

Probabilistic Programming Languages

Guillaume Baudart, Hugo Paquet, Xavier Rival, Christine Tasson

MPRI 2025-2026

MPRI PPL

Course content

- Language design (probabilistic construction, static analysis, compilation)
- Inference methods (Monte Carlo methods, symbolic inference, variational inference)
- Semantics (measure theory, quasi-Borel spaces, integrable cones)
- Static analysis of inference algorithms (variational inference)

Evaluation: $\max(\text{exam}, (\text{project} + \text{exam}) / 2)$

- Article-based project
- Final exam



Guillaume Baudart
Inria / IRIF



Hugo Paquet
Inria / DIENS



Xavier Rival
Inria / DIENS



Christine Tasson
ISAE Supaero



<https://github.com/mpri-probprog/probprog-25-26>

Introduction

MPRI-PPL

Probabilistic programming languages

General purpose programming languages extended with probabilistic constructs

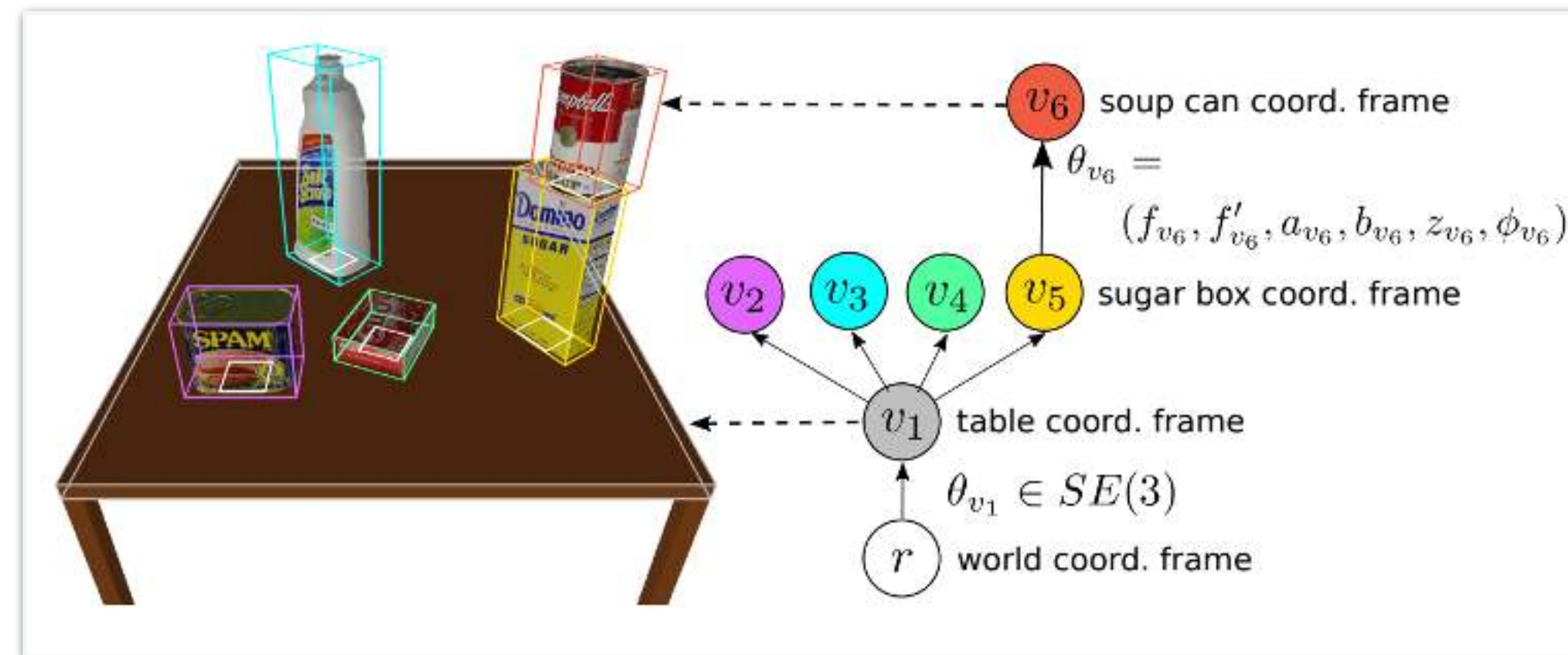
- **sample**: draw a sample from a distribution
- **assume**, **factor**, **observe**: condition the model on inputs (e.g., observed data)
- **infer**: compute the posterior distribution of a model given the inputs

Multiple examples:

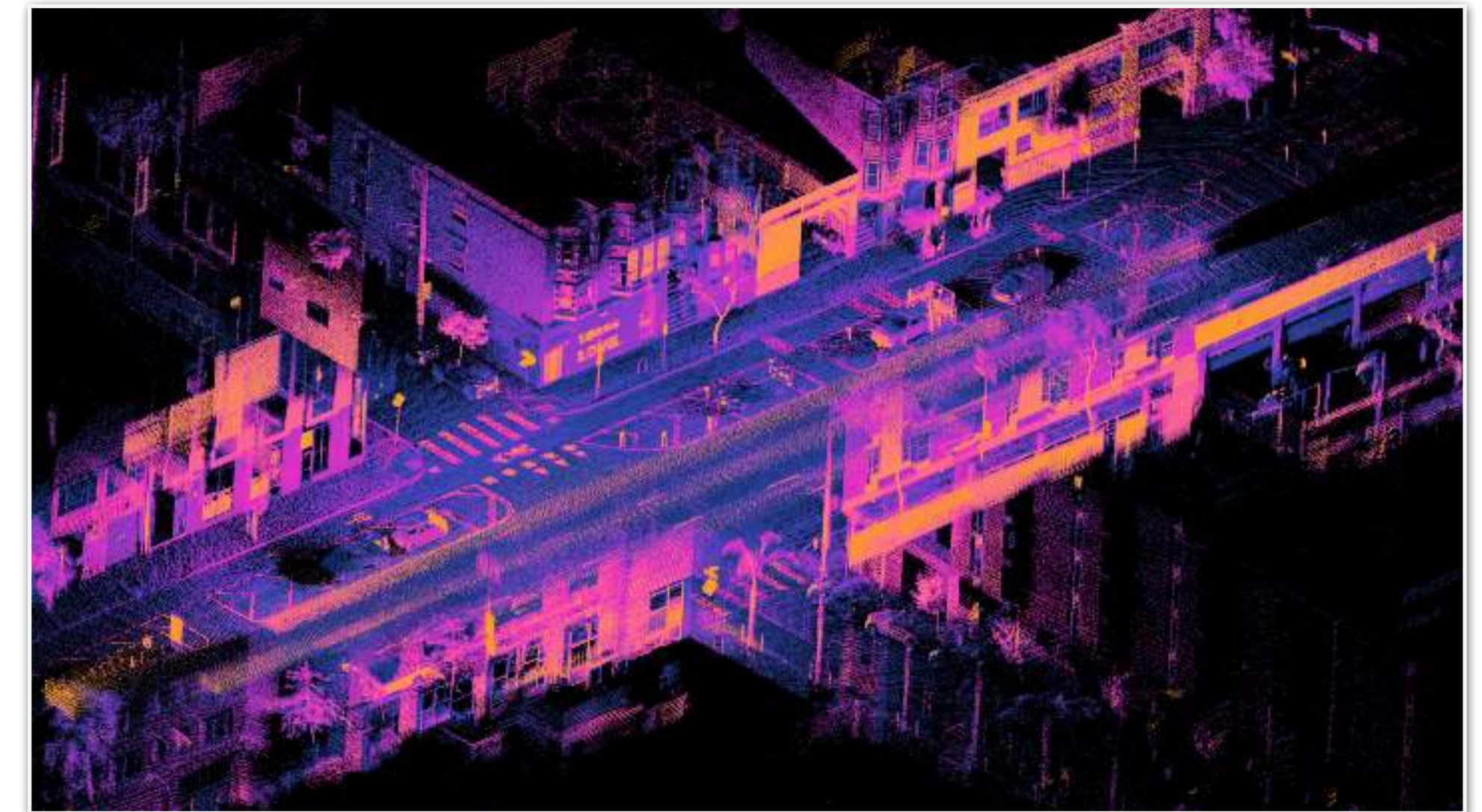
- Church, Anglican (lisp, clojure), 2008
- WebPPL (javascript), 2014
- Pyro/NumPyro (python), 2017/2019
- Gen (julia), 2018
- ProbZelus (Zelus), 2019
- ...

