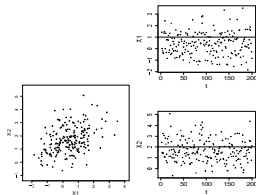


<div> <div>Monte Carlo Methods</div> <div> 732A90 Computational Statistics Krzysztof Bartoszek (krzysztof.bartoszek@liu.se)  5 11 2019 (P42) Department of Computer and Information Science Linköping University </div> </div>	<div> <div>What is the area of the unit circle?</div> <div> <pre>f.circArea&lt;-function(N){   m.xy&lt;-cbind(runif(N),runif(N))   4*sum(apply(m.xy,1,function(xy){xy[1]^2+xy[2]^2&lt;1}))/N }</pre> <div> </div> </div> </div>	<div> <div>Monte Carlo methods: outline</div> <div> <ul style="list-style-type: none"> <li>Monte Carlo methods are a class of computational algorithms that use repeated random sampling to compute their results.</li> <li>Monte Carlo methods for random number generation <ul style="list-style-type: none"> <li>Metropolis–Hastings algorithm</li> <li>Gibbs sampler</li> </ul> </li> <li>Monte Carlo methods for statistical inference <ul style="list-style-type: none"> <li>Estimate integrals (we already did!)</li> <li>Variance estimation</li> <li>Variance reduction: importance sampling, control variates</li> </ul> </li> </ul> </div> </div>	<div> <div>Markov Chain Monte Carlo</div> <div> <div>Previous lecture: Generate</div> <ul style="list-style-type: none"> <li>univariate distributions (inverse CDF, acceptance/rejection)</li> <li>multivariate normal</li> </ul> <div>but general multivariate distribution?</div> <div>MCMC</div> </div> </div>
<div> <div>Bayesian inference: Recap</div> <div> <p>A dataset <math>D</math> is obtained by sampling from a distribution <math>f(\cdot \theta)</math>. How to estimate <math>\theta</math>?</p> <ul style="list-style-type: none"> <li><b>Frequentists:</b> <math>\theta</math> is an unknown but fixed parameter, compose likelihood <math>L(D \theta)</math> and find <math>\hat{\theta}</math> that maximizes it.</li> <li><b>Bayesians:</b> <math>\theta</math> is a random variable with <b>prior</b> probability law <math>p(\theta)</math> before observing <math>D</math></li> <li>After observing <math>D</math>, Bayes' theorem gives</li> </ul> <math display="block">p(\theta D) = \frac{p(D \theta)p(\theta)}{p(D)} = \frac{p(D \theta)p(\theta)}{\int p(D \theta)p(\theta)d\theta}</math> </div> </div>	<div> <div>Bayesian inference: Recap</div> <div> <div> <math display="block">p(\theta D) = \frac{p(D \theta)p(\theta)}{p(D)} = \frac{p(D \theta)p(\theta)}{\int p(D \theta)p(\theta)d\theta}</math> </div> <div> <p>We know: <math>p(D \theta)</math> (the model), <math>p(\theta)</math> (the prior) We need: simulate from <math>p(\theta D)</math> (the posterior)</p> <ul style="list-style-type: none"> <li>General (multivariate) type distribution</li> <li>Integral can be impossible to compute</li> <li>MCMC solves this</li> <li>Not needed (given <math>D</math> it is constant)</li> </ul> </div> </div> </div>	<div> <div>Markov Chains: Recap</div> <div> <ul style="list-style-type: none"> <li>A Markov chain is a sequence <math>X_0, X_1, \dots</math> of random variables such that the distribution of the next value depends only on the current one (and parameters).</li> <li><math>P(X_{t+1} X_t)</math> is called a <b>transition kernel</b>. Assume it does not depend on <math>t</math> (<b>time homogeneous</b>).</li> <li>A Markov chain is <b>stationary</b>, with stationary distribution <math>\Phi</math>, if <math>\forall_k X_k \sim \Phi</math></li> <li>One shows (not trivial in general) that under <i>certain</i> conditions a Markov chain will converge to the stationary distribution in the limit.</li> </ul> </div> </div>	<div> <div>Markov Chains: Example</div> <div> <math display="block">X(t+1) = e^{-1}X(t) + \epsilon, \epsilon \sim \mathcal{N}(0, \frac{1}{2} \cdot (1 - e^{-2}))</math> <div>Discard first <math>K-1</math> samples: <b>burn-in</b> period</div> </div> </div>
