



Machine Learning + Statistics: Better Together Than Apart

Daniela Huppenkothen

DIRAC Institute, University of Washington



[https://github.com/dhuppenkothen/
cargese2018_tutorials](https://github.com/dhuppenkothen/cargese2018_tutorials)

All models are **wrong**,
but some are **useful**.

– George Box



Iterative Understanding of Nature



Iterative Understanding of Nature

- Experimental Design



Iterative Understanding of Nature

- Experimental Design
 - Data Collection
-



Iterative Understanding of Nature

- Experimental Design
 - Data Collection
 - Model Formulation
-



Iterative Understanding of Nature

- Experimental Design
 - Data Collection
 - Model Formulation
 - Model Criticism
-



Iterative Understanding of Nature

- Experimental Design
 - Data Collection
 - Model Formulation
 - Model Criticism
-

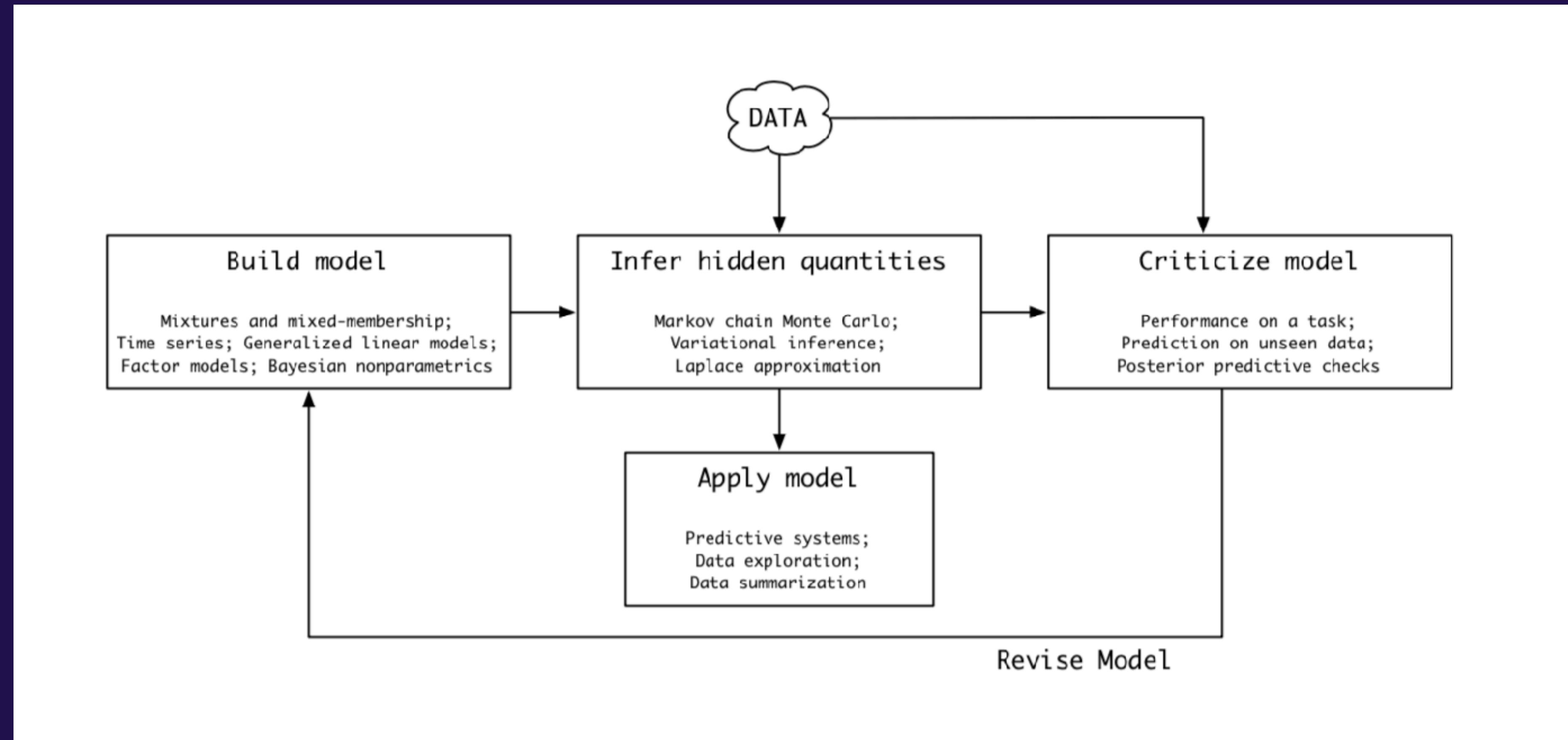




Iterative Understanding of Nature

- Experimental Design
- Data Collection
- Model Formulation
- Model Criticism





Blei (2014)

Building + Revising Models

Building + Revising Models

Part III

Building + Revising Models

Part III

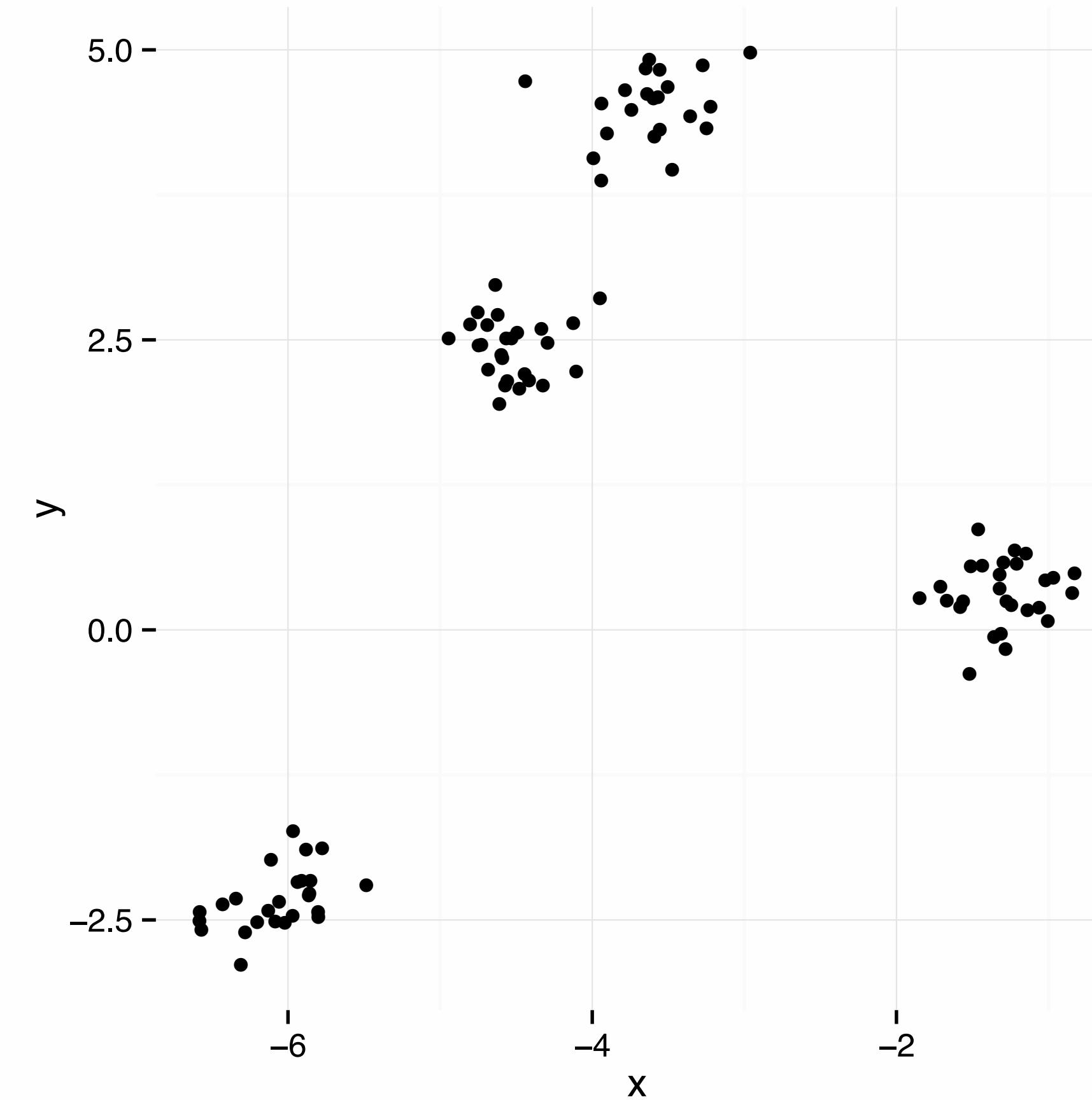
(sort of)

