

Introduction to probabilistic programming

David Duvenaud and James Lloyd

University of Cambridge

Thanks to

Daniel M Roy (Cambridge)

HOW TO WRITE A BAYESIAN MODELING PAPER

1. Write down a generative model in an afternoon
2. Get 2 grad students to implement inference for a month
3. Use technical details of inference to pad half of the paper

CAN WE DO BETTER?

Example: Graphical Models

Application Papers

1. Write down a graphical model
2. Perform inference using general-purpose software
3. Apply to some new problem

Inference papers

1. Identify common structures in graphical models (e.g. chains)
2. Develop efficient inference method
3. Implement in a general-purpose software package

Modeling and inference have been disentangled

Not all models are graphical models

What is the largest class of models available?

Probabilistic Programs

- ▶ A probabilistic program (PP) is any program that can depend on random choices. Can be written in any language that has a random number generator.
- ▶ You can specify any computable prior by simply writing down a PP that generates samples.
- ▶ A probabilistic program implicitly defines a distribution over its output.

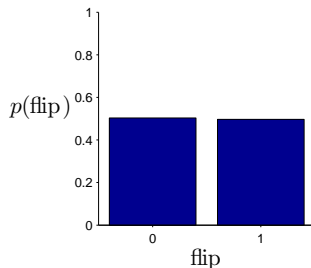
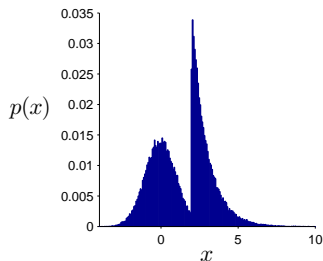
AN EXAMPLE PROBABILISTIC PROGRAM

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2    % Random draw from Gamma(1,1)
4 else
5     x = randn        % Random draw from standard Normal
6 end
```

AN EXAMPLE PROBABILISTIC PROGRAM

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2    % Random draw from Gamma(1,1)
4 else
5     x = randn        % Random draw from standard Normal
6 end
```

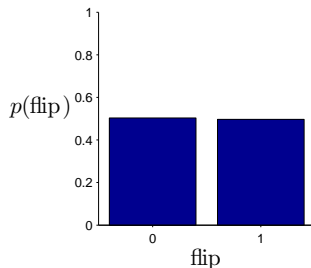
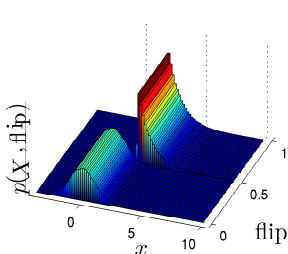
Implied distributions over variables



AN EXAMPLE PROBABILISTIC PROGRAM

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2    % Random draw from Gamma(1,1)
4 else
5     x = randn        % Random draw from standard Normal
6 end
```

Implied distributions over variables



PROBABILISTIC PROGRAMMING: CONDITIONING

Once we've defined a prior, what can we do with it?

The stochastic program defines joint distribution $P(D, N, H)$

- ▶ D to be the subset of variables we observe (condition on)
- ▶ H the set of variables we're interested in
- ▶ N the set of variables that we're not interested in, (so we'll integrate them out).

We want to know about $P(H|D)$

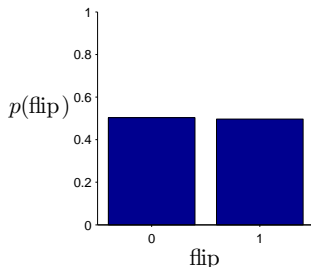
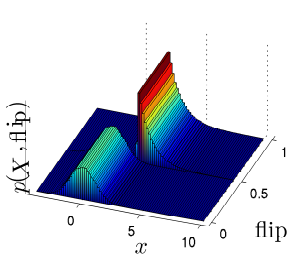
Probabilistic Programming

- ▶ Usually refers to doing conditional inference when a probabilistic program specifies your prior.
- ▶ Could also be described as automated inference given a model specified by a generative procedure.

AN EXAMPLE PROBABILISTIC PROGRAM: CONDITIONING

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2    % Random draw from Gamma(1,1)
4 else
5     x = randn        % Random draw from standard Normal
6 end
```

