

PROBABILISTIC PROGRAMMING

MODERN TOOLS FOR PROBABILISTIC MODELING AND INFERENCE

Christian Findenig, Thomas Wedenig

Graz University of Technology, Austria
Institute of Theoretical Computer Science

1. Motivation
2. Markov-Chain Monte Carlo
3. MCMC - Demo
4. Variational Inference
5. Variational Inference Demo
6. Practical considerations
7. Conclusions

MOTIVATION

Example

- Suppose you have a **coin**



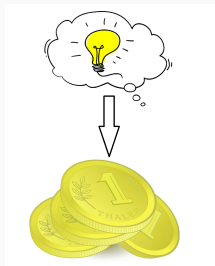
- You are unsure about its **fairness**
 - $z \in [0, 1]$
- You collect **data** by tossing it
 - $\mathbf{x} = [1, 0, 0, \dots, 1]$

Frequentist Approach

- **Maximum Likelihood!**
- Likelihood $p(x_i | z) = \text{Bernoulli}(z)$
- Tosses are i.i.d.:

$$p(\mathbf{x} | z) = \prod_{i=1}^n p(x_i | z)$$

- Find z^* that maximizes $p(\mathbf{x} | z)$
 - $z^* = \text{np.mean}(\mathbf{x})$

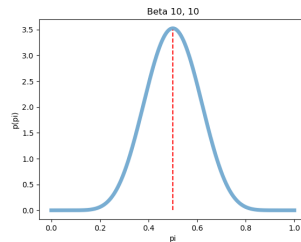


Prior

$$p(z) = \text{Beta}(10, 10)$$

Likelihood

$$p(\mathbf{x} | z) = \prod_{i=1}^n p(x_i | z)$$



Posterior

$$p(z | \mathbf{x}) = \frac{p(\mathbf{x} | z)p(z)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | z)p(z)}{\int p(\mathbf{x} | z')p(z') dz'}$$

How do we compute this?

Let's do the math!

$$\begin{aligned}
 p(z \mid \mathbf{x}) &= \frac{p(\mathbf{x} \mid z)p(z)}{\int p(\mathbf{x} \mid z')p(z') dz'} \\
 &= \frac{\prod_{i=1}^n z^{x_i} (1-z)^{1-x_i} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} z^{a-1} (1-z)^{b-1}}{\int \prod_{i=1}^n z'^{x_i} (1-z')^{1-x_i} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} z'^{a-1} (1-z')^{b-1} dz'} \\
 &= \frac{\prod_{i=1}^n z^{x_i} (1-z)^{1-x_i} z^{a-1} (1-z)^{b-1}}{\int \prod_{i=1}^n z'^{x_i} (1-z')^{1-x_i} z'^{a-1} (1-z')^{b-1} dz'} \\
 &= \frac{\prod_{i=1}^n z^{x_i} (1-z)^{1-x_i} z^{a-1} (1-z)^{b-1}}{\Gamma(a + \sum_{i=1}^n x_i) \Gamma(b + \sum_{i=1}^n (1-x_i)) \Gamma^{-1}(a+b+n)} \\
 &= \text{Beta} \left(a + \sum_{i=1}^n x_i, b + n - \sum_{i=1}^n x_i \right)
 \end{aligned}$$

- We were able to **analytically** solve the integral $p(\mathbf{x})$!
- Does this work for *any* choice for prior and likelihood?
 - Unfortunately, no! 😞

In general, we're not so lucky

- Modeling should be **flexible**
 - \implies **inference becomes difficult**
- Let's use Probabilistic Programming

Abstraction

Computer Science is all about abstraction.

- People want to **specify a model**
 - e.g., prior and likelihood
- Inference is done **by a tool**

(Informal) Definition

Probabilistic programming is about
doing statistics using the tools of computer science.

Tooling

Tooling is important.

Hypothesis

Deep Learning Revolution would have been
impossible without auto-diff.