Wednesday: Bayesian Hierarchical Models

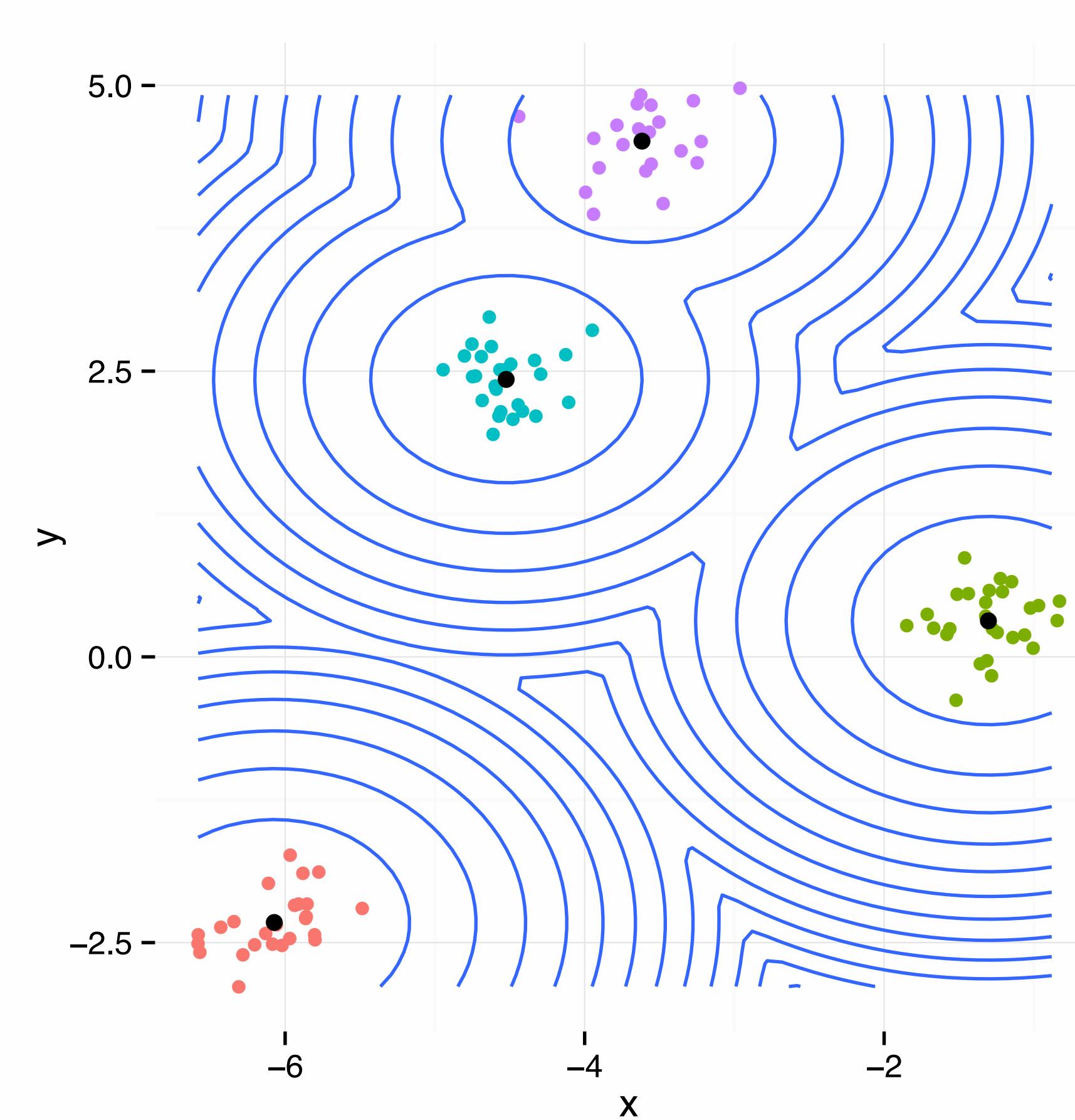
Yesterday: Machine Learning

Today: Both!

Gaussian Mixture Models



Blei (2014)



(1) The Generative Process

1. Draw mixture proportions $\theta \sim \text{Dirichlet}(\alpha)$
2. For each mixture component k , draw $\mu_k \sim \mathcal{N}(0, \sigma^2)$
3. For each data point i :
 - a. Draw mixture assignment $z_i | \theta \sim \text{Discrete}(\theta)$
 - b. Draw data point $x_i | z_i, \mu \sim \mathcal{N}(\mu_{z_i}, 1)$

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$



cluster
proportions

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

cluster
proportions

Gaussian
means

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

cluster
proportions

Gaussian
means

cluster
assignments

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

cluster
proportions

Gaussian
means

data
points
cluster
assignments

(2) The Joint Probability Distribution

Dirichlet

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

cluster
proportions

Gaussian
means

data
points
cluster
assignments

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

cluster proportions

Gaussian means

cluster assignments

data points

Dirichlet

Normal

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

Dirichlet Normal Discrete

cluster proportions data points Gaussian means cluster assignments

The diagram illustrates the joint probability distribution of parameters θ , cluster means μ , cluster assignments z , and data points x , given hyperparameters σ_0^2 and α . The distribution is factored into three components: a Dirichlet prior on θ , a product of K Normal priors on μ , and a product of N likelihood terms involving z and x . The labels 'cluster proportions', 'data points', 'Gaussian means', and 'cluster assignments' are positioned below the equation, with arrows pointing to the corresponding terms: 'cluster proportions' to θ , 'data points' to x , 'Gaussian means' to μ , and 'cluster assignments' to z .

(2) The Joint Probability Distribution

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N (p(z_i | \theta) p(x_i | z_i, \mu, z_i))$$