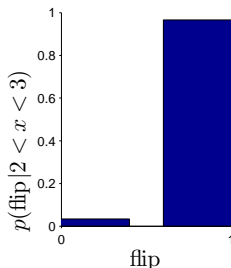
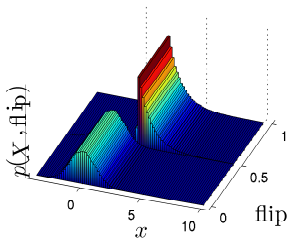
Implied distributions over variables



AN EXAMPLE PROBABILISTIC PROGRAM: CONDITIONING

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2    % Random draw from Gamma(1,1)
4 else
5     x = randn        % Random draw from standard Normal
6 end
```

Implied distributions over variables



CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

This produces samples over the execution trace

e.g.

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»True

This produces samples over the execution trace

e.g.

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»True
»2.7

This produces samples over the execution trace

e.g. (True, 2.7),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

This produces samples over the execution trace

e.g. (True, 2.7),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»True

This produces samples over the execution trace

e.g. (True, 2.7),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»True
»3.2

This produces samples over the execution trace

e.g. (True, 2.7),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

This produces samples over the execution trace

e.g. (True, 2.7),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»True

This produces samples over the execution trace

e.g. (True, 2.7),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»True
»2.1

This produces samples over the execution trace

e.g. (True, 2.7), (True, 2.1),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

This produces samples over the execution trace

e.g. (True, 2.7), (True, 2.1),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»False

This produces samples over the execution trace

e.g. (True, 2.7), (True, 2.1),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5                                »False
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end                                              »-1.3
```

This produces samples over the execution trace

e.g. (True, 2.7), (True, 2.1),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

This produces samples over the execution trace

e.g. (True, 2.7), (True, 2.1),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end
```

»False

This produces samples over the execution trace

e.g. (True, 2.7), (True, 2.1),

CAN WE DEVELOP GENERIC INFERENCE FOR ALL PPs?

Rejection sampling

1. Run the program with a fresh source of random numbers
2. If condition D is true, record H as a sample, else ignore the sample
3. Repeat

Example

```
1 flip = rand < 0.5                                »False
2 if flip
3     x = randg + 2
4 else
5     x = randn
6 end                                              »2.3
```

This produces samples over the execution trace

e.g. (True, 2.7), (True, 2.1), (False, 2.3),...

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace

- ▶ (True, 2.3)

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace
 - ▶ (True, 2.3)
2. Change one random decision, discarding subsequent decisions
 - ▶ (False,)

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace
 - ▶ (True, 2.3)
2. Change one random decision, discarding subsequent decisions
 - ▶ (False,)
3. Sample subsequent decisions
 - ▶ (False, -0.9)

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace
 - ▶ (True, 2.3)
2. Change one random decision, discarding subsequent decisions
 - ▶ (False,)
3. Sample subsequent decisions
 - ▶ (False, -0.9)
4. Accept with appropriate (RJ)MCMC acceptance probability
 - ▶ Reject, does not satisfy observation (i.e. likelihood is zero)

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace

- ▶ (True, 2.3)

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace
 - ▶ (True, 2.3)
2. Change one random decision, discarding subsequent decisions
 - ▶ (True, 2.9)

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace
 - ▶ (True, 2.3)
2. Change one random decision, discarding subsequent decisions
 - ▶ (True, 2.9)
3. Sample subsequent decisions
 - ▶ Nothing to do

CAN WE BE MORE EFFICIENT?

Metropolis-Hastings

1. Start with a trace
 - ▶ (True, 2.3)
2. Change one random decision, discarding subsequent decisions
 - ▶ (True, 2.9)
3. Sample subsequent decisions
 - ▶ Nothing to do
4. Accept with appropriate (RJ)MCMC acceptance probability
 - ▶ Accept, maybe

PP VIA METROPOLIS-HASTINGS - NOTATION

Evaluating a program results in a sequence of random choices

$$\begin{aligned}x_1 &\sim p_{t_1}(x_1). \\x_2 &\sim p_{t_2}(x_2|x_1). \\x_3 &\sim p_{t_3}(x_3|x_2, x_1). \\x_k &\sim p_{t_k}(x_k|\underbrace{x_1, \dots, x_{k-1}}_{\text{execution trace}}).\end{aligned}$$

The density/probability of a particular evaluation is then

$$p(x_1, \dots, x_K) = \prod_{k=1}^K p_{t_k}(x_k|x_1, \dots, x_{k-1}).$$

We then perform MH over the the execution trace $x = (x_1, \dots, x_K)$

MH OVER EXECUTION TRACES

1. Select a random decision in the execution trace x
 - ▶ e.g. x_k
2. Propose a new value
 - ▶ e.g. $x'_k \sim K_{t_k}(x'_k|x_k)$
3. Run the program to determine all subsequent choices $(x'_l : l > k)$, reusing current choices where possible
4. Propose moving from the state (x_1, \dots, x_K) to $(\underbrace{x_1, \dots, x_{k-1}}_{\text{old choices}}, \underbrace{x'_k, \dots, x'_{K'}}_{\text{new choices}})$
5. Accept the change with the appropriate MH acceptance probability

$$\frac{K_{t_k}(x_k|x'_k) \prod_{i=k}^{K'} p_{t'_i}(x'_i|x_1, \dots, x_{k-1}, x'_k, \dots, x'_{i-1})}{K_{t_k}(x'_k|x_k) \prod_{i=k}^K p_{t_i}(x_i|x_1, \dots, x_{k-1}, x_k, \dots, x_{i-1})}$$

DEMO: REGRESSION WITH AN INTERESTING PRIOR

NONPARAMETRIC MODELS

- ▶ If we can sample from the prior of a nonparametric model using finite resources with probability 1, then we can perform inference automatically using the techniques described thus far
- ▶ We can sample from a number of nonparametric processes/models with finite resources (with probability 1) using a variety of techniques
 - ▶ Gaussian processes via marginalisation
 - ▶ Dirichlet processes via stick breaking
 - ▶ Indian Buffet processes via urn schemes
- ▶ Active research to produce finite sampling algorithms for other nonparametric processes (e.g. hierarchical beta processes, negative binomial process)

EXAMPLE: MIXTURE OF GAUSSIANS

Generative model

$$(\mu_k)_{k=1\dots K} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$$

$$(\pi_k)_{k=1\dots K} \sim \text{Dir}(\alpha/K)$$

$$\Theta := \sum_{k=1}^K \pi_k \delta_{\mu_k}$$

$$(\theta_n)_{n=1\dots N} \stackrel{\text{iid}}{\sim} \Theta$$

$$(x_i)_{i=1\dots N} \sim \mathcal{N}(\theta_n, 1)$$

(Pseudo) MATLAB code

