

6/19/19

W2a-1

MCMC Sampling III

Today we will briefly follow-up on the evidence discussion and then talk mostly about Hamiltonian Monte Carlo (HMC) and its implementation in pymc3.

- First: Computational possibilities for evidence \rightarrow (T2a-8)
 - Examples of Information Criteria. \rightarrow briefly given time.
 - For further information BDA3 (Bayesian Data Analysis 3 by Gelman et al.) is a good reference.

AIC: Akaike Information Criteria

- Essentially Frequentist \rightarrow relies on likelihood
- Want to minimize: $-2 \log p(D | \hat{\theta}_{MLE}) + 2k$
 - \uparrow Likelihood evaluated at maximum likelihood value of parameters (peak of likelihood)
 - \nwarrow # of free parameters
- So it has ingredients: improved likelihood balanced by penalty for ^{extra} parameters. No priors.
- Not regarded well by Bayesians

BIC: Bayesian Information Criteria

- Gaussian approximation to Bayesian evidence in limit of large amount of data
- minimize BIC = $-2 \log p(D | \hat{\theta}_{MLE}) + k \ln N$
 - \nwarrow # fitted pts \nwarrow # data pts
- Assumes Occam razor penalty is negligible.

DIC: Deviance Information Criteria

- replace $\hat{\theta}_{MLE}$ by $\hat{\theta}_{Bayes}$ \leftarrow maximum of posterior ^{average θ over posterior}
- use effective # of parameters $p_{DIC} = 2 \log p(D | \hat{\theta}_{Bayes}) - E[\log p(D | \theta)]$
- So k replaced by data-based bias correction: $DIC \equiv -2 \log p(D | \hat{\theta}_{Bayes}) + 2 p_{DIC}$

widely applicable

WAIC: Favored by BDA-3 as more fully Bayesian

Sampling: $2 \sum_{i=1}^N \left(\log \frac{1}{S} \sum_{s=1}^S p(y_i | \theta^s) \right) - \frac{1}{S} \sum_{s=1}^S \log p(D_i | \theta^s)$

averages over posterior distribution

6/19/19

W2A-2

Visualization of MCMC Sampling Revisited

- We return to the excellent set of interactive demos by Chi Feng at <https://chi-feng.github.io/mcmc-demo/> and their adaptation by Richard McElreath at <http://eleventh.org/blog/2017/11/28/build-a-better-markov-chain/>.
- This blog piece very strongly advocates abandoning Metropolis-Hastings sampling in favor of Hamiltonian Monte Carlo (HMC), our topic for the rest of today.
 - First recall the Random Walk MH
 - 1) make a random proposal for new parameter values (step in parameter space \rightarrow arrow)
 - 2) Accept or reject proposal based on a Metropolis criterion.
 - Diffusion (random walk) so not efficient in exploring the space and needs special tuning to avoid too many rejections.
 - Donut shape is common in higher dimensions and it is hard to explore.
- "Better living through physics" \Rightarrow HMC simulation
 - The description is that we map our parameter vector to a particle in an n -dimensional space. The surface is an m -dimensional bowl with shape given by minus-log target distribution (e.g. posterior). See log Gaussian.
 - Treat as frictionless. "Flick" particle in random direction, so it flows across the bowl.
 - See the simulation: little gray arrow is flick. After travels some distance, decide whether to accept. Most are within high probability region, so high percentage accepted.
 - Chains can get far from starting point \Rightarrow efficient exploration of full shape.
 - More calculation but fewer samples \Rightarrow trade off wins.
 - Check the donut \Rightarrow looks very good!

6/19/19

- There is a further improvement called NUTS - no U-turn sampler. The idea is to address the problem that HMC needs to be told how many steps to take before another flick,

ask class:
What is the
drawback of
paths in
simulation?

- too few steps \Rightarrow samples too similar
- too many steps \Rightarrow also too similar
- tuning by hand is difficult with complex distribution,

• Solution is NUTS.

- adaptively finds a good number of steps
- simulates in both directions to figure out when the path turns around (U-turn) and stops it.
- Other adaptive features - see documentation.

- Note that NUTS still has trouble with multimodal targets \Rightarrow can explore each high probability area, but has trouble going between.

6/9/19

look at
these pictures

→ Hamiltonian Monte Carlo Explained by Alex Rogozhnikov
and MCMC Using Hamiltonian Dynamics by Radford Neal.

The basic idea is to translate a pdf for the distribution desired to a potential energy function and then add a momentum variable — fictitious! In Markov chain at each iteration resamples the momentum, creates a proposal using Hamiltonian dynamics, and then does a Metropolis update.

Recall Hamiltonian dynamics, now applied to d -dimensional position vector q and a d -dimensional momentum vector p
 \Rightarrow 2D phase space Hamiltonian is $H(q, p)$.