More and more, incorporating new ideas:

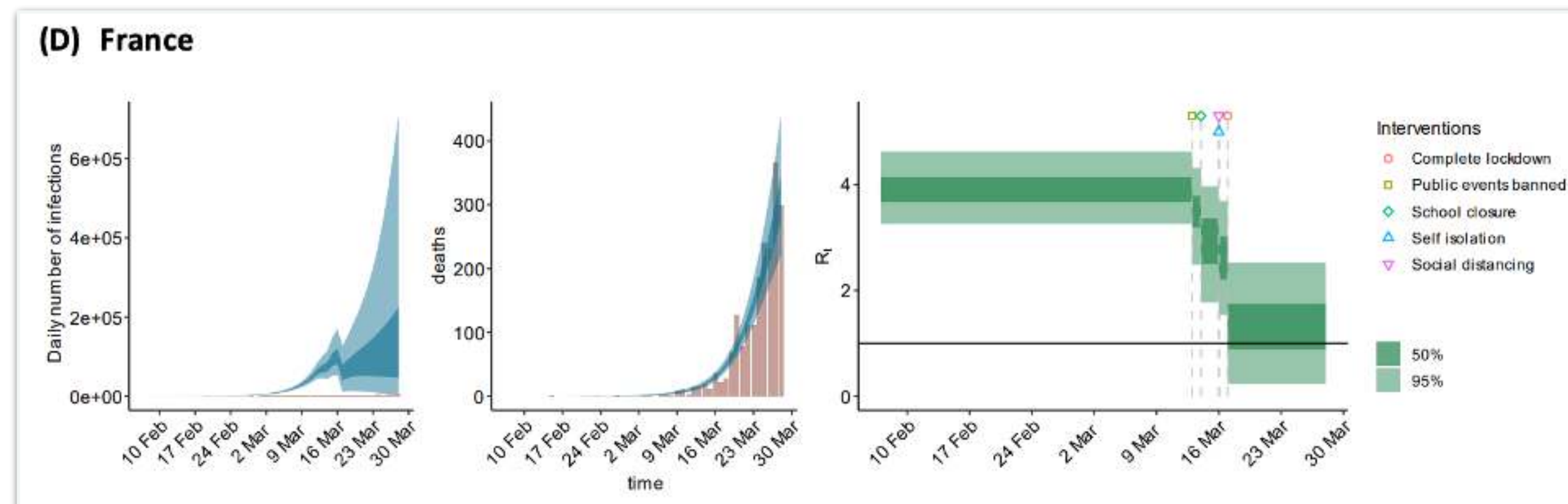
- New inference techniques, e.g., stochastic variational inference (SVI)
- Interaction with neural nets (deep probabilistic programming)



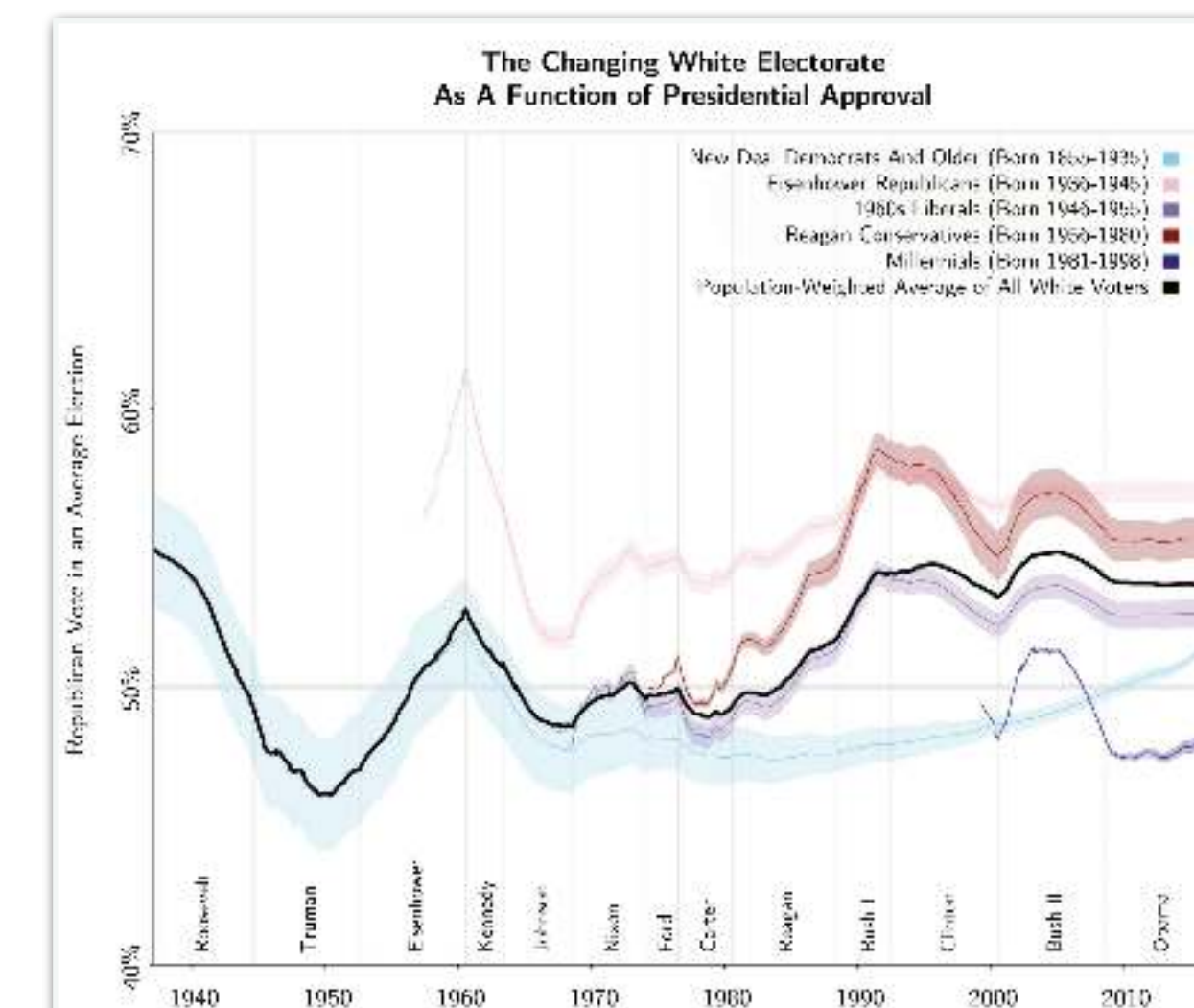
Gothoskar et al. *3dp3: 3d scene perception via probabilistic programming*
NeurIPS 2021