Dirichlet Normal Discrete Normal

cluster proportions Gaussian means data points cluster assignments

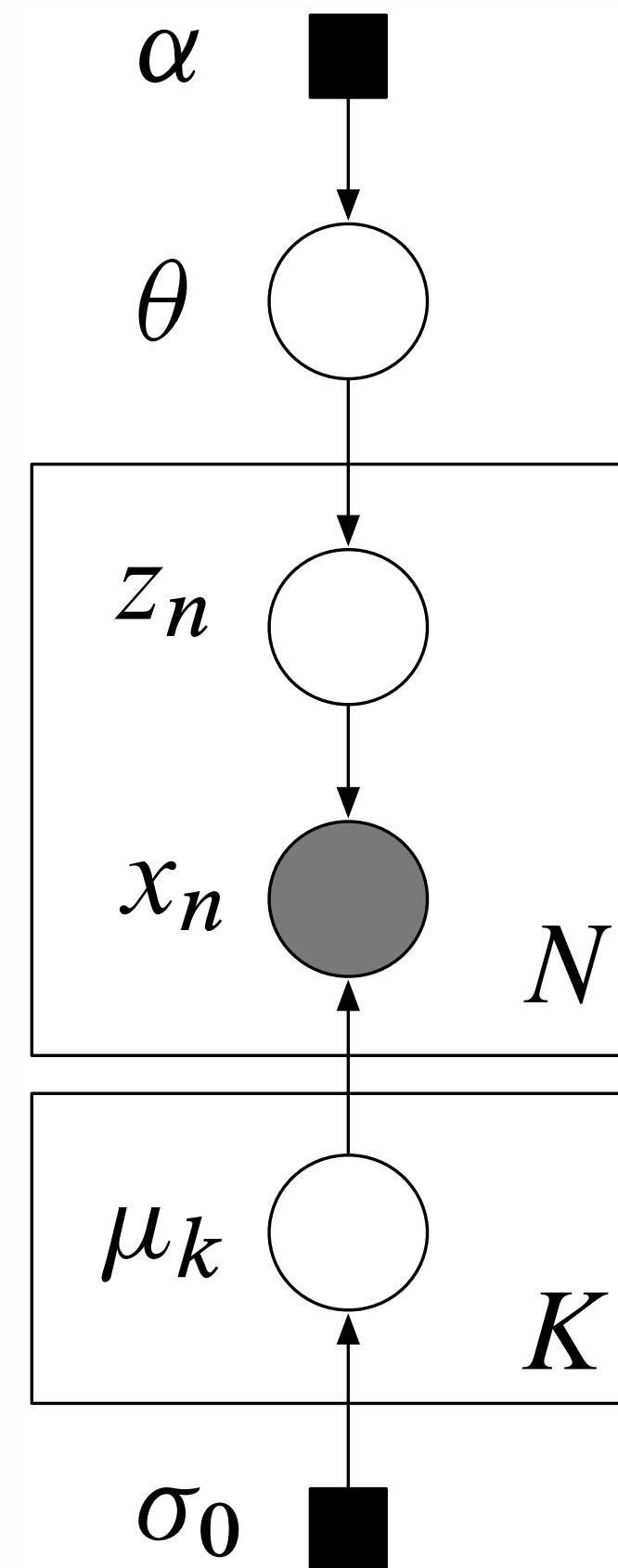
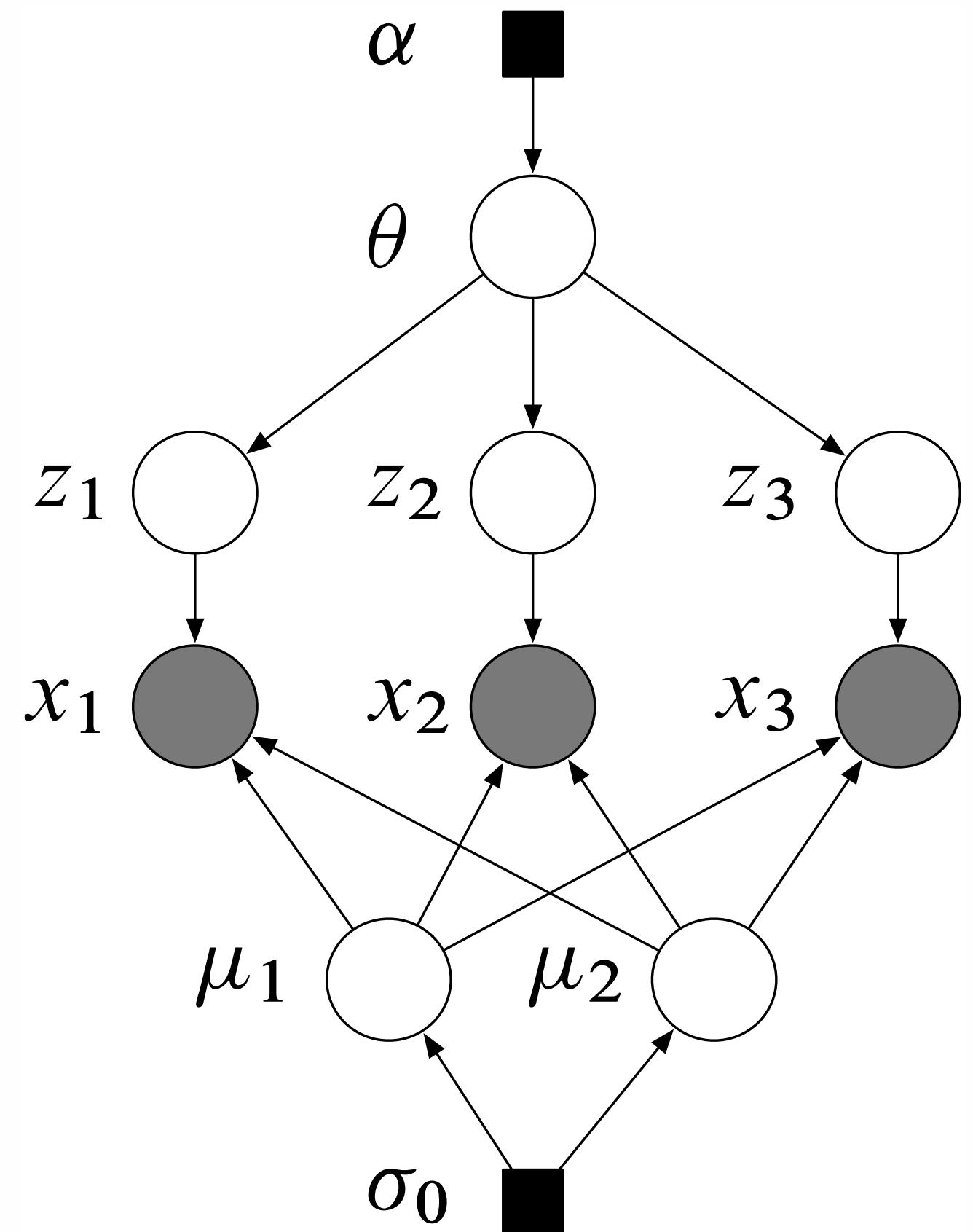
The diagram illustrates the joint probability distribution of parameters θ , cluster proportions z , and data points x . The distribution is factorized into four components: a Dirichlet prior for θ , Gaussian priors for μ , discrete priors for z , and a likelihood term involving both z and x . Arrows point from the labels to their corresponding terms in the equation.

(3) The Probabilistic Graphical Model

Exercise: Can you draw a graphical model for the Gaussian Mixture Model?