Equations of motion describe time evolution:

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i} \quad \text{and} \quad \frac{dp_i}{dt} = -\frac{\partial H}{\partial q_i} \quad i=1, \dots, d$$

so these map states at t to states at $t+s$.

We take the form of H to be $H(q, p) = U(q) + K(p)$

- potential energy $U(q)$ is minus the log probability density of the distribution for q we seek to sample.
- $K(p)$ is kinetic energy $K(p) = p^T M^{-1} p / 2$ and M is symmetric, positive "mass matrix" typically diagonal and den m_{ii}
 - This is minus the log probability density (plus constant) of zero mean Gaussian with covariance matrix M .

What are we going to do with this? We consider a canonical distribution

$$P(q, p) = \frac{1}{Z} e^{-H(q, p)/T} = \frac{1}{Z} e^{-U(q)/T} e^{-K(p)/T}$$

so q and p are independent. We are interested in q : p is fake to make things work. Usually $U(q)$ is a posterior: $-\log[p(q|D)p(q)]$

\downarrow
0, for example.

6/19/19

W2a-5

Two steps of HMC algorithm:

- 1) New values for the momentum variables are randomly drawn from their Gaussian distribution, independent of current position values.
 - This means p_i has mean zero and variance M_i ; if M diagonal,
 - q isn't changed, p is from correct conditional distribution given q , so canonical joint distribution invariant

- 2) Proposal from Hamiltonian dynamics for new state,

Simulate from (q, p) with L steps of size $\epsilon \Rightarrow$ parameters.

At end, ^{the momenta are flipped sign} a new proposed state (q^*, p^*) accepted with probability
$$\min\left[1, e^{-H(q^*, p^*) + H(q, p)}\right] = \min\left[1, e^{-U(q^*) + U(q) - K(p^*) + K(p)}\right]$$

- momentum flip makes proposal symmetrical, but not done in practice
- So probability density for (q, p) (almost) unchanged because energy is conserved, but in terms of q we get a very different probability density.

- You can show that HMC leaves canonical distribution invariant ^{because} a detailed balance holds, which is what we need. It will also be ergodic \Rightarrow doesn't get stuck in subset of state space but samples all.
 - Some subtlety about periodic orbit but ok,

Essential Features:

- Reversibility needed so that desired distribution is invariant,
- Conservation of Hamiltonian (which is the energy)
- Volume preservation — preserves volume in (q, p) space — This is Liouville's Theorem, (If we take a cluster of points and follow them, the volume occupied is unchanged.)
 \Rightarrow This is critical because a change in volume would mean we would have to make a nontrivial adjust to the proposal (the normalization Z would change.)

6/19/19

W2a-5

These requirements are satisfied by exact equation, but we are approximating differential equations. This requires a symplectic integration.

Ordinary Runge-Kutta-type ODE solver won't work. We need something like leapfrog:

$$p_i(t+\epsilon/2) = p_i(t) - (\epsilon/2) \frac{\partial U}{\partial q_i}(q(t)) \quad \leftarrow \text{half step}$$

$$q_i(t+\epsilon) = q_i(t) + \epsilon p_i(t+\epsilon/2) / m_i \quad \leftarrow \text{use intermediate } p$$

$$p_i(t+\epsilon) = p_i(t+\epsilon/2) - (\epsilon/2) \frac{\partial U}{\partial q_i}(q(t+\epsilon)) \quad \leftarrow \text{other half step}$$

using p at t or $t+\epsilon$ wouldn't be symmetric.

- See pictures in Fig 1 of Neal.

• Benefits of HMC seen in Figures 3-6

- Many tuning features are automatic with pymc3, so let's dive in and take a look.

Edges not
fine reversal
marked

6/19/19

Wade-7

Lightning intro to PyMC3

- Plan for this afternoon: your job is to go through some of the notebooks and make changes to see how it works. Later we'll have a notebook for you to step through the lighthouse problem.
- Bayes2019/Topics/mcmc-sampling/PyMC3/

• Start with PyMC3_intro.ipynb (A)

• Then do PyMC3_docs-getting-started.ipynb (B)

The posterior for

(A) Initial problem is sampling to find μ , the mean of a distribution, given data generated according to a normal distribution (actual μ is 200).

• Things to try: changing the true μ , sampling sigma as well.

• The standard deviation is initially fixed at 1.

• In the code one specifies priors, then the likelihood. By Bayes' Theorem this predicts the posterior for μ :

$$p(\mu|D, \sigma) \propto p(D|\mu, \sigma) p(\mu|\mu_0, \sigma_0)$$

distribution to sample

likelihood
specify last

$\sim N(0, 1)$
specify first

D → observed

$$p(D|\mu, \sigma) \sim N(\mu, \sigma)$$

"Cores"
is # of
chains
to run in parallel