SLAM: Simultaneous localization and mapping
https://en.wikipedia.org/wiki/Simultaneous_localization_and_mapping#/media/File:Ouster_OS1-64_lidar_point_cloud_of_intersection_of_Folsom_and_Dore_St,_San_Francisco.png



Report 13: Estimating the number of infections and the impact of non-pharmaceutical interventions on COVID-19 in 11 European countries, Imperial College COVID-19 Response Team, March 2020



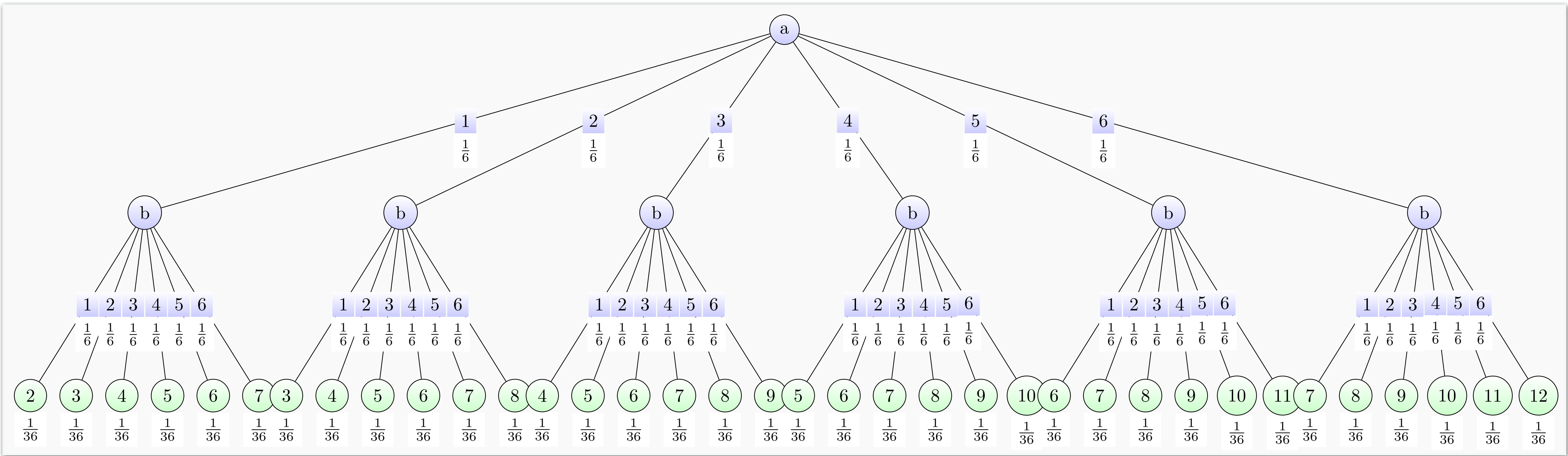
Ghitza et al. The Great Society, Reagan's Revolution, and Generations of Presidential Voting AJPS 2022

Examples: $\mu\text{-ppl}$

What is this?

```
def dice() → int:  
    a = sample(RandInt(1, 6), name="a")  
    b = sample(RandInt(1, 6), name="b")  
    return a + b
```

The sum of the values output by the two independent fair dice.



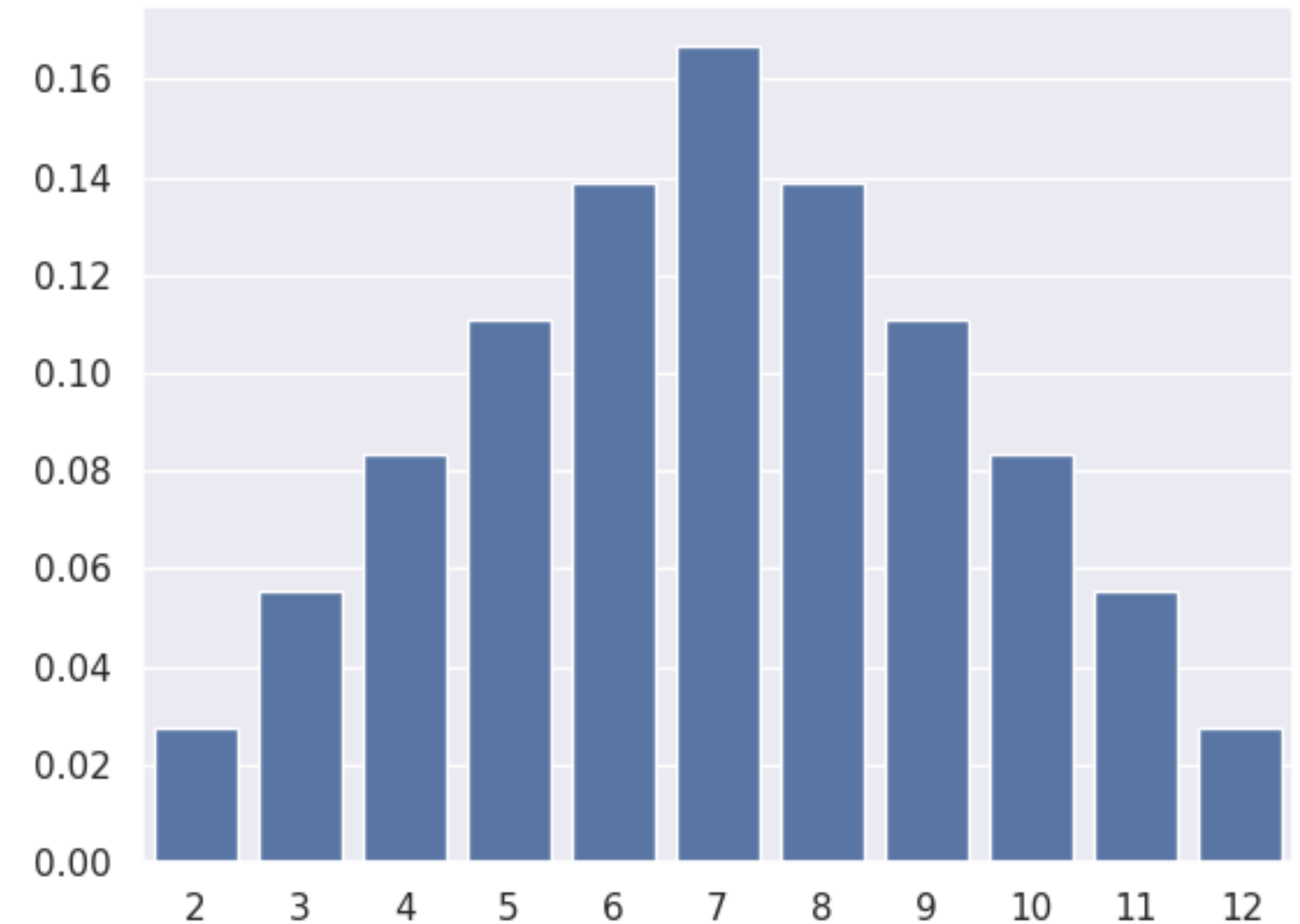
What is this?

```
def dice() → int:
    a = sample(RandInt(1, 6), name="a")
    b = sample(RandInt(1, 6), name="b")
    return a + b

with Enumeration():
    dist: Categorical[int] = infer(dice)
    viz(dice)
```

