# Bayesian Statistics

## Noric Couderc

Lund University
Department of Computer Science

7th December 2022

# About Me



My name is Noric

- ▶ 5th year PhD student
- ▶ From: Digne les Bains, France
- ▶ Topics: Dynamic program analysis, performance engineering
- ▶ Supervisors: Christoph Reichenbach, Emma Söderberg

My job

## My job
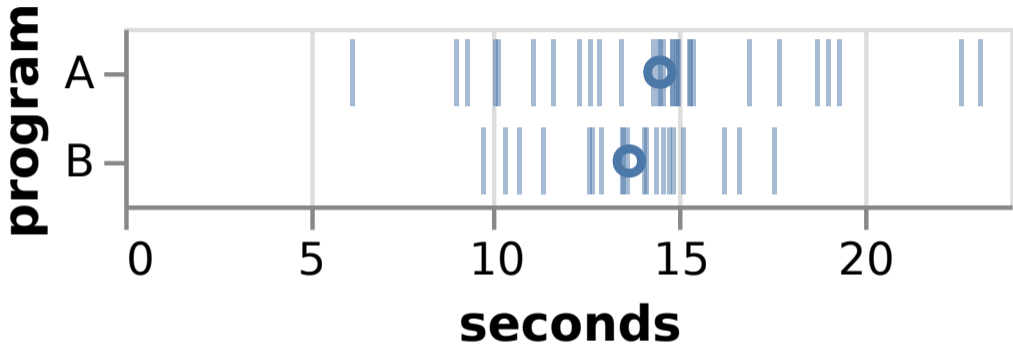
1. Run programs

### My job

1. Run programs
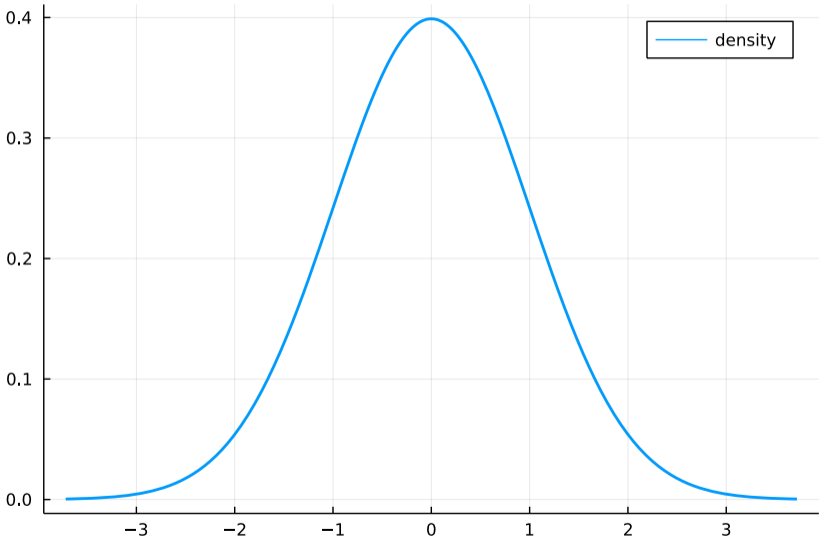2. Count stuff ⏱

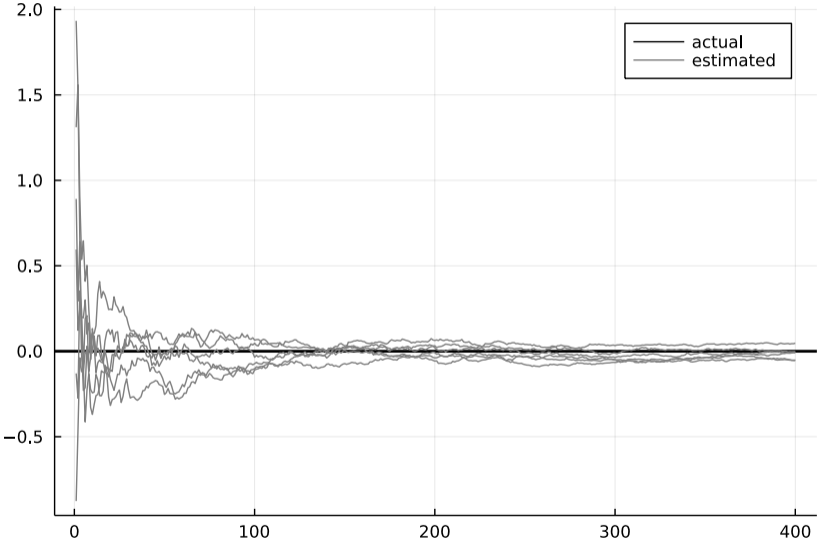# Performance Example

Figure: Which program is faster?
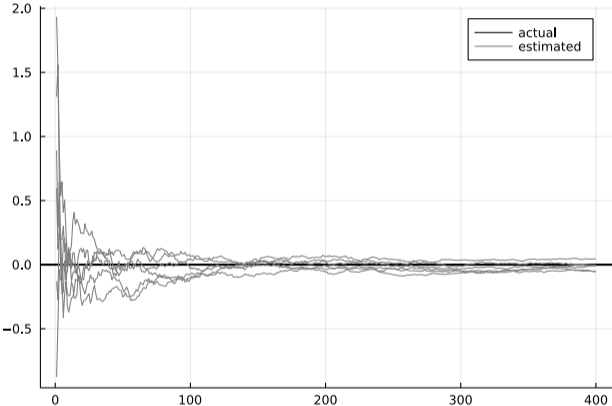
Can we trust the means we compute?
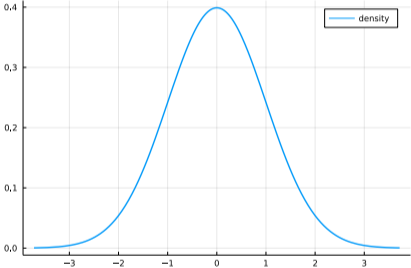
🧪 Simulate and test 🧪

# Normal Distribution

# Sampling

# Sampling

# Statistics
## Working with numbers you **can't trust**

Figure: No Man's Sky by Hello Games

fixed

Simulation

Inference

fixed

- ▶ **Prediction:** You know the inputs and the function, but not the outputs
- ▶ **Inference:** You know the outputs and the function, but not the inputs
- ▶ **Learning:** You know the inputs and the outputs, but not the function

# Bayesian inference

Make fuzzy estimates of things

# The Crush Problem

# The Crush Problem



- You observed 💐.
- What is $P(\text{❤️}|\text{💐})$?

# Maximum Likelihood Estimation

# Likelihood

$P(\text{observed data}|\text{hypothesis})$

$P(\text{💐}|\text{❤️})$