- Flexible and rich **model specification**
- Should encourage **prototyping** and **iterative model development**
- \Rightarrow Let's use a **programming language**!

```
import pymc as pm

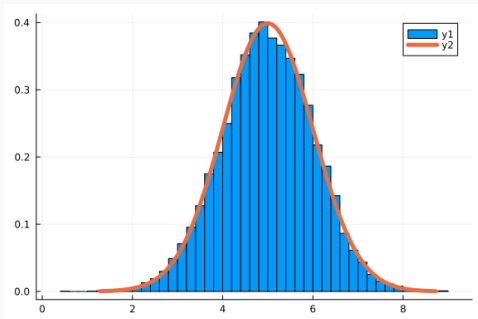
model = pm.Model()
x = [1, 1, 1, 0, 1, 1, 1, 0, 1]
with model:
    # Prior p(z)
    z = pm.Beta("z", alpha=10, beta=10)

    # Likelihood p(x | z)
    x_likelihood = pm.Bernoulli("x_likelihood", p=z, observed=x)
```

- Build on top of **auto-diff frameworks**
 - PyTorch, JAX, Tensorflow, ...
- Support for
 - Deep Generative Modeling
 - Gaussian/Dirichlet Processes
 - Discrete Latent Variables
 - ...

- **Markov Chain Monte Carlo**
 - Generate samples from the true posterior
- **Variational Inference**
 - Approximate the posterior with a tractable distribution

MARKOV-CHAIN MONTE CARLO



Metropolis Hasting

Prior, Likelihood & Posterior

$$p(z) = \text{Beta}(10, 10)$$

$$p(\mathbf{x} | z) = \prod_{i=1}^n p(x_i | z)$$

$$p(z | \mathbf{x}) \propto p(\mathbf{x} | z) \cdot p(z)$$

The algorithm

1. Propose z

$$z_{i+1} \sim g(z_{i+1} | z_i)$$

2. Ratio

$$R = \frac{p(z_{i+1} | \mathbf{x})}{p(z_i | \mathbf{x})}$$

3. z_{i+1} with *Bernoulli*($\min(R, 1)$), else z_i

Ratio within posterior

$$R = \frac{p(z_{i+1} | \mathbf{x})}{p(z_i | \mathbf{x})} = \frac{p(\mathbf{x} | z_{i+1}) \cdot p(z_{i+1})}{p(\mathbf{x} | z_i) \cdot p(z_i)}$$