$$\begin{aligned} \llbracket \text{dice} \rrbracket &:= \mathbb{P}(\text{dice}() = k) \\ &= \sum_{a=1}^6 \sum_{b=1}^6 \frac{1}{36} \mathbb{1}_{\{a+b=k\}} \end{aligned}$$

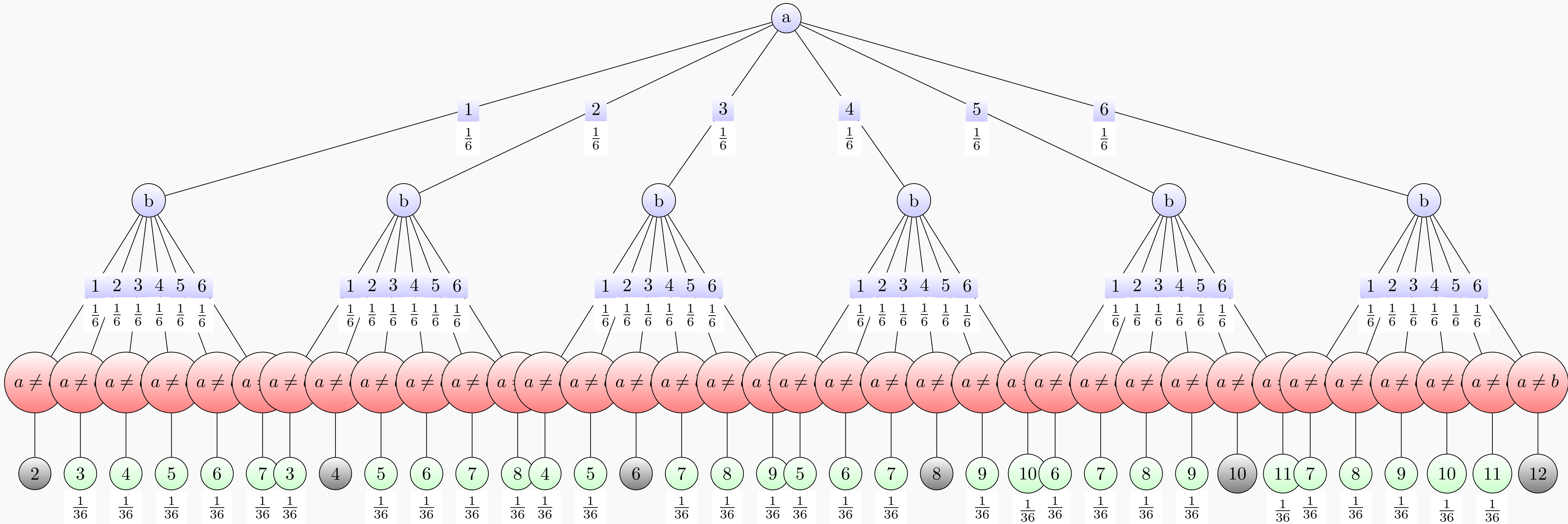
The sum of the values output by the two independent fair dice.



What is this?

```
def hard_dice() → int:  
    a = sample(RandInt(1, 6), name="a")  
    b = sample(RandInt(1, 6), name="b")  
    assume (a ≠ b)  
    return a + b
```

*The sum of two independent fair dice,
given that the output value are distinct.*



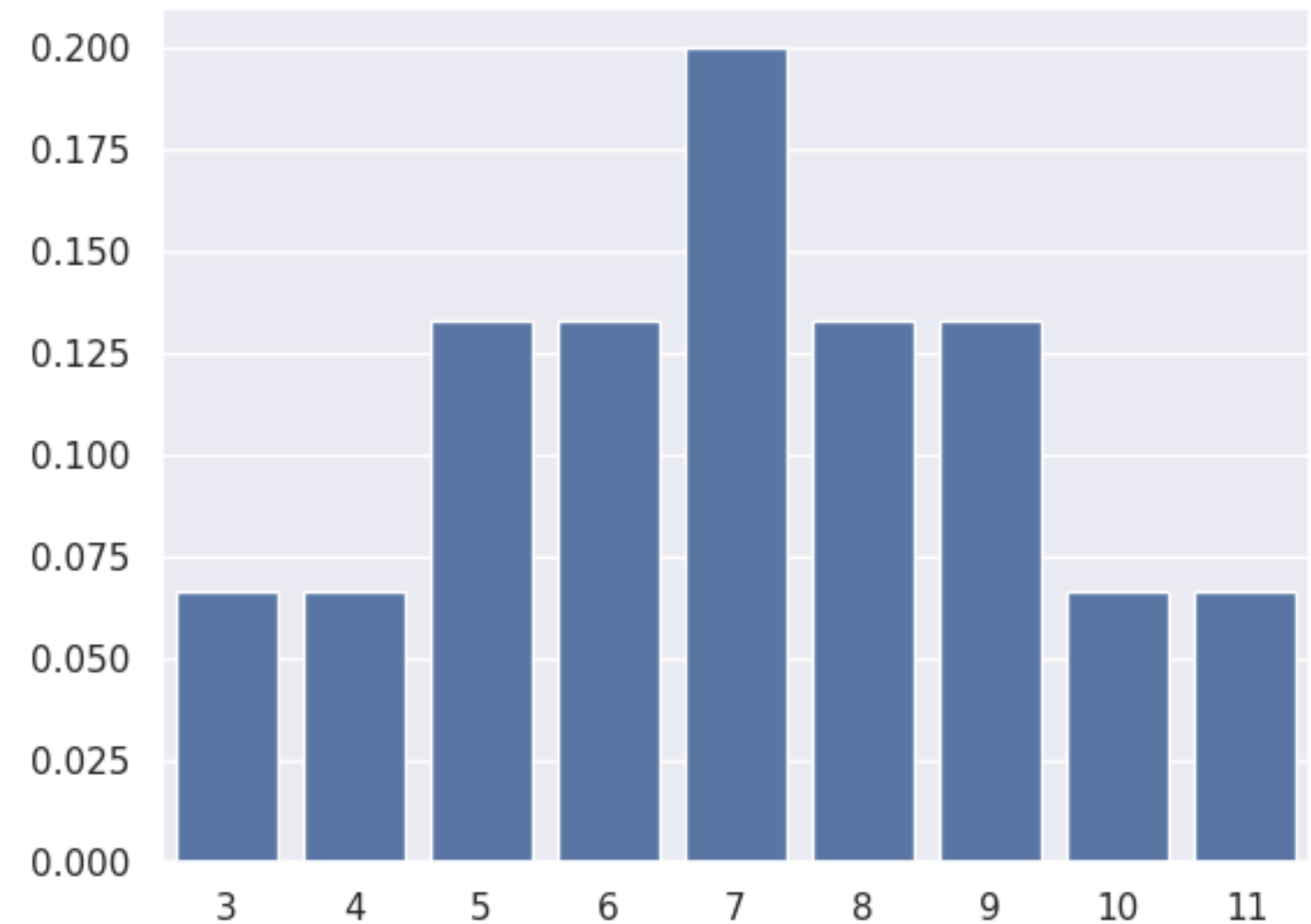
What is this?

```
def hard_dice() → int:
    a = sample(RandInt(1, 6), name="a")
    b = sample(RandInt(1, 6), name="b")
    assume (a ≠ b)
    return a + b

with Enumeration():
    dist: Categorical[int] = infer(dice)
    viz(dice)
```