# Maximum Likelihood

$$P(\text{💐}|\text{❤️}) = 0.3$$
$$P(\text{💐}|\text{💔}) = 0.02$$

# Maximum Likelihood

$$P(\text{💐}|\text{❤️}) = 0.3$$
$$P(\text{💐}|\text{💔}) = 0.02$$

▶ P(💐|❤️) > P(💐|💔)

# Maximum Likelihood

$$P(\text{💐}|\text{❤️}) = 0.3$$
$$P(\text{💐}|\text{💔}) = 0.02$$

▶ P(💐|❤️) > P(💐|💔)
▶ **Conclusion:** ❤️

# Maximum Likelihood

$$P(\text{💐} \mid \text{❤️}) = 0.3$$
$$P(\text{💐} \mid \text{🧡}) = 0.15$$
$$P(\text{💐} \mid \text{💛}) = 0.075$$
$$P(\text{💐} \mid \text{🤍}) = 0.02$$

$$P(\text{💐}|\text{❤️}) \neq P(\text{❤️}|\text{💐})$$

# Maximum Likelihood Estimation

- ▶ Flexible!
- ▶ Returns *one single value*
- ▶ What if little data?

# Bayesian Statistics

What does a probability **mean**?

# Frequentist and Bayesian

### Frequentist

If I flip a coin for an **infinite number of times**, half of the flips will turn heads.

### Bayesian

I do not know which side the coin I will see, so I say they have **equal chances**.

# Bayes' rule

$$P(B|A) \times P(A) \;=\; P(A|B) \times P(B)$$

# Bayes' rule

$$P(A \cap B) \qquad\qquad P(B \cap A)$$

$$P(B|A) \times P(A) \;=\; P(A|B) \times P(B)$$

# Bayes' Rule

$$P(\heartsuit) \times P(\text{💐}|\heartsuit) = P(\text{💐}) \times P(\heartsuit|\text{💐})$$

# Bayes' Rule

$$P(\heartsuit \cap \text{💐}) \qquad\qquad P(\text{💐} \cap \heartsuit)$$

$$P(\heartsuit) \times P(\text{💐}|\heartsuit) \ = \ P(\text{💐}) \times P(\heartsuit|\text{💐})$$