<div> <div>MCMC: Example</div> <div> <p><b>Linear regression</b> with residual normally/student/etc. distributed</p> <math display="block">Y = \beta X + \epsilon</math> <p>How to find credible interval for <math>\beta</math> if we know <math>\text{Var}[\epsilon] = \sigma^2</math>?</p> <ul style="list-style-type: none"> <li>Obtain <math>P(\beta Y, X)</math> by drawing from <math>P(Y X, \beta)P(\beta)</math> in a <b>clever way</b>.</li> <li>The prior ?</li> <li>Use the MCMC sample to obtain quantiles.</li> </ul> <p>Normal residual: analytical solution</p> </div> </div>	<div> <div>Metropolis–Hastings algorithm</div> <div> <p>We have</p> <ul style="list-style-type: none"> <li>A PDF <math>\pi(x)</math> that we want to sample from.</li> <li>A <b>proposal distribution</b> <math>q(\cdot X_t)</math> that has a <b>regular</b> form w.r.t. to <math>\pi(\cdot)</math> E.g. <math>q(\cdot X_t)</math> is normal with mean <math>X_t</math> and given variance</li> <li><b>Regular</b> form: suffices that the proposal has the same support as <math>\pi</math>.</li> </ul> </div> </div>	<div> <div>Metropolis–Hastings Sampler</div> <div> <math display="block">\alpha(X_t, Y) = \min \left\{ 1, \frac{\pi(Y)q(X_t Y)}{\pi(X_t)q(Y X_t)} \right\}</math> <pre> 1: Initialize chain to <math>X_0</math>, <math>t = 0</math> 2: <b>while</b> <math>t &lt; t_{\max}</math> <b>do</b> 3:   Generate a candidate point <math>Y \sim q(\cdot X_t)</math> 4:   Generate <math>U \sim \text{Unif}(0, 1)</math> 5:   <b>if</b> <math>U &lt; \alpha(X_t, Y)</math> <b>then</b> 6:     <math>X_{t+1} = Y</math> 7:   <b>else</b> 8:     <math>X_{t+1} = X_t</math> 9:   <b>end if</b> 10:  <math>t = t + 1</math> 11: <b>end while</b> </pre> </div> </div>	<div> <div>Metropolis–Hastings Sampler: Properties</div> <div> <ul style="list-style-type: none"> <li>Informally: “The chain <math>(X_t)_{t=0}^\infty</math> will converge to <math>\pi(\cdot)</math>.”</li> <li>The chain might not move sometimes.</li> <li>The values of the chain are dependent.</li> <li>If <math>q(X_t Y) = q(Y X_t)</math> (i.e. symmetric proposal) we get <b>Random-walk Monte Carlo</b>: <math display="block">\alpha(X_t, Y) = \min \left\{ 1, \frac{\pi(Y)}{\pi(X_t)} \right\}</math> </li> </ul> </div> </div>
<div> <div>Choice of proposal distribution</div> <div> <ul style="list-style-type: none"> <li>In Random–Walk Monte Carlo</li> </ul> <p>If <math>\pi(Y) \geq \pi(X)</math>, the chain moves to the next point, otherwise only with some probability.</p> </div> </div>	<div> <div>Choice of proposal dist.: target: <math>\pi(\cdot) = \mathcal{N}(0, 1)</math></div> <div> 732A90.ComputationalStatisticsVT2019.Lecture04codeSlide14.R </div> </div>	<div> <div>Choice of proposal distribution</div> <div> <p><math>q</math> normal with sd: <math>\text{prop} = 0.5, 0.1</math> and <math>20</math></p> </div> </div>	<div> <div>Gibbs sampler: alternative to Metropolis–Hastings</div> <div> <p>We want to generate from a distribution on <math>\mathbb{R}^d</math>.</p> <pre> 1: Initialize chain to <math>X_0 = (X_{0,1}, \dots, X_{0,d})</math>, <math>t = 0</math> 2: <b>while</b> <math>t &lt; t_{\max}</math> <b>do</b> 3:   <b>for</b> <math>i = 1, \dots, d</math> <b>do</b> 4:     Generate <math display="block">X_{t+1,i} \sim f(\cdot X_{t+1,1}, \dots, X_{t+1,i-1}, X_{t+1,i+1}, \dots, X_{t,d})</math> 5:   <b>end for</b> 6:   <math>t = t + 1</math> 7: <b>end while</b> </pre> </div> </div>

Gibbs sampler	Gibbs sampler: target: $d$ -dim $\mathcal{N}(\mu, \Sigma)$	Gibbs sampler: Example (code: see R scripts)	Convergence monitoring