(3) The Probabilistic Graphical Model

(3) The Probabilistic Graphical Model



$\sim \text{Dirichlet}_K(\alpha)$

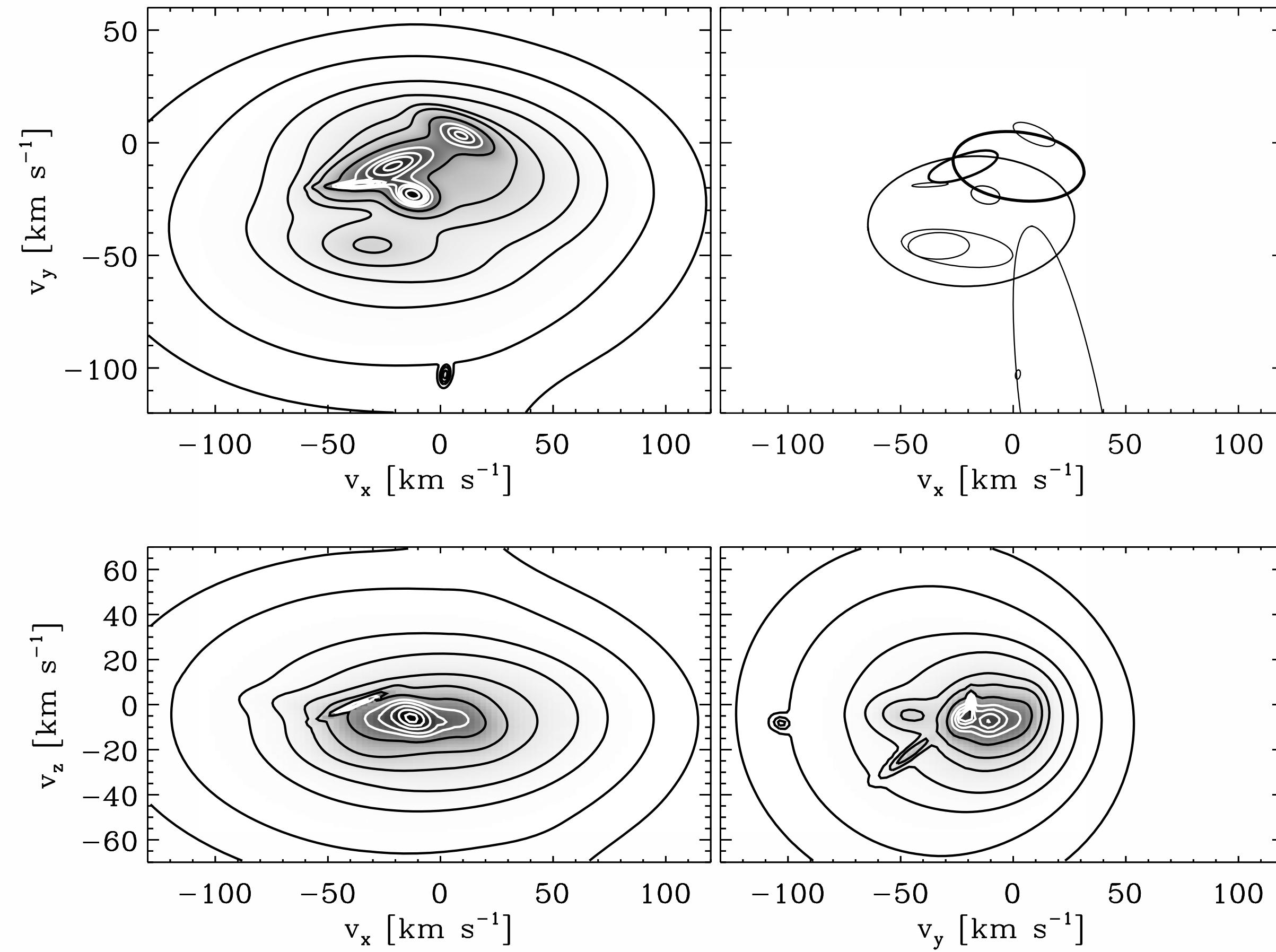
$\sim \text{Discrete}(\theta)$

$\sim \mathcal{N}(\mu_{z_n}, \sigma^2)$

$\sim \mathcal{N}(0, \sigma_0^2)$

Measurement Uncertainties? Missing Data?

Extreme Deconvolution

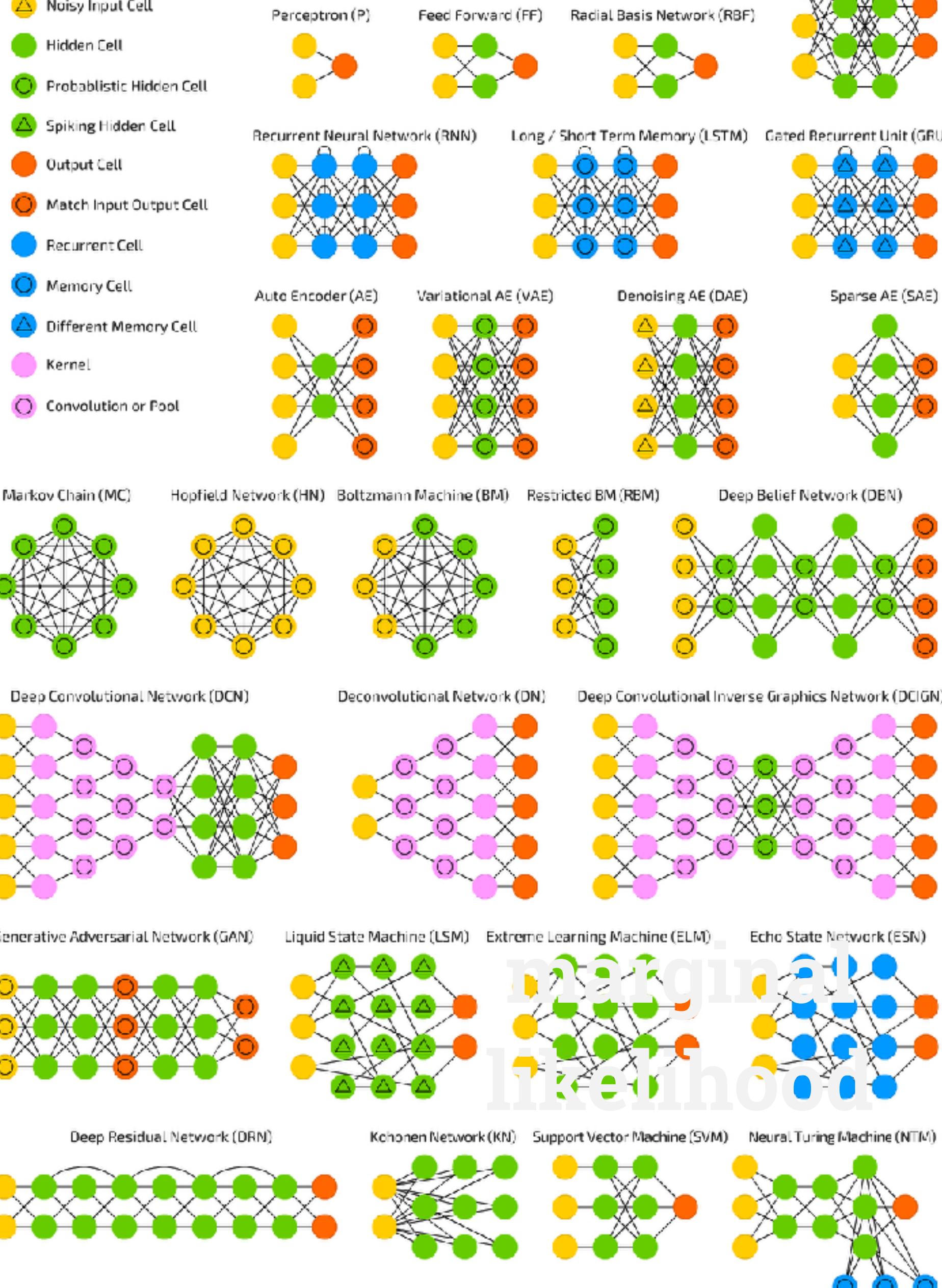


Neural Networks + Probabilistic Models

Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool





<https://playground.tensorflow.org>

Using Neural Networks as Emulators

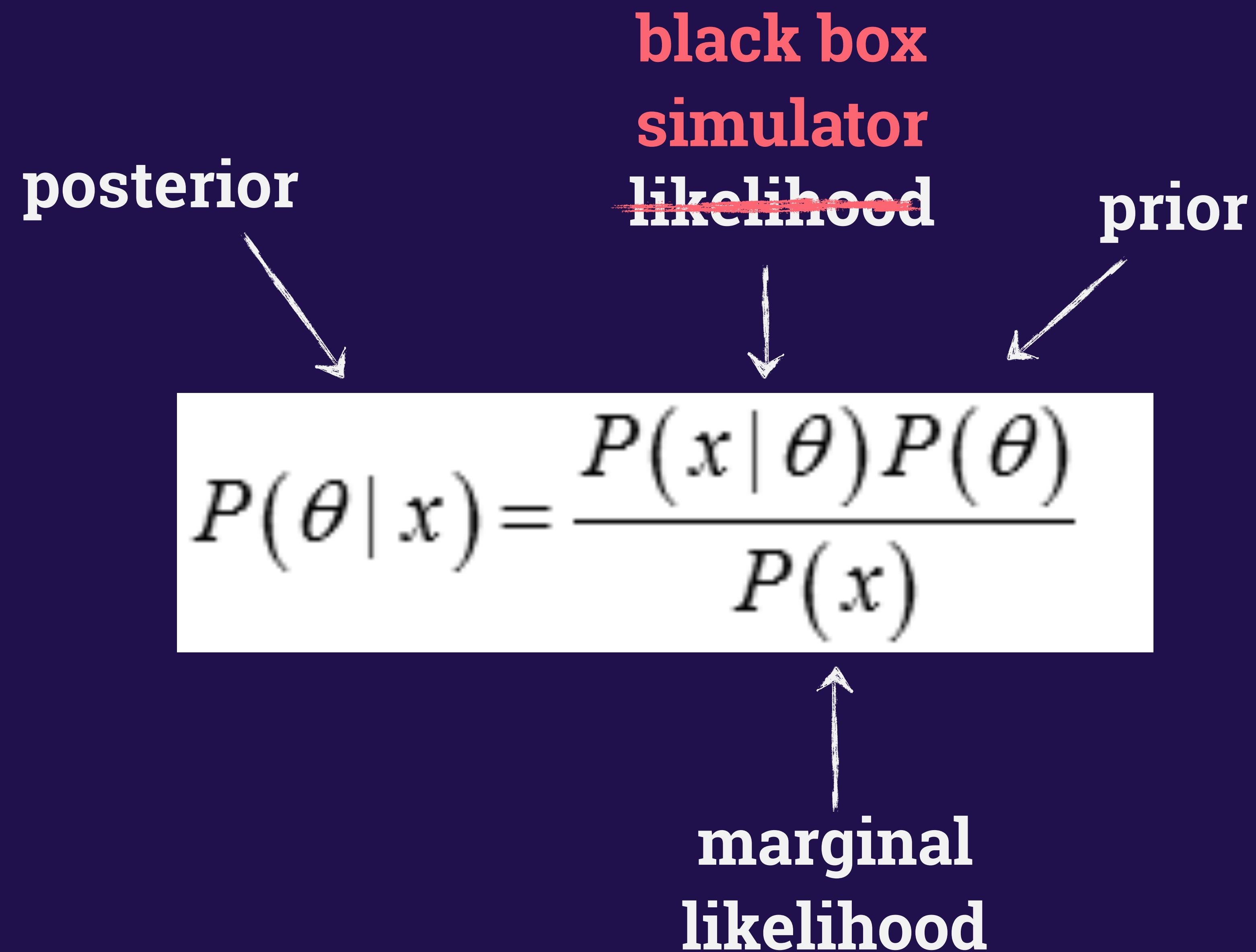
$$P(\theta | x) = \frac{P(x | \theta) P(\theta)}{P(x)}$$

