

$$p(z \rightarrow z')$$

$$p(z' \rightarrow z)$$

$$\underbrace{p^*(z)T(z, z')}_{\text{check}} = \underbrace{p^*(z')T(z', z)}_{\text{put (2) into (1)}}$$

is sufficient (but not necessary) for $p^*(z)$ to be invariant (to check, put (2) into (1)). (A Markov chain that respects the detailed balance is called **reversible**.)

- A Markov chain is **ergodic** if it converges to the invariant distribution irrespective of the choice of the initial conditions. The invariant distribution is then called **equilibrium**.
- An ergodic Markov chain can have only one equilibrium distribution.
- Why is it working? Choose the transition probabilities T to satisfy the detailed balance for our goal distribution $p(z)$.

We want

$$p^*(z) = \sum_{z'} T(z', z)p^*(z').$$

(11.39)

$\forall z$

(11.40)

repeat many steps
 $p^*(z) \rightarrow p(z)$.

Desired: Avoid rejection

The Metropolis-Hastings Algorithm

- Generalisation of the Metropolis algorithm for nonsymmetric proposal distributions q_k .
- At step τ , draw a sample z^* from the distribution $q_k(z|z^{(\tau)})$ where k labels the set of possible transitions.
- Accept with probability

Metropolis requires $q_k(z^{(t)}|z^*) = q_k(z^*|z^{(t)})$

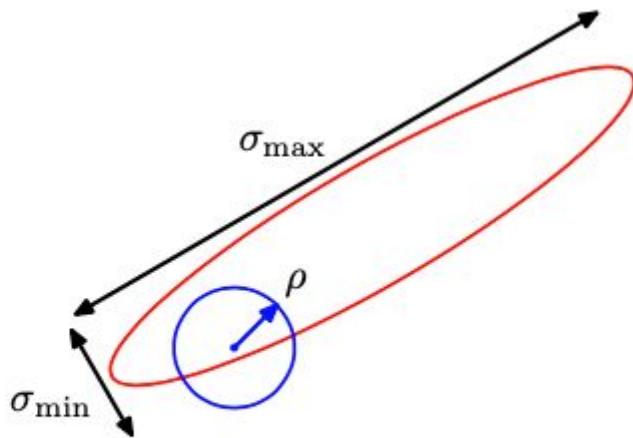
$$A_k(z^*, z^{(\tau)}) = \min \left(1, \frac{\tilde{p}(z^*) q_k(z^{(\tau)}|z^*)}{\tilde{p}(z^{(\tau)}) q_k(z^*|z^{(\tau)})} \right). \quad (11.44)$$

- Choose a symmetric proposal distribution $q(z_A|z_B) = q(z_B|z_A)$.

- Choice of proposal distribution critical.
 - Common choice : Gaussian centered on the current state.
 - small variance \rightarrow high acceptance rate, but slow walk through the state space; samples not independent
 - large variance \rightarrow high rejection rate

Figure 11.10

Schematic illustration of the use of an isotropic Gaussian proposal distribution (blue circle) to sample from a correlated multivariate Gaussian distribution (red ellipse) having very different standard deviations in different directions, using the Metropolis-Hastings algorithm. In order to keep the rejection rate low, the scale ρ of the proposal distribution should be on the order of the smallest standard deviation σ_{\min} , which leads to random walk behaviour in which the number of steps separating states that are approximately independent is of order $(\sigma_{\max}/\sigma_{\min})^2$ where σ_{\max} is the largest standard deviation.



$$A_k(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min \left(1, \frac{\tilde{p}(\mathbf{z}^*) q_k(\mathbf{z}^{(\tau)} | \mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)}) q_k(\mathbf{z}^* | \mathbf{z}^{(\tau)})} \right)$$

Metropolis-Hastings: why does it work?