```
mu = randn(K,1)
pi = dirichlet(K, alpha/K)

for n = 1:N
    theta = mu(mnrnd(1,pi))
    x(n) = theta + randn
end
```

EXAMPLE: INFINITE MIXTURE OF GAUSSIANS

Limit of generative model is a DP

$$(\mu_k)_{k=1\dots K} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$$

$$(\pi_k)_{k=1\dots K} \sim \text{Dir}(\alpha/K)$$

$$\Theta := \sum_{k=1}^K \pi_k \delta_{\mu_k} \xrightarrow{K \rightarrow \infty} \Theta \sim \text{DP}(\alpha, \mathcal{N}(0, 1))$$

Avoiding infinity

- ▶ Θ is now infinitely complex, and can only be represented approximately with finite resources
- ▶ However, we can sample a finite number of samples $(\theta_n)_{n=1\dots N}$ from some unknown Θ in finite time (with probability 1) using a stick-breaking algorithm

EXAMPLE: INFINITE MIXTURE OF GAUSSIANS

MATLAB stick breaking construction

```
1 sticks = [];  
2 atoms = [];  
3 for i = 1:n  
4     p = rand;  
5     while p > sum(sticks)  
6         sticks(end+1) = (1-sum(sticks)) * betarnd(1, alpha);  
7         atoms(end+1) \ = randn;  
8     end  
9     theta(i) = atoms(find(cumsum(sticks)>=p, 1, 'first'));  
10 end  
11 x = theta' + randn(n, 1);
```


DEMOS: NONPARAMETRIC MODELS

ADVANCED AUTOMATIC INFERENCE

- ▶ Now that we have separated inference and model design, can use any inference algorithm.
- ▶ Free to develop inference algorithms independently of specific models.
- ▶ Once graphical models identified as a general class, many model-agnostic inference methods:
 - ▶ Belief Propagation
 - ▶ Pseudo-likelihood
 - ▶ Mean-field Variational
 - ▶ MCMC
- ▶ What generic inference algorithms can we implement for more expressive generative models?

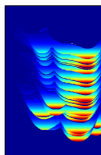
ADVANCED AUTOMATIC INFERENCE: METROPOLIS-HASTINGS

- ▶ Bher, MIT-Church
(Goodman, Mansinghka, Roy, Bonawitz and Tenenbaum, 2008)
 - ▶ (Automatic inference in Scheme)
- ▶ Stochastic Matlab
 - ▶ Lightweight Implementations of Probabilistic Programming Languages
Via Transformational Compilation
(Wingate, Stuhlmüller, Goodman, 2011)

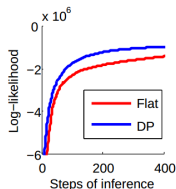
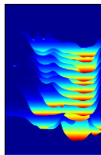
A layer model of rock porosity