<ul style="list-style-type: none"><li>At each iteration inside the <code>for</code> loop univariate random numbers are generated.</li><li>Only one element is updated.</li><li><b>WE NEED TO KNOW THE CONDITIONAL MARGINAL DISTRIBUTIONS.</b></li><li>Convergence may be slow.</li><li>Can be useful in high dimensions (i.e. proposal density may be difficult to find in another way).</li></ul>	732A90.ComputationalStatisticsVT2019.Lecture04codeSlide18.R	<p>Generate from</p> $\mathcal{N}\left(\frac{1}{2}\mathbf{2}^T, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$ 	<ul style="list-style-type: none"><li>When should we stop the chain? When are we (nearly) at the stationary distribution?</li><li>Typically such a sample is generated to make further inference.</li></ul>

Convergence monitoring: Gelman–Rubin method	Gibbs sampler	MC for inference	MC for inference
<p>We want to estimate <math>\psi(\theta)</math>.</p> <ul style="list-style-type: none"> <li>Generate <math>k</math> sequences of length <math>n</math> with different starting points.</li> <li>Compute between- and within- sequence variances:</li> </ul> $B = \frac{n}{k-1} \sum_{i=1}^k (\bar{y}_i - \bar{\bar{y}})^2 \quad W = \sum_{i=1}^k \frac{s_i^2}{k} \quad s_i^2 = \sum_{j=1}^n \frac{(\bar{y}_{ij} - \bar{y}_i)^2}{n-1}$ <ul style="list-style-type: none"> <li>Overall variance estimate: <math>\text{Var}[\psi] \approx \frac{n-1}{n} W + \frac{1}{n} B</math></li> <li>Gelman–Rubin factor:</li> </ul> $\sqrt{R} = \sqrt{\frac{\text{Var}[\psi]}{W}}$ <ul style="list-style-type: none"> <li>Values much larger than 1 indicate lack of convergence</li> <li>See <code>?coda::gelman.diag</code></li> </ul>	<pre>library(coda) f1&lt;-mcmc.list();f2&lt;-mcmc.list();n&lt;-100;k&lt;-20 X1&lt;-matrix(rnorm(n*k),ncol=k,nrow=n) X2&lt;-X1+(apply(X1,2,cumsum)*matrix(rep(1:n,k),ncol= k)^2)) for (i in 1:k){f1[[i]]&lt;-as.mcmc(X1[,i]);f2[[i]]&lt;-as .mcmc(X2[,i])} print(gelman.diag(f1)) # Potential scale reduction factors: # Point est. Upper C.I. # 1.  0.999 1.01  print(gelman.diag(f2)) # Potential scale reduction factors: # Point est. Upper C.I. # 1.  1.82 2.38</pre>	<ul style="list-style-type: none"> <li>Estimation of a definite integral</li> </ul> $\theta = \int_D f(x)dx \quad \left( \text{recall } \pi = \int_D 1dx \right)$ <ul style="list-style-type: none"> <li>Decompose into:</li> </ul> $f(x) = g(x)p(x) \quad \text{where } \int_D p(x)dx = 1$ <ul style="list-style-type: none"> <li>Then, if <math>X \sim p(\cdot)</math></li> </ul> $\theta = \mathbb{E}[g(X)] = \int_D g(x)p(x)dx$ <ul style="list-style-type: none"> <li></li> </ul> $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n g(x_i), \quad \forall x_i \sim p(\cdot)$	<ul style="list-style-type: none"> <li>Decomposition is not unique, some will be better (lower variance) others worse. <math>p(x) \propto f(x)</math>: minimal</li> <li>Can we easily generate from <math>p(\cdot)</math>?</li> <li>Bayesian inference: use MCMC samples from <math>p(\theta D)</math> to obtain a point estimator</li> </ul> $\theta^* = \int \theta p(\theta D) \approx \frac{1}{n} \sum_{i=1}^n \theta_i$ <ul style="list-style-type: none"> <li><math>\hat{\theta}</math> depends on <math>n</math> and <math>g(X)</math>, how variable will it be?</li> </ul> $\widehat{\text{Var}}[\hat{\theta}] = \frac{1}{n(n-1)} \sum_{i=1}^n \left( g(x_i) - \overline{g(x)} \right)^2$ <ul style="list-style-type: none"> <li>MCMC: estimator biases as chain correlated, use longer chain and batch mean instead of <math>x_i</math>.</li> </ul>

Summary
<ul style="list-style-type: none"> <li>Generating data from a general multivariate distribution</li> <li>Markov Chain Monte Carlo: Metropolis–Hastings algorithm, Gibbs sampling</li> <li>Convergence: Gelman–Rubin method</li> <li>Estimation of integral</li> </ul>