- Transition probability of this Markov chain is

$$T(z, z') = q_k(z' | z) A_k(z', z)$$

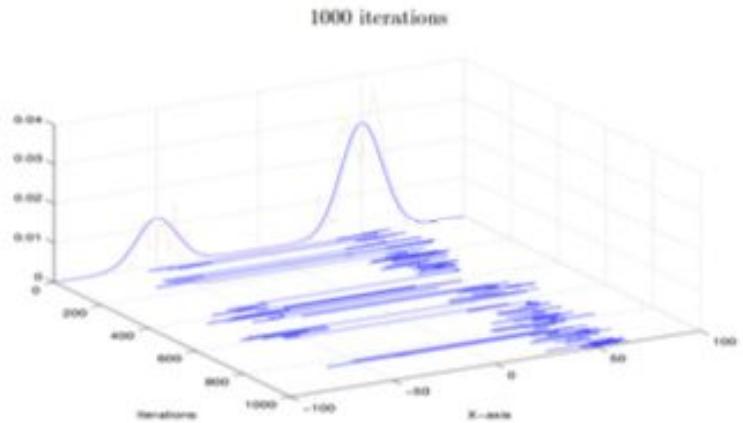
- Prove that $p(z)$ is the invariant distribution if the detailed balance holds

$$p(z) T(z, z') = T(z', z) p(z').$$

- Using the symmetry $\min(a, b) = \min(b, a)$ it can be shown that the detailed balance holds

$$\begin{aligned} p(\mathbf{z}) q_k(\mathbf{z} | \mathbf{z}') A_k(\mathbf{z}', \mathbf{z}) &= \min(p(\mathbf{z}) q_k(\mathbf{z} | \mathbf{z}'), p(\mathbf{z}') q_k(\mathbf{z}' | \mathbf{z})) \\ &= \min(p(\mathbf{z}') q_k(\mathbf{z}' | \mathbf{z}), p(\mathbf{z}) q_k(\mathbf{z} | \mathbf{z}')) \\ &= p(\mathbf{z}') q_k(\mathbf{z}' | \mathbf{z}) A_k(\mathbf{z}, \mathbf{z}') \end{aligned} \tag{11.45}$$

Metropolis-Hastings for sampling from mixture of Gaussians:



<http://www.cs.ubc.ca/~arnaud/stat535/slides10.pdf>

- With a random walk q we may get stuck in one mode.
- We could have **proposal be mixture** between random walk and “mode jumping”.

Gibbs Sampling

- Goal: sample from a joint distribution $p(\mathbf{z}) = p(z_1, \dots, z_M)$
- How? sample one variable from the distribution conditioned on all the other variable
- Example : given $p(z_1, z_2, z_3)$
- At step τ we have samples $z_1^{(\tau)}$, $z_2^{(\tau)}$ and $z_3^{(\tau)}$.
- Get samples for the next step $\tau + 1$

$$z_1^{(\tau+1)} \sim p(z_1^{(\tau)} | z_2^{(\tau)}, z_3^{(\tau)})$$

$$z_2^{(\tau+1)} \sim p(z_2^{(\tau)} | z_1^{(\tau+1)}, z_3^{(\tau)})$$

$$z_3^{(\tau+1)} \sim p(z_3^{(\tau)} | z_1^{(\tau+1)}, z_2^{(\tau+1)})$$

Gibbs sampling - why does it work?

- ① $p(\mathbf{z})$ is invariant of each of the Gibbs sampling steps and hence of the whole Markov chain : a) $p(z_i | \{\mathbf{z}_{\setminus i}\})$ is invariant because the marginal distribution $p(\mathbf{z}_{\setminus i})$ does not change, and b) by definition each step samples from $p(z_i | \{\mathbf{z}_{\setminus i}\})$.
- ② Ergodicity: sufficient that none of the conditional distribution is zero anywhere.
- ③ Gibbs sampling is a Metropolis-Hastings sampling in which each step is accepted.