# Bayes' Rule

$$P(\heartsuit | \text{💐}) = \frac{P(\heartsuit) \times P(\text{💐} | \heartsuit)}{P(\text{💐})}$$

# Bayes' Rule

$$P(\heartsuit \mid \text{💐}) = \frac{P(\heartsuit) \times P(\text{💐} \mid \heartsuit)}{P(\text{💐})}$$

# Bayes' Rule

Prior

$$P(\heartsuit | \text{💐}) = \frac{P(\heartsuit) \times P(\text{💐} | \heartsuit)}{P(\text{💐})}$$

# Bayes' Rule

Prior

Likelihood

$$P(\heartsuit | \text{💐}) = \frac{P(\heartsuit) \times P(\text{💐} | \heartsuit)}{P(\text{💐})}$$

# Bayes' Rule

Prior

Likelihood

Posterior $\longrightarrow$ $P(\heartsuit|\text{💐}) = \dfrac{P(\heartsuit) \times P(\text{💐}|\heartsuit)}{P(\text{💐})}$

# Bayes' Rule



Prior

Likelihood

Posterior $\longrightarrow$ $P(\heartsuit|\text{💐}) = \dfrac{P(\heartsuit) \times P(\text{💐}|\heartsuit)}{P(\text{💐})}$

Normalization

# Bayes' Rule

$$P(\textcolor{red}{\heartsuit}|\textcolor{green}{\maltese}) \propto P(\textcolor{red}{\heartsuit}) \times P(\textcolor{green}{\maltese}|\textcolor{red}{\heartsuit})$$

Bayesian Statistics

1. Express your *beliefs* in math (priors)
2. Test them against the evidence (likelihood)
3. Update your beliefs (posteriors)

- ► Flexible!
- ► Returns distributions
- ► Works with little data

# Probabilistic Programming & Bayesian Methods for Hackers

An intro to Bayesian methods and probabilistic programming from a computation/understanding-first, mathematics-second point of view.
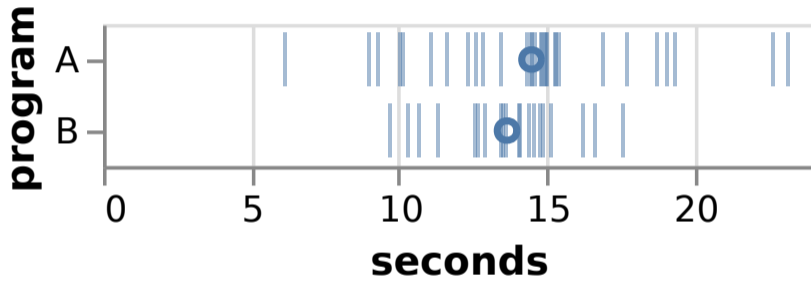
## Prologue

The Bayesian method is the natural approach to inference, yet it is hidden from readers behind chapters of slow, mathematical analysis. The typical text on Bayesian inference involves two to three chapters on probability theory, then enters what Bayesian inference is. Unfortunately, due to mathematical intractability of most Bayesian models, the reader is only shown simple, artificial examples. This can leave the user with a *so-what* feeling about Bayesian inference. In fact, this was the author's own prior opinion.

After some recent success of Bayesian methods in machine-learning competitions, I decided to investigate the subject again. Even with my mathematical background, it took me three straight-days of reading examples and trying to put the pieces together to understand the methods. There was simply not enough literature bridging theory to practice. The problem with my misunderstanding was

# Comparing two means

# Mathematical Model

$$A_i \sim \text{Normal}(\mu_A, \sigma_A)$$
$$B_i \sim \text{Normal}(\mu_B, \sigma_B)$$
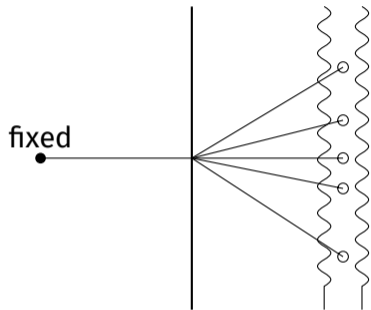$$\mu_A \sim \text{Normal}(10.0, 5.0)$$
$$\mu_B \sim \text{Normal}(10.0, 5.0)$$
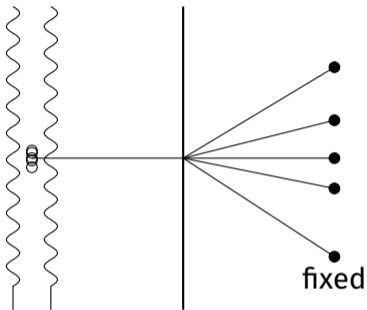$$\sigma_A \sim \text{Exponential}(1.0)$$
$$\sigma_B \sim \text{Exponential}(1.0)$$