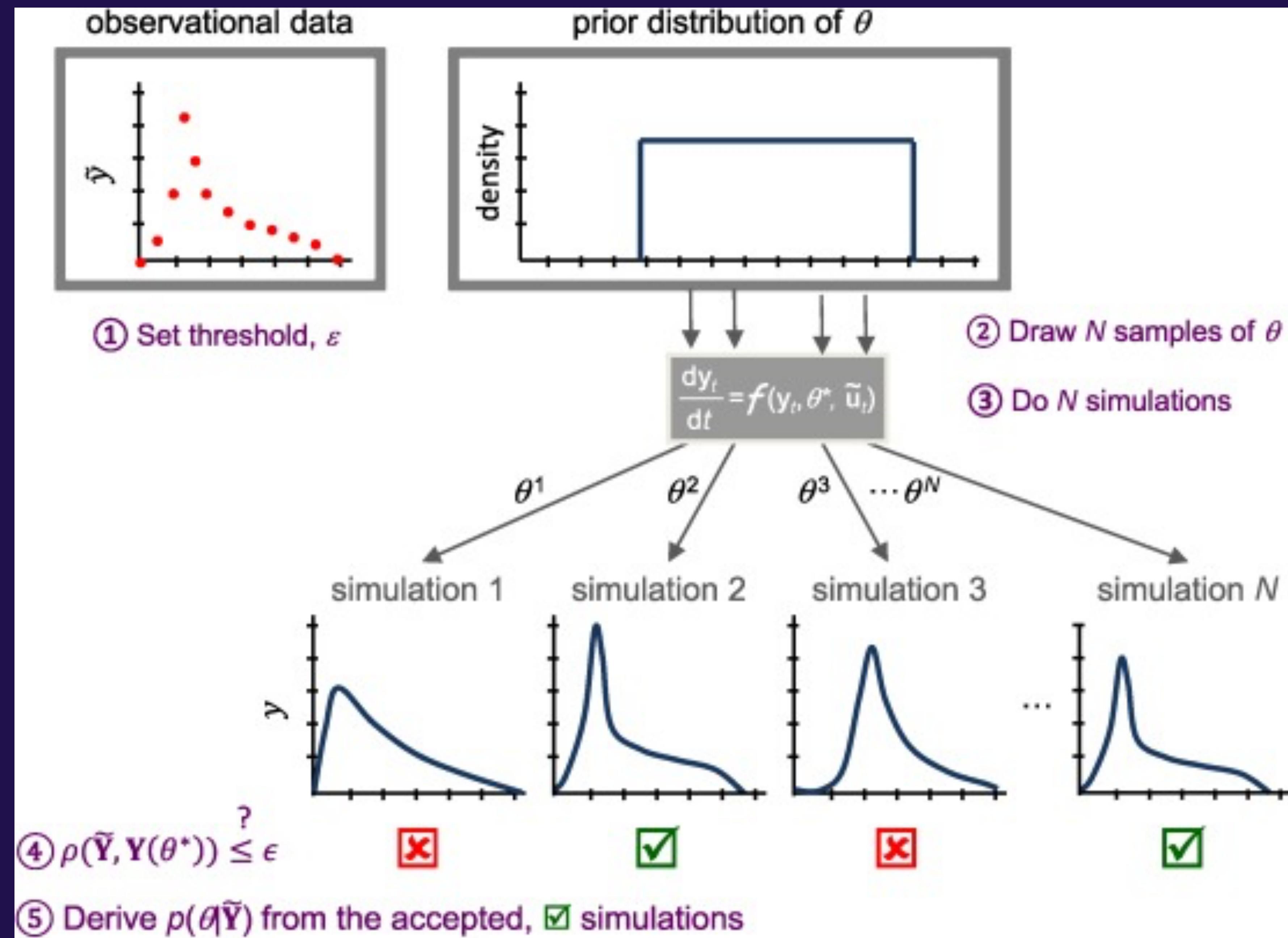
posterior

likelihood

prior

marginal likelihood







observational data

prior distribution of θ

① Set threshold, ε

② Draw N samples of θ

③ Do N simulations

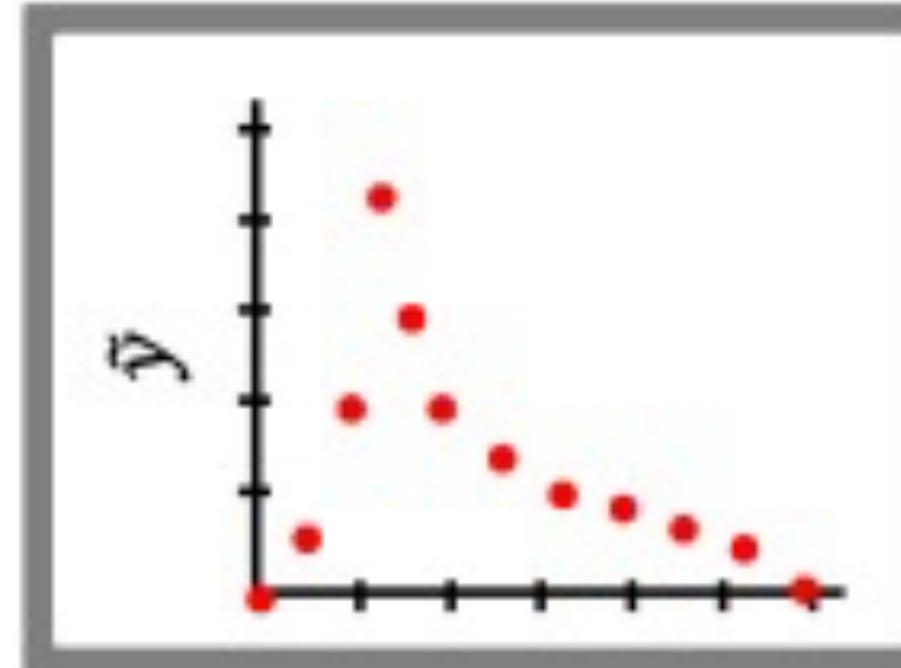
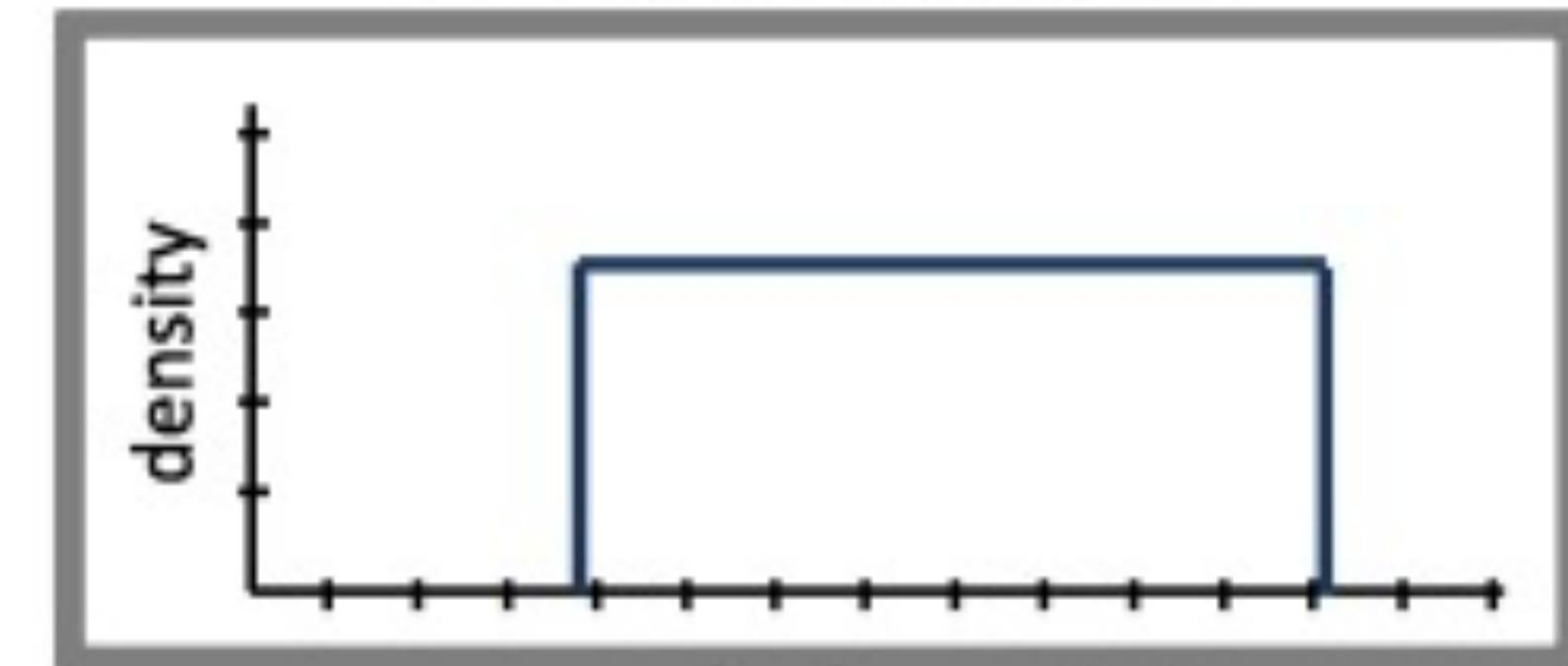
$\frac{dy_t}{dt} = f(y_t, \theta^*, \tilde{u}_t)$

$\theta^1, \theta^2, \dots, \theta^N$

