#### Statistics 360/601 – Modern Bayesian Theory

Alexander Volfovsky

Lecture 7 - Sept 20, 2018

# Monte Carlo

### Monte Carlo approximation

Want to compute

$$\mathsf{E}[ heta|y] = \int heta \mathsf{p}( heta|y) d heta$$

without worrying about actual integration...

- Let  $\theta$  be a parameter of interest
- Let y<sub>1</sub>,..., y<sub>n</sub> be numerical values of a sample from a distribution p(y<sub>1</sub>,..., y<sub>n</sub>|θ)
- Let θ<sup>1</sup>,..., θ<sup>S</sup> ∼ p(θ|y<sub>1</sub>,..., y<sub>n</sub>) be iid samples from the posterior.
- We estimate our posterior quantity of interest as  $1/S \sum g(\theta^i)$

#### Full on example

- Data: birthrates and education information for women.
- Question: are there any differences in the bithrates of women without and with bachelor's degrees?
- ▶ 111 women without bachelors gave birth to 217 children.
- ▶ 44 women with bachelors degrees gave birth to 66 children.
- Model: Poisson with a mean parameter θ<sub>1</sub> for women without a bachelors, and θ<sub>2</sub> with a bachelors.
- Prior:  $\theta_1, \theta_2 \sim gamma(a = 2, b = 1)$
- ▶ Posterior:  $p(\theta_1 | \sum_{i=1}^{111} Y_{i,1} = 217)$  is gamma(219, 11) and  $p(\theta_2 | \sum_{i=1}^{44} Y_{i,2} = 66)$  is gamma(68, 45).
- ▶ Want to compute  $p(\theta_1 > \theta_2 | \sum Y_{i1} = 217, \sum Y_{i2} = 66)$  and  $p(\tilde{Y}_1 > \tilde{Y}_2 | \sum Y_{i1} = 217, \sum Y_{i2} = 66)$ .

Magic answers to our questions:

$$p(\theta_1 > \theta_2 | \sum Y_{i1} = 217, \sum Y_{i2} = 66) = 0.97$$

$$p( ilde{Y}_1 > ilde{Y}_2 | \sum Y_{i1} = 217, \sum Y_{i2} = 66) = 0.48$$

As you might imagine these are some ... sums and integrals.

#### Monte Carlo to the rescue

```
> a<-2 ; b<-1
> sy1<-217 ; n1<-111
> sy2<-66 ; n2<-44
>
> theta1.mc<-rgamma(10000,a+sy1, b+n1)</pre>
> theta2.mc<-rgamma(10000,a+sy2, b+n2)</pre>
>
> y1.mc<-rpois(10000,theta1.mc)</pre>
> y2.mc<-rpois(10000,theta2.mc)</pre>
>
>
> mean(theta1.mc>theta2.mc)
[1] 0.9745
>
> mean(y1.mc>y2.mc)
[1] 0.4768
```

#### Buffon's needle

A crazy man wants to estimate  $\pi$  and all he has is an infinite carpet and a needle...



Wiki Images. Practical example: http://mste.illinois.edu/activity/buffon/

#### Buffon's needle

- ▶ Real world: throw a bunch of needles and estimate.
- Computer world: pretend to throw a bunch of needles and estimate.
- Sample  $(x, \theta)$  uniformly from  $[0, t/2] \times [0, \pi/2]$ .



#### Monte Carlo does not solve everything

# Monte Carlo failure



n

- Want to estimate  $I = \int_0^1 x^3 dx$ .
- Method 1:
  - 1. Sample  $u_1, \ldots, u_S \sim Unif(0, 1)$
  - 2. Estimate  $\hat{I} = \sum (u_i)^3 / S$
- Method 2:
  - 1. Sample  $b_1, \ldots, b_S \sim Beta(a, b)$
  - 2. Estimate  $\hat{I} = \sum ((b_i)^3/dbeta(b_i, a, b))/S$
- Which one is better?