$$\begin{aligned} \llbracket \text{hard_dice} \rrbracket (k) &= \mathbb{P}(a + b = k | a \neq b) \\ &= \frac{\sum_{a \neq b \in \{1, \dots, 6\}^2} \frac{1}{36} \mathbb{1}_{\{a+b=k\}}}{\sum_{a \neq b \in \{1, \dots, 6\}^2} \frac{1}{36}} \end{aligned}$$

The sum of two independent fair dice, given that the output value are distinct.



Bayesian programming

MPRI-PPL

Probabilistic programming

Bayesian Inference: learn parameters from data

- Latent parameter x
- Observed data y_1, \dots, y_n

$$p(x \mid y_1, \dots, y_n) = \frac{p(x) p(y_1, \dots, y_n \mid x)}{p(y_1, \dots, y_n)} \quad (\text{Bayes' theorem})$$

posterior

$$\propto p(x) p(y_1, \dots, y_n \mid x) \quad (\text{Data are constants})$$

prior

likelihood

Probabilistic constructs

- $x = \text{sample}(d)$: introduce a random variable x of distribution d
- $\text{observe}(d, y)$: condition on the fact that y was sampled from d
- $\text{infer}(m, y)$: compute posterior distribution of m given y

Notation: $x \sim \mu$

- parameter (output): `sample(mu)`
- observation (input): `observe(mu)`



Thomas Bayes (1701-1761)

Probabilistic programming

Bayesian Inference: learn parameters from data

- Latent parameter x
- Observed data y_1, \dots, y_n

$$p(x \mid y_1, \dots, y_n) = \frac{p(x) p(y_1, \dots, y_n \mid x)}{p(y_1, \dots, y_n)} \quad (\text{Bayes' theorem})$$

posterior

$$\propto p(x) p(y_1, \dots, y_n \mid x) \quad (\text{Data are constants})$$

prior

likelihood

```
def model(y1, ..., yn):  
    x = sample prior  
    observe (likelihood(x), (y1, ..., yn))  
    return x
```

```
infer(model, (y1, ..., yn))
```

A Bayesian model is described by a program



Thomas Bayes (1701-1761)