- x Parameter
- x Observed

fixed

Simulation

Inference

fixed

```julia
function generate_samples(mean_a, mean_b, stdev_a, stdev_b)
    n_a = 30
    n_b = 20

    As = zeros(n_a)
    Bs = zeros(n_b)
    for i in 1:n_a
        As[i] = rand(Normal(mean_a, stdev_a))
    end

    for i in 1:n_b
        Bs[i] = rand(Normal(mean_b, stdev_b))
    end

    return (As,Bs)
end
```

```julia
@model function means_model(As, Bs)
    # Priors
    mean_a ~ Normal(10.0, 5.0)
    mean_b ~ Normal(10.0, 5.0)

    stdev_a ~ Exponential(1.0)
    stdev_b ~ Exponential(1.0)

    # Likelihood
    for i in eachindex(As)
        As[i] ~ Normal(mean_a, stdev_a)
    end

    for i in eachindex(Bs)
        Bs[i] ~ Normal(mean_b, stdev_b)
    end
end
```
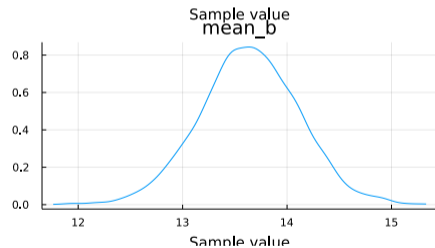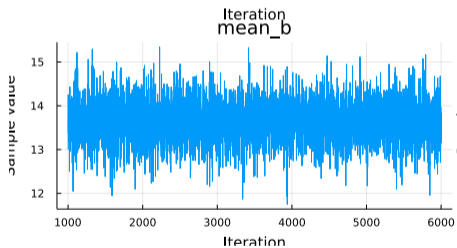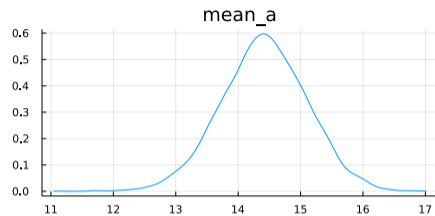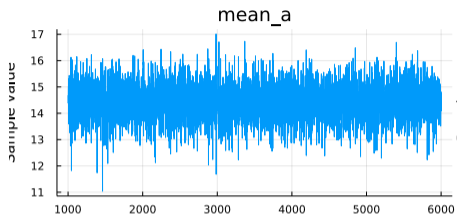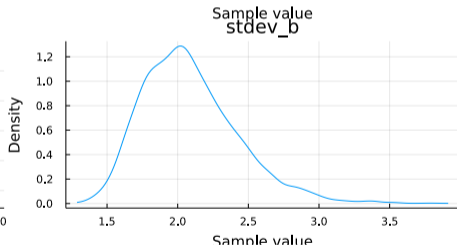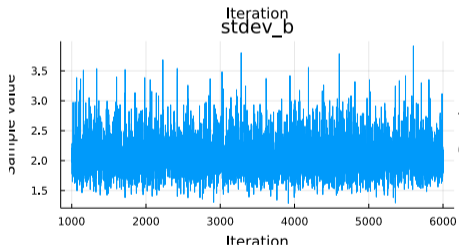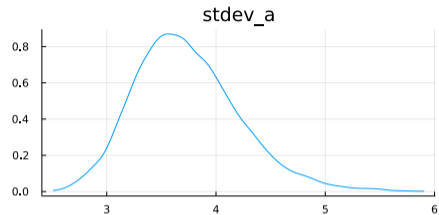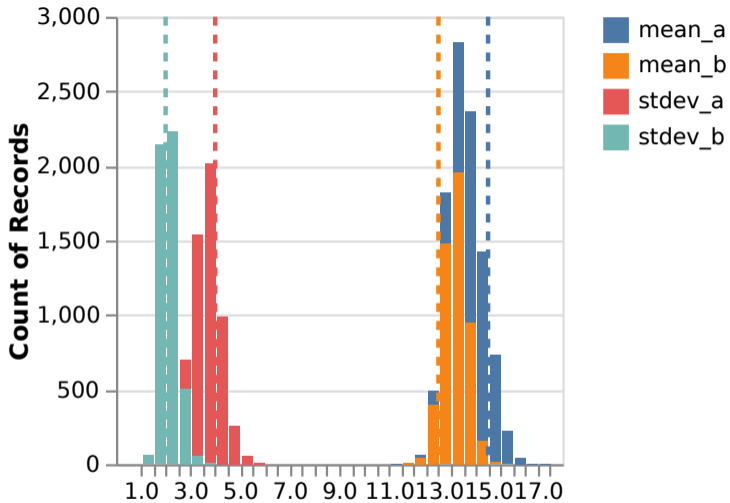
```
# Create model (pass the data to it)
model = means_model(samples_a, samples_b)

# Run bayesian inference
chain = sample(model, NUTS(0.65), 5000)

# Get posterior samples for a variable in model
chain[:mean_a]
```

```
diff = chain[:mean_b] - chain[:mean_a]
```