```
function dp_render_rock( data )  
    num_layers = poissrnd();  
    dp_sample_layer = dpmem( @sample_layer, 1.0 );  
    for i=1:num_layers  
        layer(i,:) = dp_sample_layer();  
    end;  
    rock_volume = render_layers( layers );  
    data = rock_volume + randn( size(rock_volume) );  
    return;
```

True porosity



Inferred porosity



ADVANCED AUTOMATIC INFERENCE: HMC

- ▶ Automatic Differentiation in Church:
Nonstandard Interpretations of Probabilistic Programs for Efficient Inference
(Wingate, Goodman, Stuhlmuller, Siskind, 2012)

- ▶ Stan (Gelman et al)

<http://mc-stan.org/>

```
// Predict from Gaussian Process Logistic Regression
// Fixed covar function: eta_sq=1, rho_sq=1, sigma_sq=0.1
```

```
data {
  int<lower=1> N1;
  vector[N1] x1;
  int<lower=0,upper=1> z1[N1];
  int<lower=1> N2;
  vector[N2] x2;
}
transformed data {
  int<lower=1> N;
  vector[N1+N2] x;
  vector[N1+N2] mu;
  cov_matrix[N1+N2] Sigma;
  N <- N1 + N2;
  for (n in 1:N1) x[n] <- x1[n];
  for (n in 1:N2) x[N1 + n] <- x2[n];
  for (i in 1:N) mu[i] <- 0;
  for (i in 1:N)
    for (j in 1:N)
      Sigma[i,j] <- exp(-pow(x[i] - x[j],2))
        + if_else(i==j, 0.1, 0.0);
}
parameters {
  vector[N1] y1;
  vector[N2] y2;
}
model {
  vector[N] y;
  for (n in 1:N1) y[n] <- y1[n];
  for (n in 1:N2) y[N1 + n] <- y2[n];
```

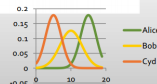
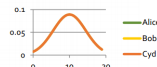
ADVANCED AUTOMATIC INFERENCE: EXPECTATION PROPAGATION

- Infer.NET (Minka, Winn, Guiver, Knowles, 2012)

- EP in graphical models:
- Now works in functional language F#:

TrueSkill in Fun

```
// prior distributions, the hypothesis
let skill() = sample (Gaussian(10.0,20.0))
let Alice,Bob,Cyd = skill(),skill(),skill()
// observe the evidence
let performance player = sample (Gaussian(player,1.0))
observe (performance Alice > performance Bob) //Alice beats Bob
observe (performance Bob > performance Cyd) //Bob beats Cyd
observe (performance Alice > performance Cyd) //Alice beats Cyd
// return the skills
Alice,Bob,Cyd
```



ADVANCED AUTOMATIC INFERENCE: VARIATIONAL

- ▶ Infer.NET has it too.
- ▶ Automated Variational Inference in Probabilistic Programming (Wingate, Weber, 2012)

Probabilistic program A

```
1: M = normal();
2: if M>1
3:   mu = complex_deterministic_func( M );
4:   X = normal( mu );
5: else
6:   X = rand();
7: end;
```

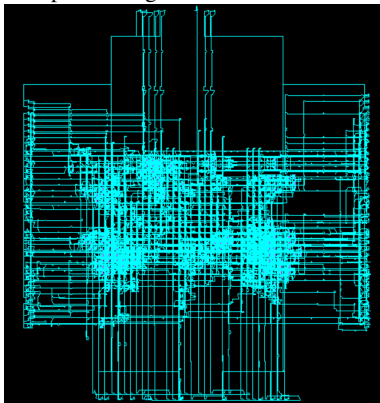
Mean-Field variational program A