Example: coin

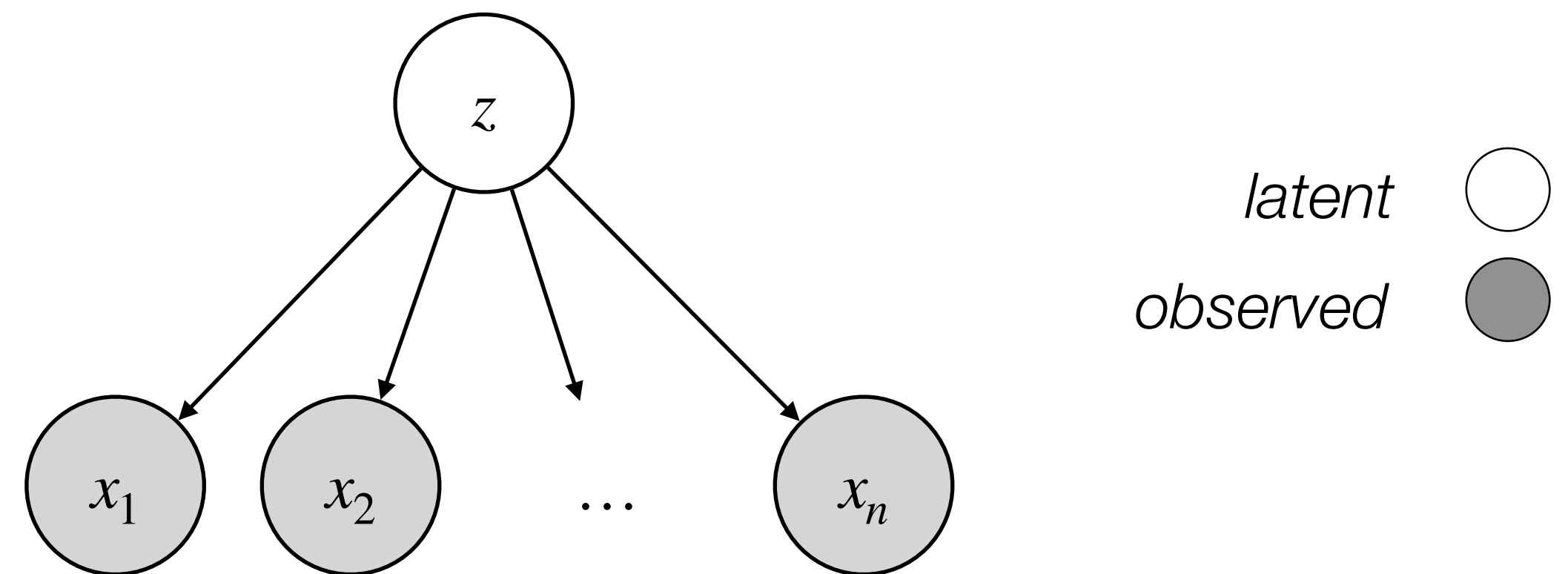
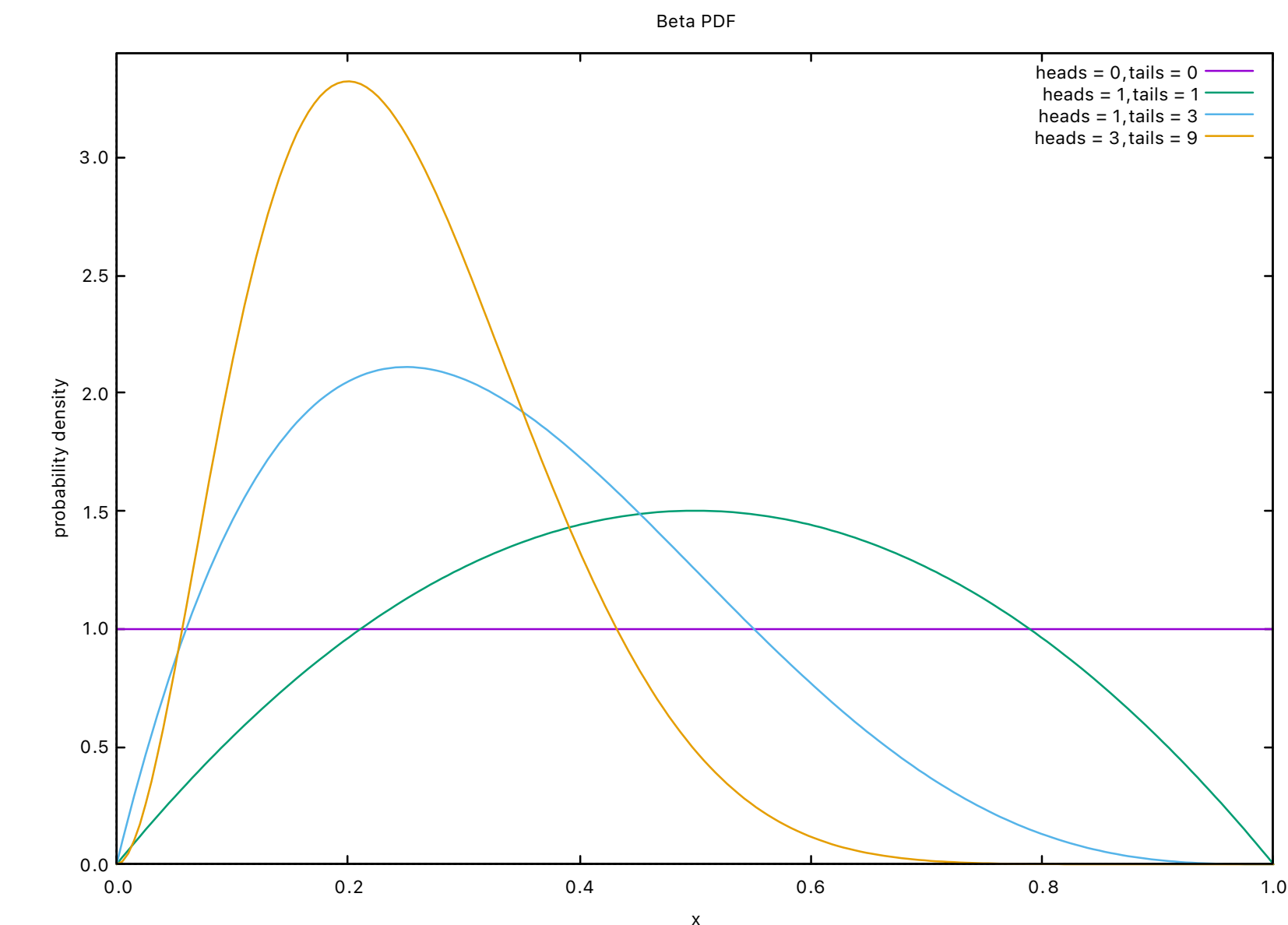


Consider a series of coin tosses

- Observations: head or tail
- Each toss is independant
- What is the probability of getting head at the next toss?

Probabilistic model

- Prior: $z \sim \text{Uniform}(0, 1)$
- Observations: for $i \in [1, n]$, $x_i \sim \text{Bernoulli}(z)$
- Posterior: $p(z \mid x_1, \dots, x_n)$?



Example: coin



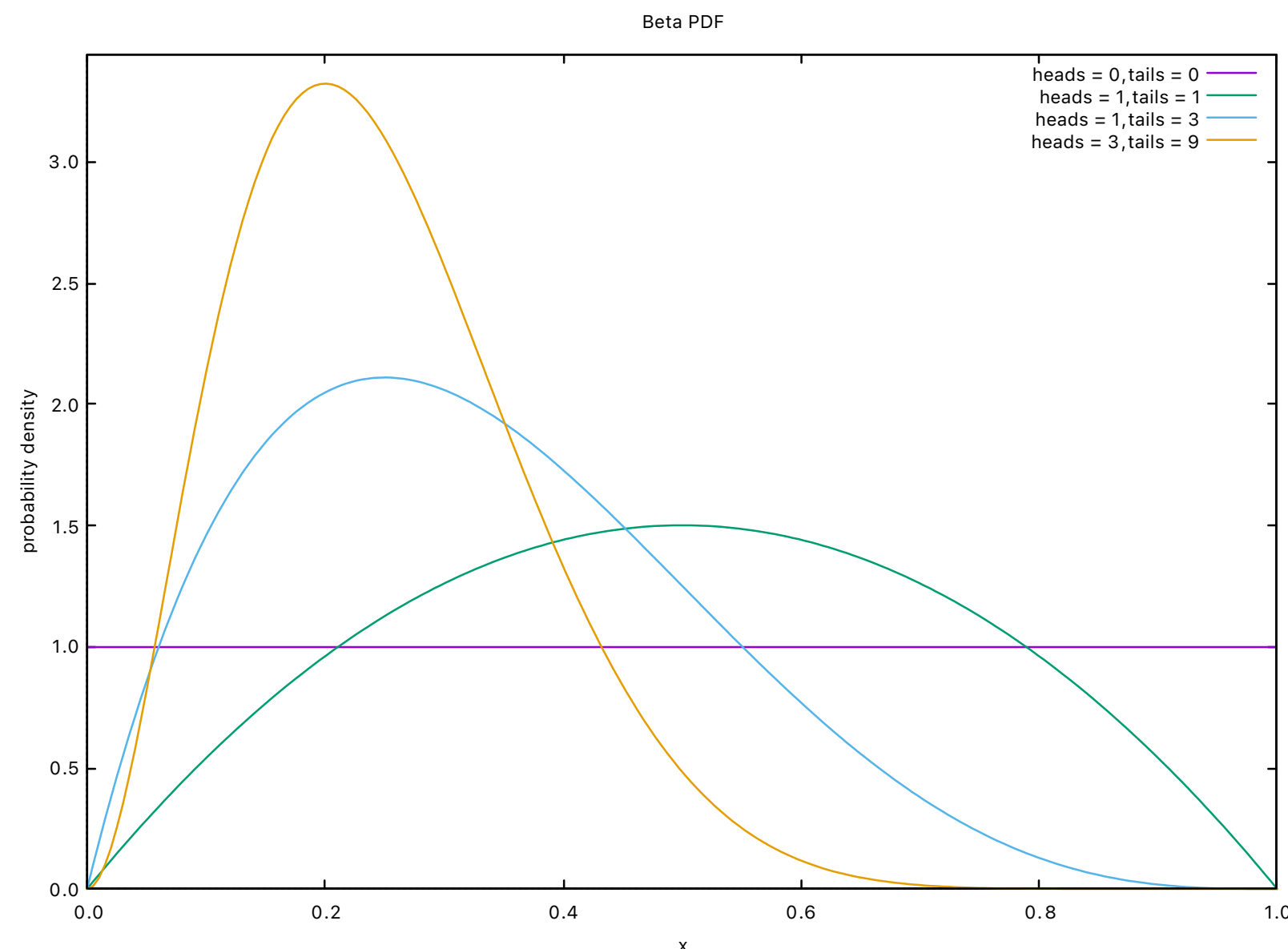
Consider a series of coin tosses

- Observations: head or tail
- Each toss is independant
- What is the probability of getting head at the next toss?

