Scientific Computing for ML and Supercloud

Jonathan Huggins

Scientific Computing for Machine Learning¹

- Challenges
 - Algorithms are stochastic
 - Correct algorithms can perform badly
 - What counts as good performance?
 - ML algorithms can be robust to bugs!
 - Fails on large datasets, in high dimensions, etc.
 - Unforeseen interactions between ML algorithms

Scientific Computing for Machine Learning¹

- Best practices
 - modular code
 - separate logically distinctive components
 - e.g., optimizer + function being optimized
 - e.g., sampler + model
 - testing
 - unit tests: correctness of small piece of code
 - integration tests: correctness of system

Example: Gibbs sampler for a mixture of Gaussians model¹

 $\pi \sim \text{Dirichlet}(\alpha)$ $\sigma_{\mu}^2 \sim \text{InverseGamma}(a_{\mu}, b_{\mu})$ $\sigma_n^2 \sim \text{InverseGamma}(a_n, b_n)$ $z_i \mid \pi \sim \text{Multinomial}(\pi)$ $\mu_{kj} \mid \sigma_{\mu}^2 \sim \text{Normal}(0, \sigma_{\mu}^2)$ $x_{ij} | z_i, \mu_{z_i,j}, \sigma_n^2 \sim \operatorname{Normal}(\mu_{z_i,j}, \sigma_n^2)$

Gibbs sampling

- Model is p(x, y, z)
- Start with state (x, y, z), then

 $x' \sim p(. | y, z)$ $y' \sim p(. | x', z)$ $z' \sim p(. | x', y')$

- Set (x, y, z) ← (x', y', z')
- Repeat

Writing MCMC code

- For MCMC, split code into
 - Distributions
 - Model
 - Sampler
- Modularity makes testing and reuse easier

<show MoG code>

 $\pi \sim \text{Dirichlet}(\alpha)$ $\sigma_{\mu}^2 \sim \text{InverseGamma}(a_{\mu}, b_{\mu})$ $\sigma_n^2 \sim \text{InverseGamma}(a_n, b_n)$ $z_i \mid \pi \sim \text{Multinomial}(\pi)$ $\mu_{kj} \mid \sigma_{\mu}^2 \sim \text{Normal}(0, \sigma_{\mu}^2)$ $x_{ij} | z_i, \mu_{z_i,j}, \sigma_n^2 \sim \text{Normal}(\mu_{z_i,j}, \sigma_n^2)$

Unit Testing

- Ideally, write the tests before writing the code
- Modular code makes writing unit tests easier
- Write and test "naive" implementation before adding an optimized implementation
- Trick for testing conditional probabilities: check that

p(x | z)/p(x' | z) = p(x, z)/p(x', z)

holds for random x, x', z (normalization terms will cancel)

class Model:

. . .

```
def joint_log_p(self, state, X):
return DirichletDistribution(self.alpha * np.ones(self.K)).log_p(state.pi) + \
MultinomialDistribution.from_probabilities(state.pi).log_p(state.z).sum() + \
self.sigma_sq_mu_prior.log_p(state.sigma_sq_mu) + \
self.sigma_sq_n_prior.log_p(state.sigma_sq_n) + \
GaussianDistribution(0., state.sigma_sq_mu).log_p(state.mu).sum() + \
GaussianDistribution(state.mu[state.z, :], state.sigma_sq_n).log_p(X).sum()
```

Unit Testing in Python

- In python, nose makes testing very easy
 - name test files test_* and test functions test_*
- also very useful: numpy.testing
 - when a test fails, provides useful debugging info

Integration Testing

- Geweke test:
 - 1. Sample iid from full model (data and latents)
 - 2. Iterate between (a) sampling new latents using MCMC algorithm and (b) resampling data from forward model
 - 3. Compare samples from steps 1 + 2 using P-P plot:
 - for CDFs F and G, plot (F(z), G(z)) for z in (-inf, inf)
- Geweke test can fail because of bug or poor mixing







Using Supercloud

- Supercloud is a high-performance cluster available through Lincoln Labs
- ssh <username>@txe1-login.mit.edu
- LLsub <job-script> -o <log-file> -J <job-name>
- See <u>http://groups.csail.mit.edu/broderick/wiki/</u> index.php?title=Supercloud
- Please add whatever you learn to the wiki!

Future topics...

- Using Cython (easy integration of C code with Python)
- Gelman–Rubin diagnostics for debugging MCMC
- Coding for GPUs (if anyone knows how)
- Autodiff