(note that the optimal choice for importance sampling has the density  $g^*(x) = 4x^3$  – obviously  $x^3/g^* = 1/4$  is the answer)







zooming in



zooming in



So much better...  $g(x) = 3x^2$ 



#### Why is there a difference?



### Rejection sampling



#### Consistent parameters??

- ▶ Let the data come from *Y* ~Geometric(*p*).
- Recall that a geometric variable has pdf

$$\Pr(Y=k)=(1-p)^k p$$

and it captures a success on the k + 1 trial after k failures.

- The mean of a geometric is (1-p)/p.
- Let the model we think the data comes from be  $Poisson(\theta)$ .
- Let the prior be a Gamma $(\alpha, \beta)$ .
- We know the posterior is  $Gamma(\alpha + \sum y_i, \beta + n)$ .

# Consistent parameters??



#### Consistent parameters??



# Posterior Predictive Checks

Lets look at the data (from the book).



Data: twice as many women with 2 children as with 1. Posterior predictive: fewer women with 2 children than with 1.

- > t.mc <- t2.mc <-NULL</pre>
- > for(s in 1:10000) {
- > theta1<-rgamma(1,a+sum(y1), b+length(y1))</pre>
- > y1.mc<-rpois(length(y1),theta1)</pre>
- > t.mc <- c(t.mc,mean(y1.mc))</pre>
- > t2.mc<-c(t2.mc,sum(y1.mc==2)/sum(y1.mc==1))
  > }





Where t(y) is the ratio of 2's to 1's in a dataset.

• Lets look at the data (Gelman, Meng and Stern).

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.
- Goal: Estimate the density of a radioactive isotope in a cross-section of the brain.

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.
- Goal: Estimate the density of a radioactive isotope in a cross-section of the brain.
- 2-d image is estaimte from gamma-ray counts in a rigng of detectors around the head.

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.
- Goal: Estimate the density of a radioactive isotope in a cross-section of the brain.
- 2-d image is estaimte from gamma-ray counts in a rigng of detectors around the head.
- n bins of counts based on positions of detectors 6 million counts.

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.
- Goal: Estimate the density of a radioactive isotope in a cross-section of the brain.
- 2-d image is estaimte from gamma-ray counts in a rigng of detectors around the head.
- n bins of counts based on positions of detectors 6 million counts.
- Bin count  $y_i$  are modeled as independent Poisson $(\theta_i)$

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.
- Goal: Estimate the density of a radioactive isotope in a cross-section of the brain.
- 2-d image is estaimte from gamma-ray counts in a rigng of detectors around the head.
- n bins of counts based on positions of detectors 6 million counts.
- Bin count  $y_i$  are modeled as independent Poisson $(\theta_i)$
- → Θ = Ag + r where g is the unknown image, A is a known linear operator and r are known corrections.

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.
- Goal: Estimate the density of a radioactive isotope in a cross-section of the brain.
- 2-d image is estaimte from gamma-ray counts in a rigng of detectors around the head.
- n bins of counts based on positions of detectors 6 million counts.
- Bin count  $y_i$  are modeled as independent Poisson $(\theta_i)$
- → Θ = Ag + r where g is the unknown image, A is a known linear operator and r are known corrections.
- A, g and r are non-negative.

- Lets look at the data (Gelman, Meng and Stern).
- Gelman (1990,92) describe positron emission tomography experiment.
- Goal: Estimate the density of a radioactive isotope in a cross-section of the brain.
- 2-d image is estaimte from gamma-ray counts in a rigng of detectors around the head.
- n bins of counts based on positions of detectors 6 million counts.
- Bin count  $y_i$  are modeled as independent Poisson $(\theta_i)$
- → Θ = Ag + r where g is the unknown image, A is a known linear operator and r are known corrections.
- ► *A*, *g* and *r* are non-negative.
- This is an easy problem without the constraint.

 Poisson noise + model problems makes exact non-negative solutions impossible.