simulation 2 simulation 3 ... simulation N

Approximate Bayesian Computation
Likelihood-Free Inference
Implicit Models

observational data

prior distribution of θ 

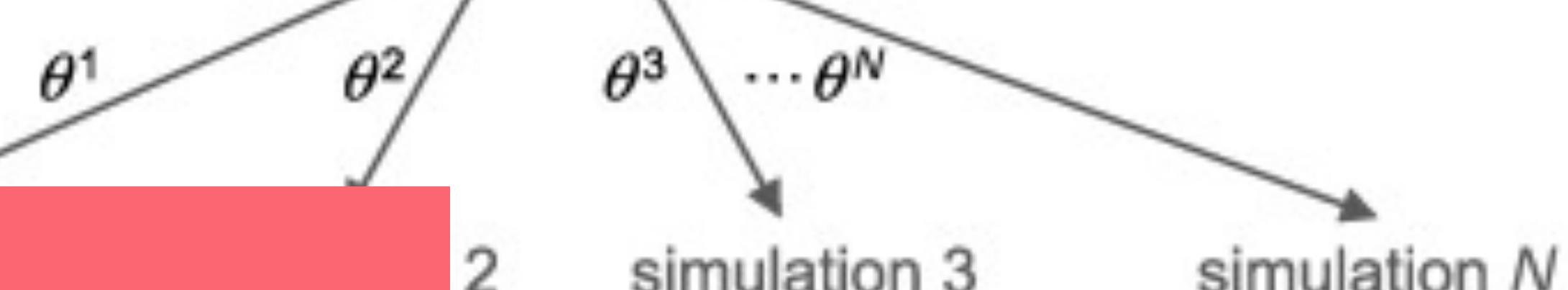
① Set threshold, ε

② Draw N samples of θ

③ Do N simulations

expensive!

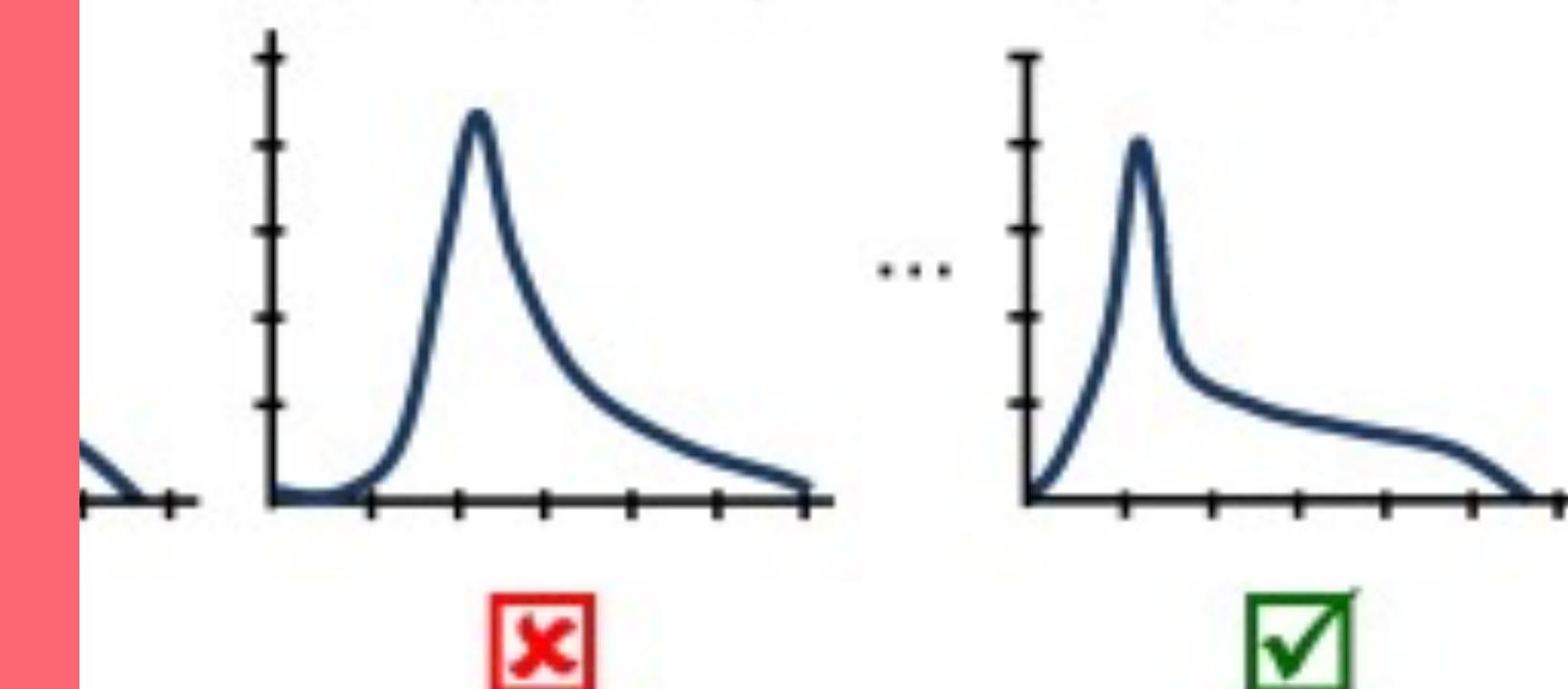
$$\frac{dy_t}{dt} = f(y_t, \theta^*, \tilde{u}_t)$$