```
1: M = normal(  $\theta_1$  );
2: if M>1
3:   mu = complex_deterministic_func( M );
4:   X = normal(  $\theta_3$  );
5: else
6:   X = rand( $\theta_4, \theta_5$ );
7: end;
```

- ▶ Learning phase: Forward sample, then stochastically update θ s to minimize expected KL from true distribution.
- ▶ Dependency of variational dist on control logic remains.

ADVANCED AUTOMATIC INFERENCE: HARDWARE

- ▶ Natively Probabilistic Computation ([Mansinghka, 2009](#))
- ▶ Lyric Semiconductor? (Error correcting codes)
- ▶ Main idea: If we know we're going to be sampling, some errors in computation can be OK.
 - ▶ Samplers can be made robust to computational error.
 - ▶ Run at low voltage on (cheap?) FPGA
- ▶ Compile from generative model to FPGA (9x9 Ising model sampler):



AUTOMATED MODELING

- ▶ Automated inference helpful for human modelers.
- ▶ Essential for machine-generated models
 - ▶ For example, approximate Solomonoff induction.
- ▶ Essential for more general version of automated Bayesian statistician.

THEORETICAL DIRECTIONS

Inference in stochastic programs opens up a new branch of computer science, new generalizations of computability:

- ▶ "Computable de Finetti measures"
(Freer, Roy, 2012)
- ▶ "Noncomputable conditional distributions"
(Ackerman, Freer, Roy, 2011)
- ▶ "Computable exchangeable sequences have computable de Finetti measures"
(Freer, Roy, 2009)
- ▶ "On the computability and complexity of Bayesian reasoning"
(Roy, 2012)

Main takeaways:

- ▶ No general computable algorithm exists for conditioning
- ▶ Representational choices important
 - ▶ i.e. stick-breaking vs CRP latent representation changes computability

COMPILER DEVELOPMENT

- ▶ 1950s: Ask a programmer to implement an algorithm efficiently: They'll write it on their own in assembly.
 - ▶ No good compilers; problem-dependent optimizations that only human expert could see.
- ▶ 1970s: Novice programmers use high-level languages and let compiler work out details, experts still write assembly.
 - ▶ Experts still write custom assembly when speed critical.
- ▶ 2000s: On most problems, even experts can't write faster assembly than optimizing compilers.
 - ▶ can automatically profile (JIT).
 - ▶ can take advantage of parallelization, complicated hardware, make appropriate choices w.r.t. caching.
 - ▶ Compilers embody decades of compiler research

INFERENCE METHODS IN THE FUTURE

- ▶ 2010: Ask a grad student to implement inference efficiently: They'll write it on their own.
 - ▶ No good automatic inference engines; problem-dependent optimizations that only human expert can see.
- ▶ 2015: Novice grad students use automatic inference engines and let compiler work out details, experts still write their own inference.
 - ▶ Experts still write custom inference when speed critical.
- ▶ 2020: On most problems, even experts can't write faster inference than mature automatic inference engines.
 - ▶ Can use parallelization, sophisticated hardware
 - ▶ Can automatically choose appropriate methods (meta-reasoning?).
 - ▶ Inference engines will embody 1 decade (!) of PP research.

HOW TO WRITE A BAYESIAN MODELING PAPER: 2010

1. Write down a generative model in an afternoon
2. Get 2 grad students to implement inference for a month
3. Use technical details of inference to pad half of the paper

HOW TO WRITE A BAYESIAN MODELING PAPER: 2020

1. Write down a generative model in an afternoon
2. ~~Get 2 grad students to implement inference for a month~~
Run automatic inference engine.
3. Use technical details of inference to pad half of the paper

HOW TO WRITE A BAYESIAN MODELING PAPER: 2020

1. Write down a generative model in an afternoon
2. ~~Get 2 grad students to implement inference for a month~~
Run automatic inference engine.
3. ~~Use technical details of inference to pad half of the paper~~
Discuss strengths and weaknesses of the model in the extra 4 pages.

HOW TO WRITE A BAYESIAN MODELING PAPER: 2020

1. Write down a generative model in an afternoon
2. ~~Get 2 grad students to implement inference for a month~~
Run automatic inference engine.
3. ~~Use technical details of inference to pad half of the paper~~
Discuss strengths and weaknesses of the model in the extra 4 pages.

Thanks!