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

• Fitting to real data: y with n = 22464.

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Fitting to real data: y with n = 22464.
- The best-fit non-negative image ĝ was not an exact fit leading to the discrepancy between y and θ̂ to be X<sup>2</sup>(y; θ̂) ≈ 30,000.

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Fitting to real data: y with n = 22464.
- ► The best-fit non-negative image ĝ was not an exact fit leading to the discrepancy between y and θ to be X<sup>2</sup>(y; θ̂) ≈ 30,000.
- Reject the model!

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Fitting to real data: y with n = 22464.
- ► The best-fit non-negative image ĝ was not an exact fit leading to the discrepancy between y and θ to be X<sup>2</sup>(y; θ̂) ≈ 30,000.
- Reject the model!
- Possible failures:

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Fitting to real data: y with n = 22464.
- ► The best-fit non-negative image ĝ was not an exact fit leading to the discrepancy between y and θ to be X<sup>2</sup>(y; θ̂) ≈ 30,000.
- Reject the model!
- Possible failures:
  - Error in the specification of A, r

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Fitting to real data: y with n = 22464.
- The best-fit non-negative image ĝ was not an exact fit leading to the discrepancy between y and θ̂ to be X<sup>2</sup>(y; θ̂) ≈ 30,000.
- Reject the model!
- Possible failures:
  - Error in the specification of A, r
  - Lack of independence

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Fitting to real data: y with n = 22464.
- The best-fit non-negative image ĝ was not an exact fit leading to the discrepancy between y and θ̂ to be X<sup>2</sup>(y; θ̂) ≈ 30,000.
- Reject the model!
- Possible failures:
  - Error in the specification of A, r
  - Lack of independence
  - super-Poisson variance in the counts

- Poisson noise + model problems makes exact non-negative solutions impossible.
- Use an estimate  $\hat{g}$  and capture the discrepancy between y and  $\hat{\theta} = A\hat{g}_r$ .
- Consider  $\chi^2$  discrepancy.

$$X^2(y;\hat{ heta}) = \sum rac{(y_i - \hat{ heta}_i)^2}{\hat{ heta}_i}$$

- Fitting to real data: y with n = 22464.
- The best-fit non-negative image ĝ was not an exact fit leading to the discrepancy between y and θ̂ to be X<sup>2</sup>(y; θ̂) ≈ 30,000.
- Reject the model!
- Possible failures:
  - Error in the specification of A, r
  - Lack of independence
  - super-Poisson variance in the counts
  - Error from discretizing the continuous g.

• The "critically bad" level of  $X^2(y; \hat{\theta})$  is  $n + 2\sqrt{2n} \approx 23,000$ .

- The "critically bad" level of  $X^2(y; \hat{\theta})$  is  $n + 2\sqrt{2n} \approx 23,000$ .
- We reassess and change the model.

- The "critically bad" level of  $X^2(y; \hat{\theta})$  is  $n + 2\sqrt{2n} \approx 23,000$ .
- We reassess and change the model.
- We get a  $\chi^2$  discrepancy that is 22,000 or even 20,000.

- The "critically bad" level of  $X^2(y; \hat{\theta})$  is  $n + 2\sqrt{2n} \approx 23,000$ .
- We reassess and change the model.
- We get a  $\chi^2$  discrepancy that is 22,000 or even 20,000.
- Should we just accept this new model?

- The "critically bad" level of  $X^2(y; \hat{\theta})$  is  $n + 2\sqrt{2n} \approx 23,000$ .
- We reassess and change the model.
- We get a  $\chi^2$  discrepancy that is 22,000 or even 20,000.
- Should we just accept this new model?
- No.

- The "critically bad" level of  $X^2(y; \hat{\theta})$  is  $n + 2\sqrt{2n} \approx 23,000$ .
- We reassess and change the model.
- We get a  $\chi^2$  discrepancy that is 22,000 or even 20,000.
- Should we just accept this new model?
- ► No.
- Be skeptical!