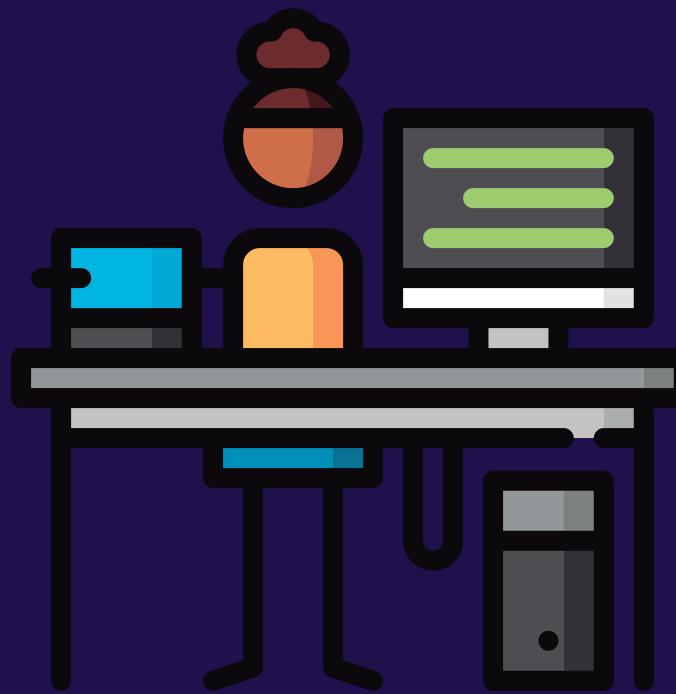
Approximate Bayesian Computation

Likelihood-Free Inference

Implicit Models



data
analysis

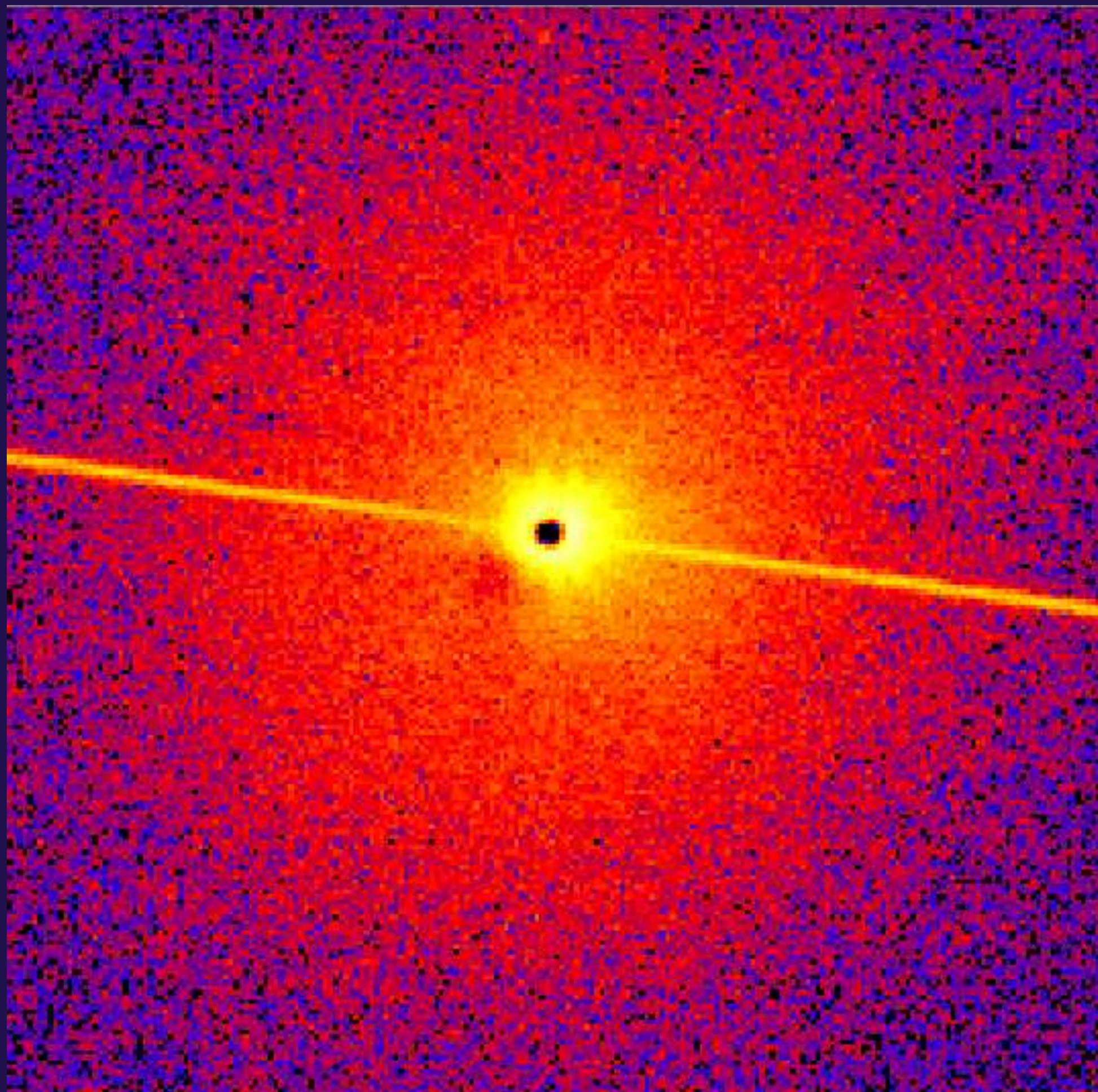


data
collection

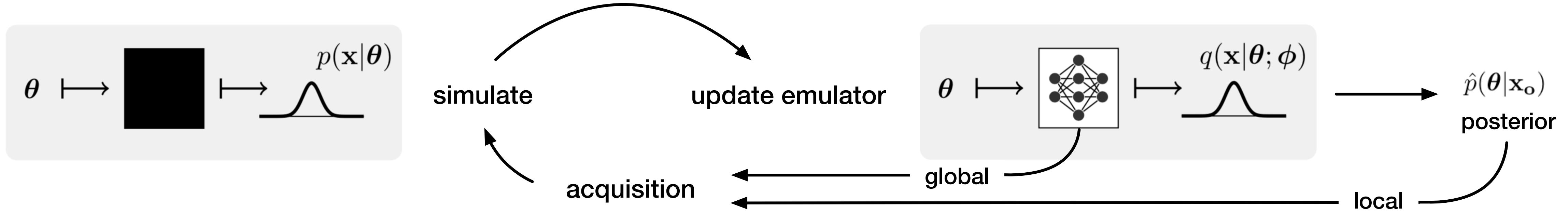


physics





Lueckmann et al (2018)



also see: Brehmer et al (2018): Mining gold from implicit models to improve likelihood-free inference,
 Brehmer et al (2018): A Guide to Constraining Effective Field Theories with Machine Learning

Probabilistic Programming

pymc3

```
# setup model
model = pm.Model()
with model:
    # cluster sizes
    p = pm.Dirichlet('p', a=np.array([1., 1., 1.]), shape=k)
    # ensure all clusters have some points
    p_min_potential = pm.Potential('p_min_potential', tt.switch(
        tt.sum(p) - 1., 0., -1000))

    # cluster centers
    means = pm.Normal('means', mu=[0, 0, 0], sd=15, shape=k)
    # break symmetry
    order_means_potential = pm.Potential('order_means_potential',
                                          tt.switch(means[1]-means[0],
                                                    + tt.switch(means[2]-means[0],
                                                               -1000, 0),
                                                    0))

    # measurement error
    sd = pm.Uniform('sd', lower=0, upper=20)

    # latent cluster of each observation
    category = pm.Categorical('category',
                               p=p,
                               shape=ndata)

    # likelihood for each observed value
    points = pm.Normal('obs',
                        mu=means[category],
                        sd=sd,
                        observed=data)
```

```
# fit model
with model:
    step1 = pm.Metropolis(vars=[p, sd, means])
    step2 = pm.ElemwiseCategorical(vars=[category], values=[0, 1, 2])
    tr = pm.sample(10000, step=[step1, step2])
```



also:

- STAN
- Edward
- Tensorflow Probability

**Exercise (if time): Write the chocolate
productivity hierarchical model from the
tutorial in a PyMC3 model**

Model Interpretation + Critiquing

(see also Blei (2014))

Assess the fitness of the model

- How good is your model in general?
- Where does it fail in particular?