Prior, Likelihood & Posterior

$$p(z) = \text{Beta}(10, 10)$$

$$p(\mathbf{x} | z) = \prod_{i=1}^n p(x_i | z)$$

$$p(z | \mathbf{x}) \propto p(\mathbf{x} | z) \cdot p(z)$$

Rejection rate

- Hamiltonian MC & NUTS sampling
- Gradient needed
- + No rejection

Metropolis Hasting

1. Propose z
2. Calculate acceptance Ratio R
3. Accept?

MCMC - DEMO

VARIATIONAL INFERENCE

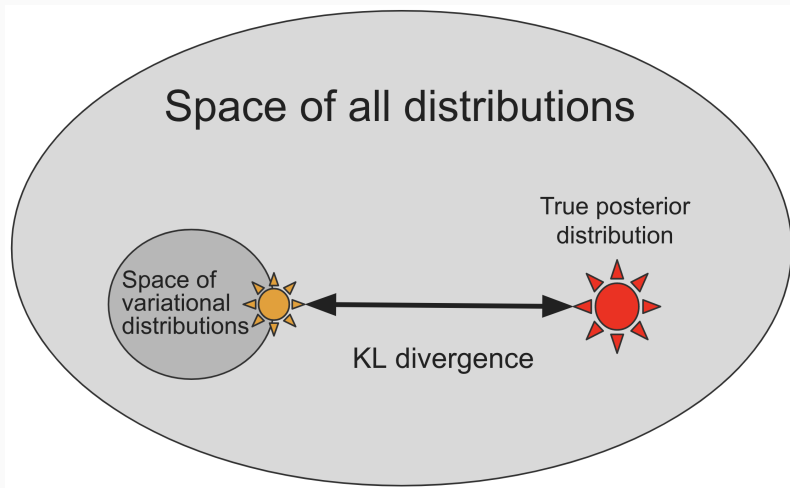
Problem

- Computing the posterior is intractable

$$p_{\theta}(\mathbf{z} \mid \mathbf{x}) = \frac{p_{\theta}(\mathbf{x} \mid \mathbf{z})p(\mathbf{z})}{\int p_{\theta}(\mathbf{x} \mid \mathbf{z}')p(\mathbf{z}') d\mathbf{z}'}$$

Variational Inference

- Consider a fixed $p_{\theta}(\mathbf{x} \mid \mathbf{z})$
- Pick a *variational* distribution $q_{\phi}(\mathbf{z})$
- Find ϕ^* such that $q_{\phi^*}(\mathbf{z})$ is close to $p_{\theta}(\mathbf{z} \mid \mathbf{x})$



https://pyro.ai/examples/intro_long.html

- When fitting q_ϕ , we would like to minimize $D_{KL}(q_\phi(\mathbf{z}) \parallel p_\theta(\mathbf{z} \mid \mathbf{x}))$
- But we can't compute the posterior ...
- Recall that

$$\begin{aligned}\text{ELBO} &= \mathbb{E}_{q_\phi(\mathbf{z})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z})] \\ &= \dots \\ &= \underbrace{\log p_\theta(\mathbf{x})}_{\text{constant w.r.t. } \phi} - D_{KL}(q_\phi(\mathbf{z}) \parallel p_\theta(\mathbf{z} \mid \mathbf{x}))\end{aligned}$$

Maximizing the **ELBO** \Leftrightarrow Minimizing D_{KL} to the true posterior

ELBO

$$\text{ELBO} = \log p_\theta(\mathbf{x}) - D_{KL}(q_\phi(\mathbf{z}) \parallel p_\theta(\mathbf{z} \mid \mathbf{x}))$$

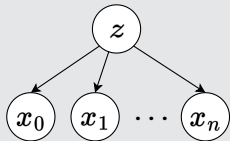
- Fit $p_\theta(\mathbf{x}, \mathbf{z})$ using maximum likelihood learning
- \Rightarrow We wish to maximize $p_\theta(\mathbf{x})$

Maximizing the **ELBO** encourages this

- We optimize p_θ and q_ϕ simultaneously
- \Rightarrow Compute $\nabla_{\theta, \phi} \text{ELBO}$

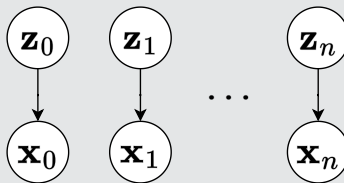
VARIATIONAL INFERENCE DEMO

Coin Example



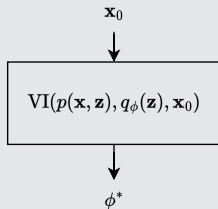
- One latent variable explains all coin flips
- We find $q_{\phi}(z \mid x_0, \dots, x_n)$

Image Example



- z_i are latent explanations of an image x_i
- Find posterior estimate $q_{\phi}(z_i \mid x_i)$ for each image

Classical VI

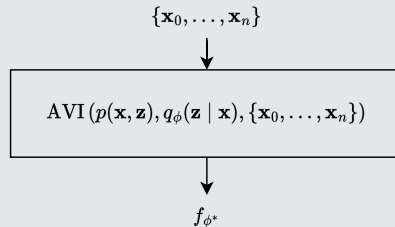


$$\phi^* = (\mu^*, \sigma^*)$$

$$q_{\phi^*}(\mathbf{z}) = \mathcal{N}(\mu^*, (\sigma^*)^2 I)$$

- $q_{\phi^*}(\mathbf{z})$ is the posterior for \mathbf{x}_0
- If we are given \mathbf{x}' , we need to re-run VI!