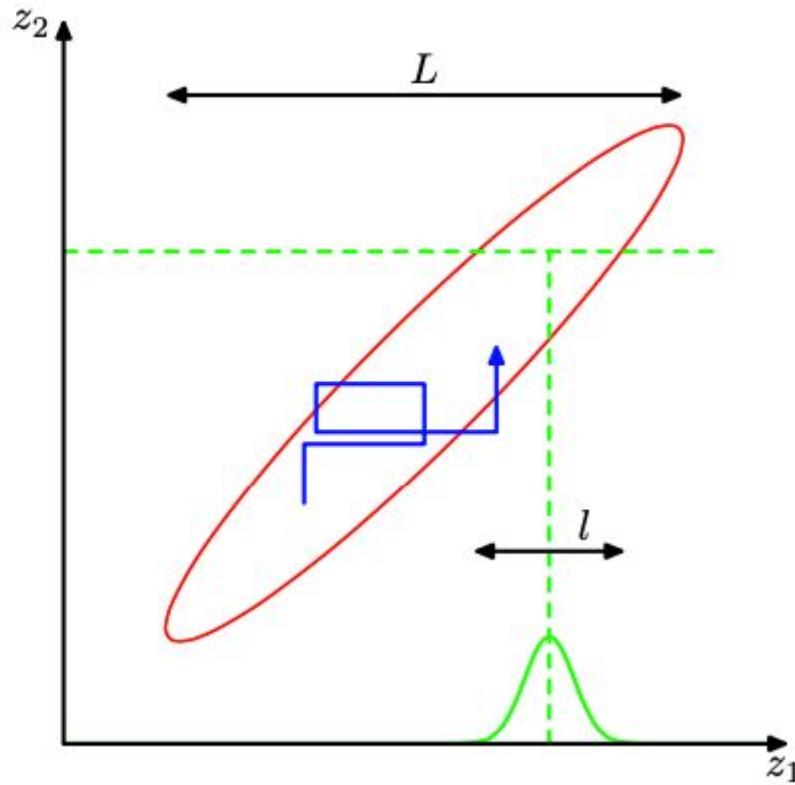
$$A(\mathbf{z}^*, \mathbf{z}) = \frac{p(\mathbf{z}^*) q_k(\mathbf{z} | \mathbf{z}^*)}{p(\mathbf{z}) q_k(\mathbf{z}^* | \mathbf{z})} = \frac{p(z_k^* | \mathbf{z}_{\setminus k}^*) p(\mathbf{z}_{\setminus k}^*) p(z_k | \mathbf{z}_{\setminus k}^*)}{p(z_k | \mathbf{z}_{\setminus k}) p(\mathbf{z}_{\setminus k}) p(z_k^* | \mathbf{z}_{\setminus k})} = 1 \quad (11.49)$$

Figure 11.11

Illustration of Gibbs sampling by alternate updates of two variables whose distribution is a correlated Gaussian. The step size is governed by the standard deviation of the conditional distribution (green curve), and is $O(l)$, leading to slow progress in the direction of elongation of the joint distribution (red ellipse). The number of steps needed to obtain an independent sample from the distribution is $O((L/l)^2)$.

initial goal $E_p(f(x))$

practically: keep one sample every L^1 samples



Gibbs Sampling

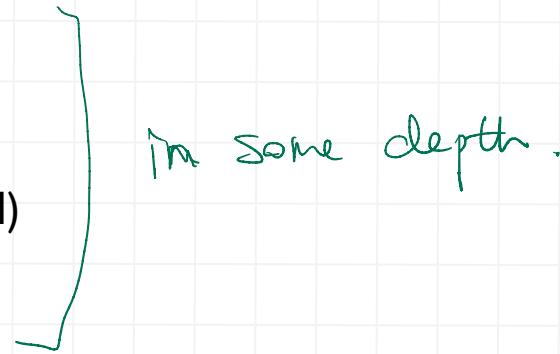
1. Initialize $\{z_i : i = 1, \dots, M\}$
2. For $\tau = 1, \dots, T$:
 - Sample $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$.
 - Sample $z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau+1)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$.
 - ⋮
 - Sample $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \dots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \dots, z_M^{(\tau)})$.
 - ⋮
 - Sample $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \dots, z_{M-1}^{(\tau+1)})$.

Sampling methods

Motivation, why? Sampling and ML

Basic sampling algorithms

- Sampling standard distributions from $U(0, 1)$
- Rejection sampling
- Importance sampling



Markov chain Monte Carlo (MCMC)

- Markov chains
- Metropolis-Hastings

Gibbs sampling