Difference mean B vs mean A

# Posterior predictive simulation

# Predictions return *distributions* too!

- ▶ **Prediction:** You know the inputs and the function, but not the outputs
- ▶ **Inference:** You know the outputs and the function, but not the inputs
- ▶ **Learning:** You know the inputs and the outputs, but not the function

## Posterior predictive simulation

Make **predictions** with some uncertain inputs

```
function generate_samples(mean_a, mean_b, stdev_a, stdev_b)
    n_a = 30
    n_b = 20

    As = zeros(n_a)
    Bs = zeros(n_b)
    for i in 1:n_a
        As[i] = rand(Normal(mean_a, stdev_a))
    end

    for i in 1:n_b
        Bs[i] = rand(Normal(mean_b, stdev_b))
    end

    return (As,Bs)
end
```
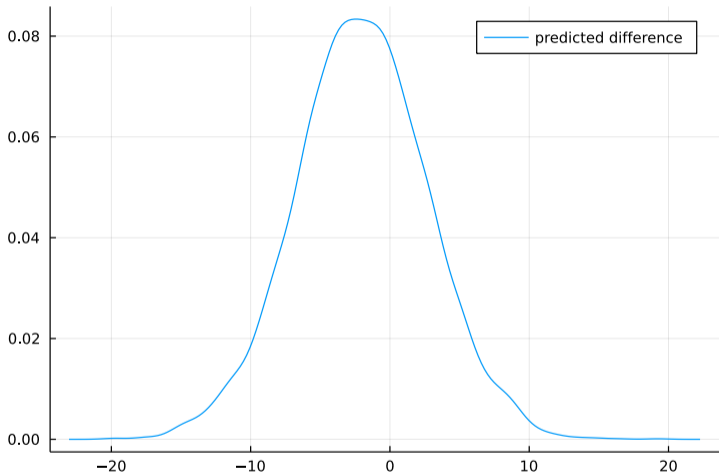
# Linear regression

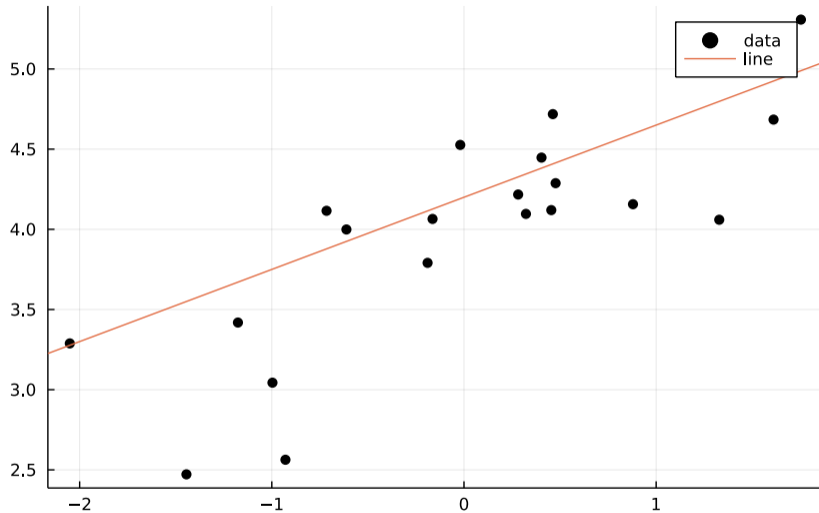Figure: Simulated data

Why is linear regression useful?