Probabilistic model

- Prior: $z \sim \text{Uniform}(0, 1)$
- Observations: for $i \in [1, n]$, $x_i \sim \text{Bernoulli}(z)$
- Posterior: $p(z \mid x_1, \dots, x_n)$?

```
def coin(obs: List[int]) → float:  
    p = sample(Uniform(0, 1))  
    for o in obs:  
        observe(Bernoulli(p), o)  
    return p
```



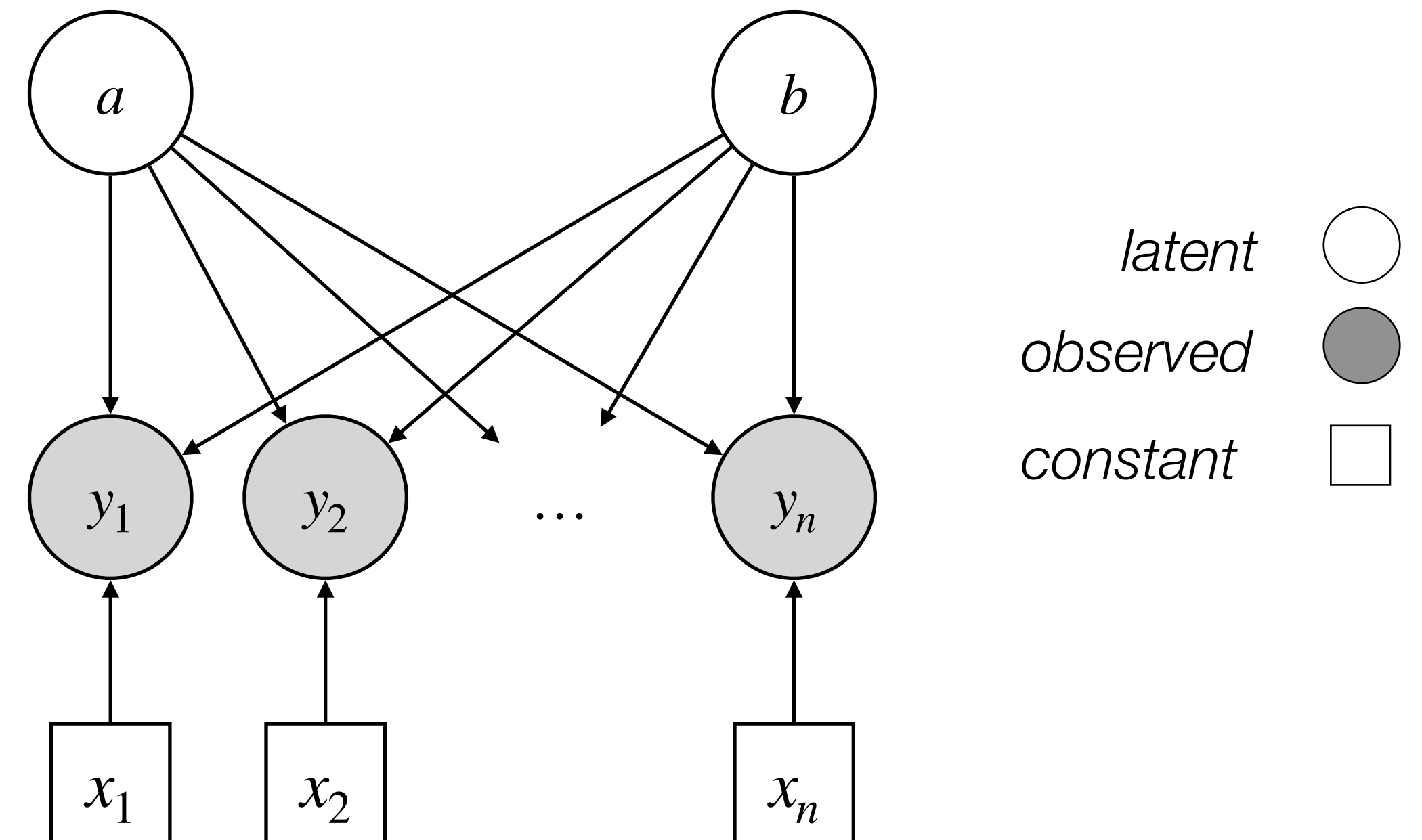
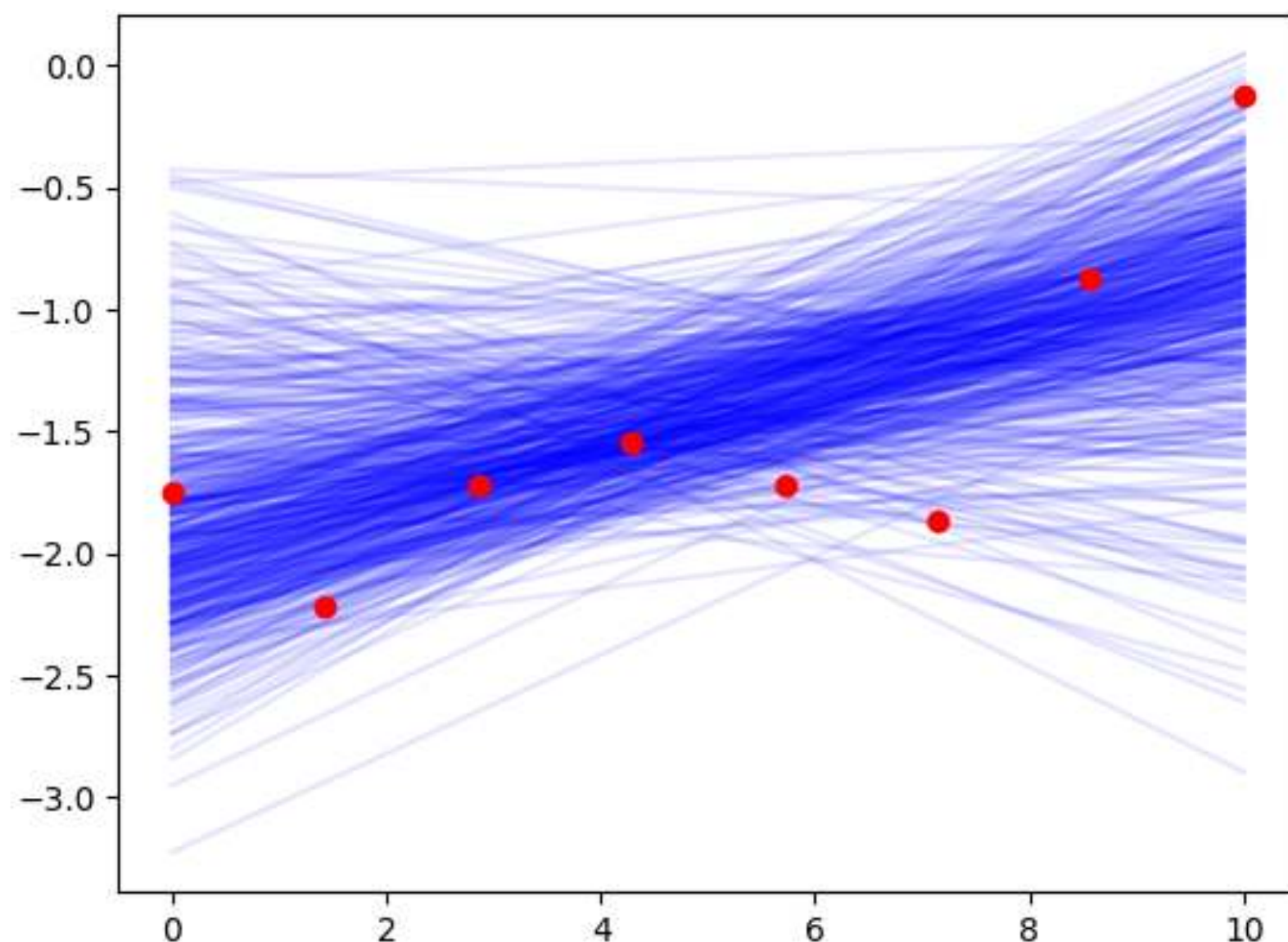
Example: linear regression