Amortized VI



$$f_{\phi^*}(\mathbf{x}) = (\mu^*, \sigma^*)$$

$$q_{\phi^*}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mu^*, (\sigma^*)^2 I)$$

- ϕ^* refer to **NN parameters**
- **Function** from any \mathbf{x} to $q_{\phi^*}(\mathbf{z} | \mathbf{x})$

PRACTICAL CONSIDERATIONS

MCMC

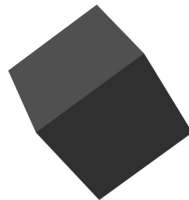
- 😍 Very **general**
- 😍 Exact in the limit (infinite time)
- 😞 Computationally **very expensive**
- 😞 Convergence **hard to diagnose**

Variational Inference

- 😍 Inference replaced by **optimization**
- 😍 More **efficient**
- 😞 How to **choose** $q_{\phi}(z)$?
- 😞 No **accuracy-computation tradeoff**



Edward



CONCLUSIONS

Probabilistic Programming is a set of tools for Bayesian Modeling & Inference

- Declarative, rich modeling system
- Abstracts away inference routines
- Works in tandem with techniques from **Deep Learning**

APPENDIX

ELBO GRADIENT ESTIMATION

ELBO

$$\text{ELBO} = \mathbb{E}_{q_{\phi}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})]$$

ELBO Gradient Estimation w.r.t. θ

$$\begin{aligned}\nabla_{\theta} \text{ELBO} &= \nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z})] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z})] \\ &\approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{z}_i) \quad \mathbf{z}_i \stackrel{\text{i.i.d.}}{\sim} q_{\phi}(\mathbf{z})\end{aligned}$$

ELBO GRADIENT ESTIMATION W.R.T. ϕ

ELBO

$$\text{ELBO} = \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(\mathbf{x}, z) - \log q_{\phi}(z)]$$

ELBO Gradient Estimation w.r.t. ϕ

$$\begin{aligned}\nabla_{\phi} \text{ELBO} &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\theta}(\mathbf{x}, z) - \log q_{\phi}(z)] \\ &= ???\end{aligned}$$

- In VAEs, we have the same problem!
- We use the **reparametrization trick** 😎

REPARAMETRIZATION TRICK

- If $q_\phi(z)$ was a 1D-Gaussian, we can produce a sample by calculating

$$z = \mu + \sigma\epsilon \quad \epsilon \sim \mathcal{N}(0, 1)$$

- We **separate** the **sample** ϵ from the **distribution parameters** $\phi = (\mu, \sigma)^T$
- In general:

$$z = g(\epsilon, \phi) \quad \epsilon \sim p(\epsilon)$$

ELBO Gradient Estimation w.r.t. ϕ

$$\begin{aligned} \nabla_\phi \mathbb{E}_{q_\phi(z)} [\log p_\theta(\mathbf{x}, z) - \log q_\phi(z)] &= \nabla_\phi \mathbb{E}_{q_\phi(z)} [f(z)] \\ &= \nabla_\phi \mathbb{E}_{p(\epsilon)} [f(g(\epsilon, \phi))] \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_\phi f(g(\epsilon, \phi))] \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_z f(z)|_{z=g(\epsilon, \phi)} \nabla_\phi g(\epsilon, \phi)] \end{aligned}$$

REPARAMETRIZATION TRICK

- Can we reparametrize all distributions q_ϕ ?
- No! 😞
- For all discrete distributions: $\nabla_\phi g(\epsilon, \phi)$ does not exist
- We can use a different estimator for ∇_ϕ ELBO (REINFORCE estimator)
 - ... which suffers from high variance
- Can be improved if we leverage dependency structure in the model