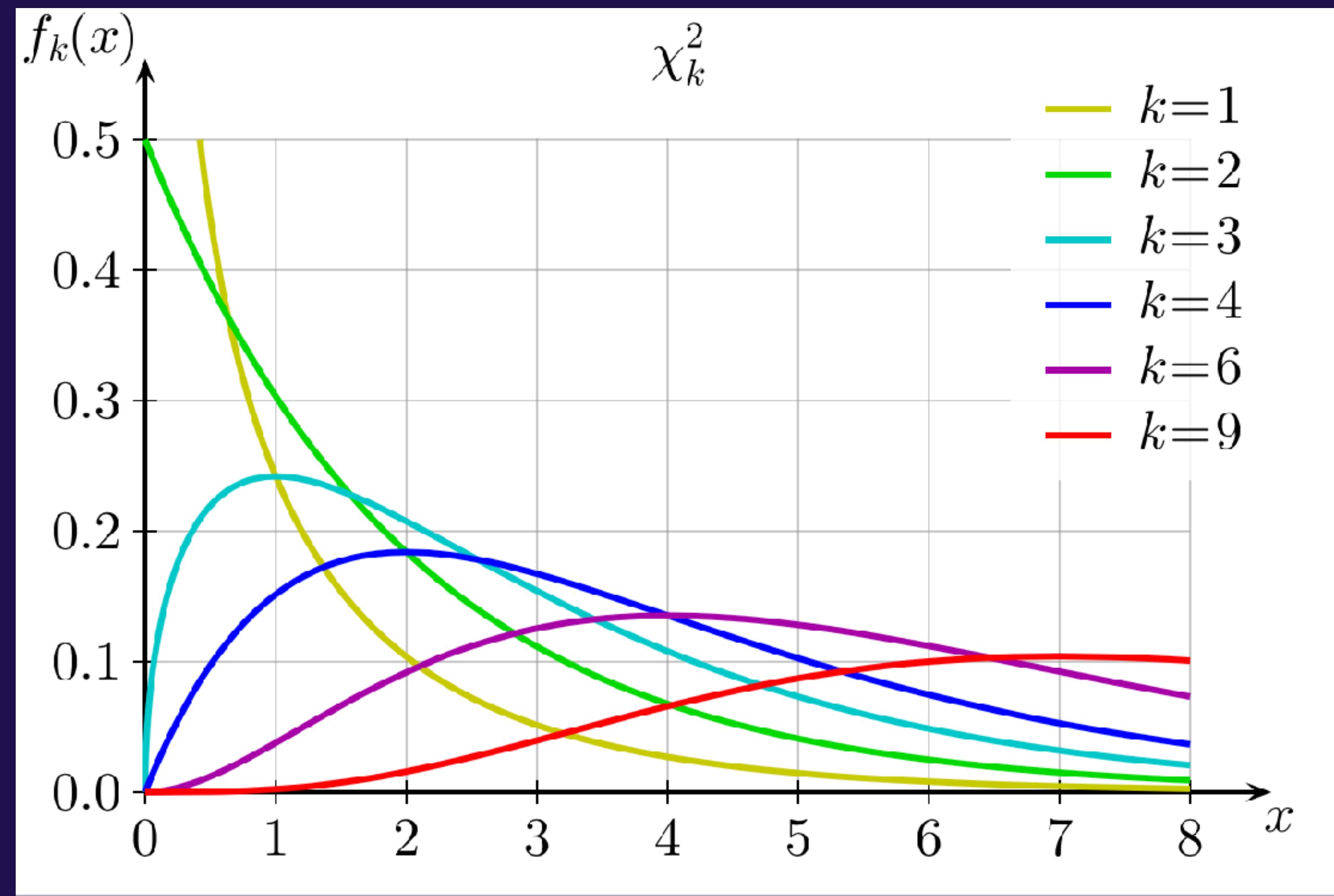
Be pragmatic about your models + inference

Goodness of Fit

$$\mathcal{L}(\theta) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - m(x_i, \theta))^2$$

distributed as χ^2 with $n-k$ degrees of freedom

Goodness of Fit



key concept: confront the model's **posterior**
predictive distribution with the **observed data**

Posterior Predictive Distribution

$$p(\theta \mid \mathbf{y}, I) \propto p(\mathbf{y} \mid \theta, I)p(\theta \mid I)$$

Posterior Predictive Distribution

$$p(\theta | \mathbf{y}, I) \propto p(\mathbf{y} | \theta, I)p(\theta | I)$$

$$p(y' | \mathbf{y}, I) = \int_{\Theta} p(y' | \theta, \mathbf{y}, I)p(\theta | \mathbf{y}, I)d\theta$$

Posterior Predictive Distribution

$$p(\theta | \mathbf{y}, I) \propto p(\mathbf{y} | \theta, I)p(\theta | I)$$

$$p(y' | \mathbf{y}, I) = \int_{\Theta} p(y' | \theta, \mathbf{y}, I)p(\theta | \mathbf{y}, I)d\theta$$

Predictive Sample Re-use (cross-validation)

- Evaluate model assumption
- can also compare different sampling algorithms!

Posterior Predictive Checks

Is my model good enough in ways that matter?

Posterior Predictive Checks

Posterior Predictive Checks

1. Draw data points from the posterior predictive distribution:

Posterior Predictive Checks

1. Draw data points from the posterior predictive distribution:
 - a. Draw parameters θ from $p(\theta | y)$

Posterior Predictive Checks

1. Draw data points from the posterior predictive distribution:
 - a. Draw parameters θ from $p(\theta | y)$
 - b. Draw replicated data y^{rep} from the sampled θ

Posterior Predictive Checks

1. Draw data points from the posterior predictive distribution:
 - a. Draw parameters θ from $p(\theta | y)$
 - b. Draw replicated data y^{rep} from the sampled θ
2. compute a summary statistic $T(y)$ that encapsulates some structure in the data:

Posterior Predictive Checks

1. Draw data points from the posterior predictive distribution:
 - a. Draw parameters θ from $p(\theta | y)$
 - b. Draw replicated data y^{rep} from the sampled θ
2. compute a summary statistic $T(y)$ that encapsulates some structure in the data: $\text{PPC} = P(T(y^{\text{rep}}) > T(y) | y))$

Posterior Predictive Checks

1. Draw data points from the posterior predictive distribution:
 - a. Draw parameters θ from $p(\theta | y)$
 - b. Draw replicated data y^{rep} from the sampled θ
2. compute a summary statistic $T(y)$ that encapsulates some structure in the data: $\text{PPC} = P(T(y^{\text{rep}}) > T(y) | y))$

$T(y)$ is adaptable to the needs of the practitioner!

Summary

Summary

- Be **principled**, but **pragmatic** about your models

Summary

- Be **principled**, but **pragmatic** about your models
- Start **simple!**

Summary

- Be **principled**, but **pragmatic** about your models
- Start **simple!**
- Explore **combinations of statistical inference and machine learning** (especially when you don't care about the parameters)

Summary

- Be **principled**, but **pragmatic** about your models
- Start **simple!**
- Explore **combinations of statistical inference and machine learning** (especially when you don't care about the parameters)
- **Critique** your models!