- ▶ simple model, but easy to extend
- ▶ shows positive/negative associations
- ▶ interpretable

$$y_i \sim \text{Normal}(\mu_i, \sigma)$$
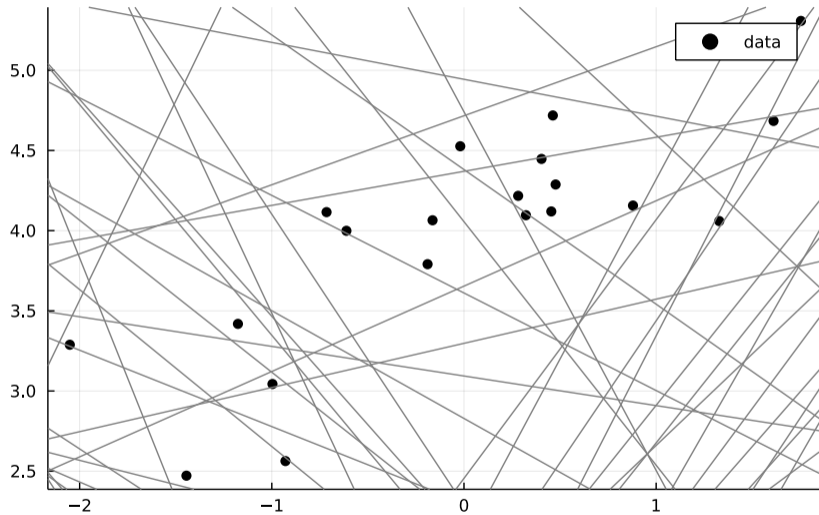$$\mu_i = \alpha x_i + \beta$$
$$\alpha \sim \text{Normal}(0, 1)$$
$$\beta \sim \text{Normal}(0, 1)$$
$$\sigma \sim \text{Exponential}(0.2)$$

- ▶ x Parameter
- ▶ x Observed

# Prior lines

# Posterior lines

Observations     Structure

Prior(s)    Likelihood

Posterior(s)

# Posterior predictive simulation

Predictions return *distributions* too!

prediction intervals

Posterior prediction

# Interactions

### Examples of interactions

► Plants need both **light** and **water** to grow.
► To have sweet coffee you need to **add sugar** and **stir** the coffee.
► **Some optimizations** only work on **some CPUs**.

```julia
@model function interaction_model(N, xs, ys, teams)
    a ~ filldist(Normal(0.0, 0.5), 2)

    b ~ filldist(Normal(0.0, 0.5), 2)

    sigma ~ Exponential(1.0)

    for i in 1:N
        if teams[i] == 1.0
            mu_1 = a[1] * xs[i] + b[1]
            ys[i] ~ Normal(mu_1, sigma)
        else
            mu_2 = a[2] * xs[i] + b[2]
            ys[i] ~ Normal(mu_2, sigma)
        end
    end
end
```

# Analysing Experiments

How to make my programs faster?

How to make my programs faster?
It **depends**

How to make my programs faster?
It **depends**
Depends on what?

Collection type ──────→ Execution time ←────── Hardware

Collection usage ↗        ↑        ↖ Compiler

Runtime

Collection type $\longrightarrow$ Execution time $\longleftarrow$ Hardware

Collection usage $\nearrow$

Warmup $\uparrow$

Compiler $\nwarrow$

JVM $\uparrow$

Collection type $\longrightarrow$ Execution time $\longleftarrow$ Hardware

JVM Warmup

Collection type $\longrightarrow$ Execution time $\longleftarrow$ Hardware

$\uparrow$

JVM Warmup

- ▶ **Machines:** Athena, Hades, Hera
- ▶ **Warmup:** Startup (cold JVM), Steady-state (warm JVM)
- ▶ **Treatments:** 0->arraylist, 1->arraymap, ...
    - ▶ Allocation site: 0, 1, 2, ...
    - ▶ Collection: hashmap, arraylist, ...

**status**

Startup · Steady-state

$$ET_j \sim \text{Normal}(\mu_j, \sigma)$$

$$\mu_j = \overline{ET} \times \alpha[M_j] \times \beta[W_j] \times \gamma[T_j]$$

$$\alpha_m \sim \text{Normal}(0.0, 0.5)$$

$$\beta_w \sim \text{Normal}(0.0, 0.2)$$

$$\gamma_t \sim \text{Normal}(0.0, 0.2)$$

$$\sigma \sim \text{Exponential}(30.0)$$

▶ x Parameter
▶ x Observed

$$ET_j \sim \mathsf{Normal}(\mu_j, \sigma)$$

$$\mu_j = \overline{ET} \times \alpha[M_j] \times \beta[W_j] \times \gamma[T_j] \longleftarrow \quad \text{Array access}$$

$$\alpha_m \sim \mathsf{Normal}(0.0, 0.5)$$

$$\beta_w \sim \mathsf{Normal}(0.0, 0.2)$$

$$\gamma_t \sim \mathsf{Normal}(0.0, 0.2)$$

$$\sigma \sim \mathsf{Exponential}(30.0)$$

▶ x Parameter
▶ x Observed

$$ET_j \sim \mathsf{Normal}(\mu_j, \sigma)$$

$$\mu_j = \overline{ET} \times \alpha[M_j] \times \beta[W_j] \times \gamma[T_j] \longleftarrow \quad \text{Array access}$$

Arrays $\longrightarrow$ $\alpha_m \sim \mathsf{Normal}(0.0, 0.5)$

$\longrightarrow$ $\beta_w \sim \mathsf{Normal}(0.0, 0.2)$

$$\gamma_t \sim \mathsf{Normal}(0.0, 0.2)$$

$$\sigma \sim \mathsf{Exponential}(30.0)$$

▶ x Parameter
▶ x Observed

$$\log(ab) = \log(a) + \log(b)$$

$$\log(ET_j) \sim \text{Normal}(\mu_j, \sigma)$$
$$\mu_j = \log(\overline{ET}) + \alpha[M_j] + \beta[W_j] + \gamma[T_j]$$
$$\alpha_m \sim \text{Normal}(0.0, 0.5)$$
$$\beta_w \sim \text{Normal}(0.0, 0.2)$$
$$\gamma_t \sim \text{Normal}(0.0, 0.2)$$
$$\sigma \sim \text{Exponential}(1.0)$$

► x Parameter
► x Observed

Figure: Effect of machine

Figure: Effect of warmup

Figure: Effect of treatment

(Discrete) Interactions

Transform an **array** of coefficients into a **matrix**

$$\log(ET_j) \sim \text{Normal}(\mu_j, \sigma)$$

$$\mu_j = \log(\overline{ET}) + \alpha[M_j] + \beta[W_j] + \gamma^M[M_j, T_j] + \gamma^W[W_j, T_j]$$

$$\alpha_m \sim \text{Normal}(0.0, 0.5)$$

$$\beta_w \sim \text{Normal}(0.0, 0.2)$$

$$\gamma^W_{w,t} \sim \text{Normal}(0.0, 0.2)$$

$$\gamma^M_{m,t} \sim \text{Normal}(0.0, 0.2)$$

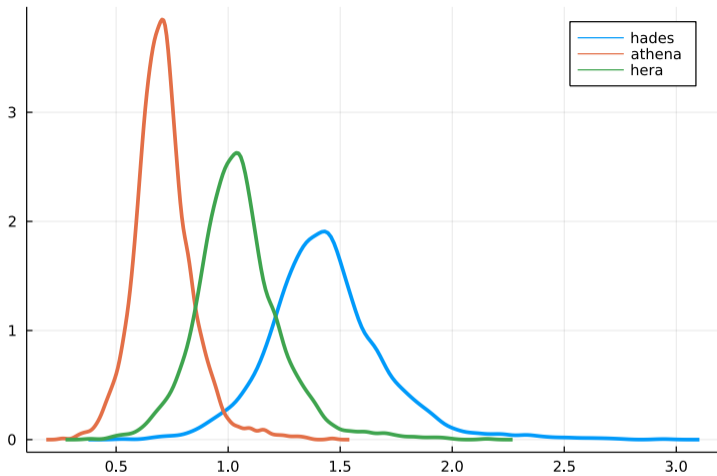$$\sigma \sim \text{Exponential}(1.0)$$

▶ x Parameter
▶ x Observed

# Hierarchical Models
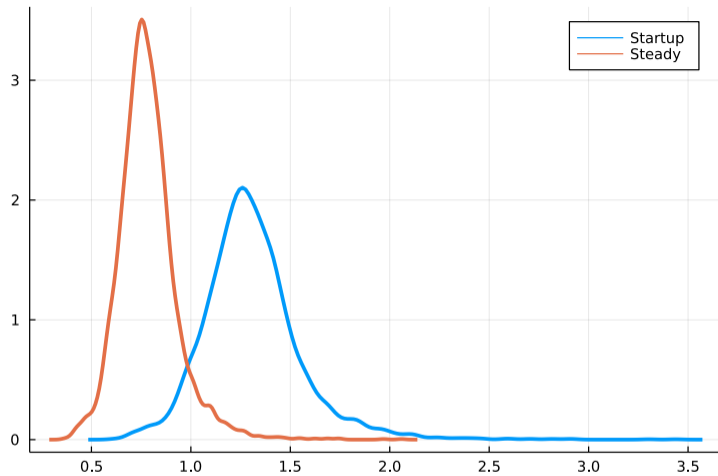
**Hyperpriors**
Priors can have parameters too

**Hyperpriors**
Priors can have parameters too
🤯

```julia
@model function interaction_model(N, xs, ys, teams)
    a ~ filldist(Normal(0.0, 0.5), 2)

    b ~ filldist(Normal(0.0, 0.5), 2)

    sigma ~ Exponential(1.0)

    for i in 1:N
        if teams[i] == 1.0
            mu_1 = a[1] * xs[i] + b[1]
            ys[i] ~ Normal(mu_1, sigma)
        else
            mu_2 = a[2] * xs[i] + b[2]
            ys[i] ~ Normal(mu_2, sigma)
        end
    end
end
```

```julia
@model function interaction_model(N, xs, ys, teams)
    sigma_slopes ~ Exponential(0.1)
    a ~ filldist(Normal(0.0, sigma_slopes), 2)

    b ~ filldist(Normal(0.0, 0.5), 2)

    sigma ~ Exponential(1.0)

    for i in 1:N
        if teams[i] == 1.0
            mu_1 = a[1] * xs[i] + b[1]
            ys[i] ~ Normal(mu_1, sigma)
        else
            mu_2 = a[2] * xs[i] + b[2]
            ys[i] ~ Normal(mu_2, sigma)
        end
    end
end
```

## Why hyperpriors

- ▶ Prevent overfitting
- ▶ Use *all* the data

# Conclusion

Bayesian inference is:

- **intuitive:** Use probability distributions, no need to learn specific vocabulary
- **flexible:**
  - Use the distributions you want
  - Use the code you want
- **honest:** If not enough data, you get the priors back
- **explicit:** All your assumptions are visible

Bayesian inference is:

- **intuitive:** Use probability distributions, no need to learn specific vocabulary
- **flexible:**
    - Use the distributions you want
    - Use the code you want
- **honest:** If not enough data, you get the priors back
- **explicit:** All your assumptions are visible
- Sometimes **slow**

My Stats Workflow

1. Generate some fake data (makes you think about what causes what)
2. Make a statistical model
3. Test if the model can recover parameters
4. Use it on experimental data

Toys to Play With

- ▶ Stan (Standalone DSL)
- ▶ Turing (Julia)
- ▶ Pyro (Python, supported by Uber)
- ▶ Bean Machine (Python, supported by Meta)
- ▶ Tensorflow Probability (Python, supported by Google)
- ▶ PyMC3 (Python)

## Books to Read

- ▶ Bayesian Methods for Hackers
- ▶ Statistical Rethinking (McElreath)
- ▶ Causal Inference In Statistics, a primer (Pearl, Glymour, Jewell)

Thanks!

# Extra Slides

# Prior Sensitivity Analysis

# The Idea

Try out different priors, and check how predictions change!

# Causal Inference

# Correlation isn't causation
(But then what is causation?)

# Causation

X causes Y if Y relies on X for its value
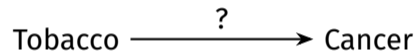
How can correlation appear?

- ▶ X causes Y
- ▶ Y causes X
- ▶ X and Y share common cause Z (Confounding)
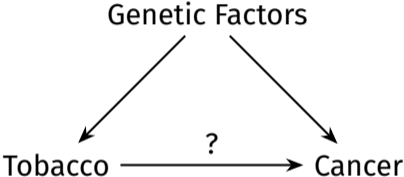- ▶ Selection bias (Colliders)

Direct causes

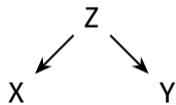$$X \longrightarrow Y$$

$$X \longleftarrow Y$$

# Confounding

$$\text{Tobacco} \xrightarrow{\quad ? \quad} \text{Cancer}$$

# Confounding

Confounder (or Fork)

Z

X          Y

# Collider

$$X \rightarrow Z \leftarrow Y$$