Consider a series of observations

- Observation: point (x, y)
- Each point is independent from others
- Find the distribution of possible regressions

Probabilistic model

- Prior: $a \sim \mathcal{N}(0, 1)$, and $b \sim \mathcal{N}(0, 1)$
- Observations: for $i \in [1, n]$, $y_i \sim \mathcal{N}(a \times x_i + b, \sigma)$
- Posterior: $p(a, b \mid (x_1, y_1) \dots, (x_n, y_n))$?



Example: linear regression

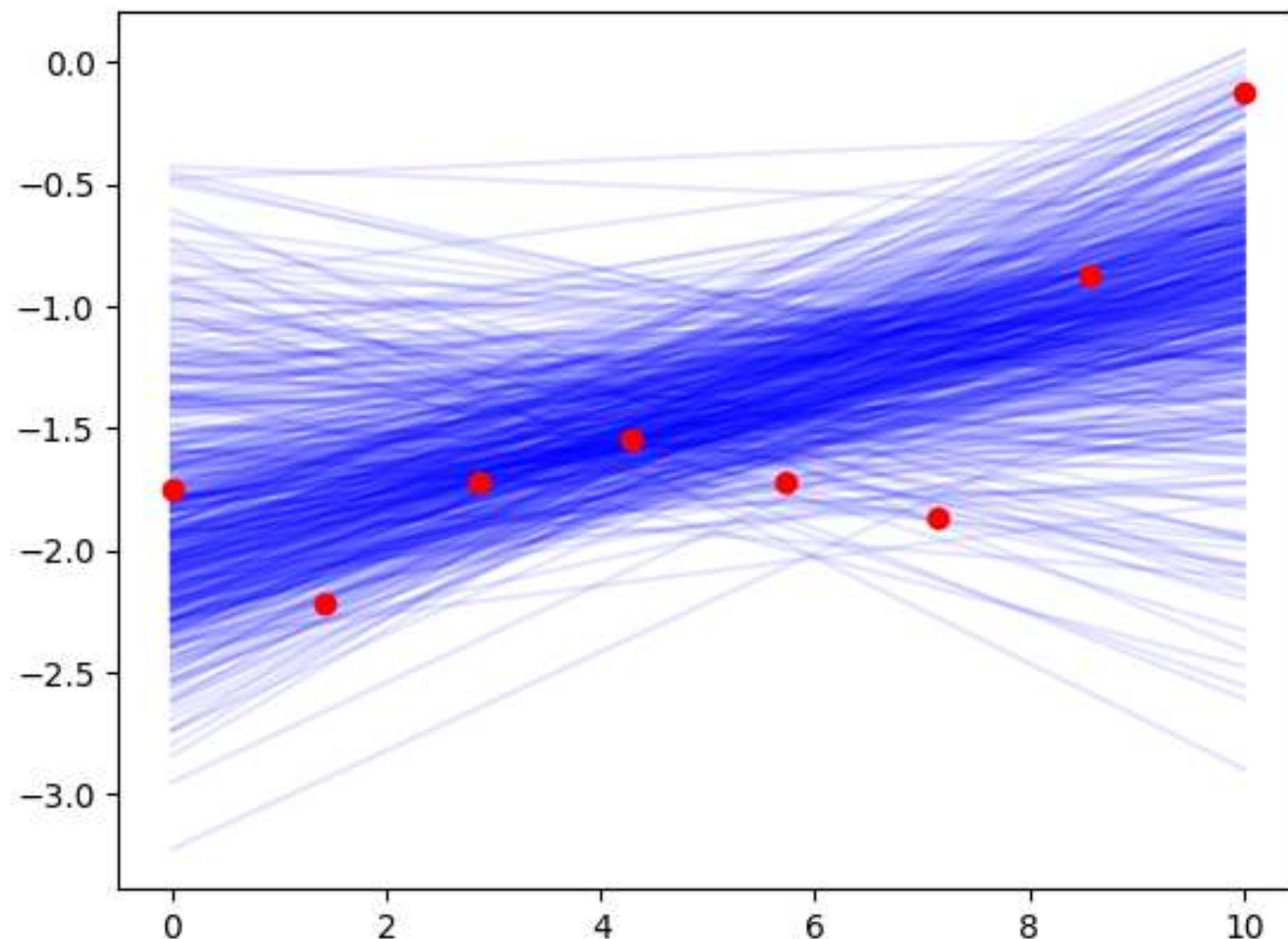
Consider a series of observations

- Observation: point (x, y)
- Each point is independent from others
- Find the distribution of possible regressions

Probabilistic model

- Prior: $a \sim \mathcal{N}(0, 1)$, and $b \sim \mathcal{N}(0, 1)$
- Observations: for $i \in [1, n]$, $y_i \sim \mathcal{N}(a \times x_i + b, \sigma)$
- Posterior: $p(a, b \mid (x_1, y_1) \dots, (x_n, y_n))$?

```
def model(x_obs, y_obs):  
    a = sample(Gaussian(0, 5))  
    b = sample(Gaussian(0, 5))  
    sigma = sample(Uniform(0, 2))  
    for (xo, yo) in zip(x_obs, y_obs):  
        observe(Gaussian(a*xo + b, sigma), yo)  
    return [a, b, sigma]
```



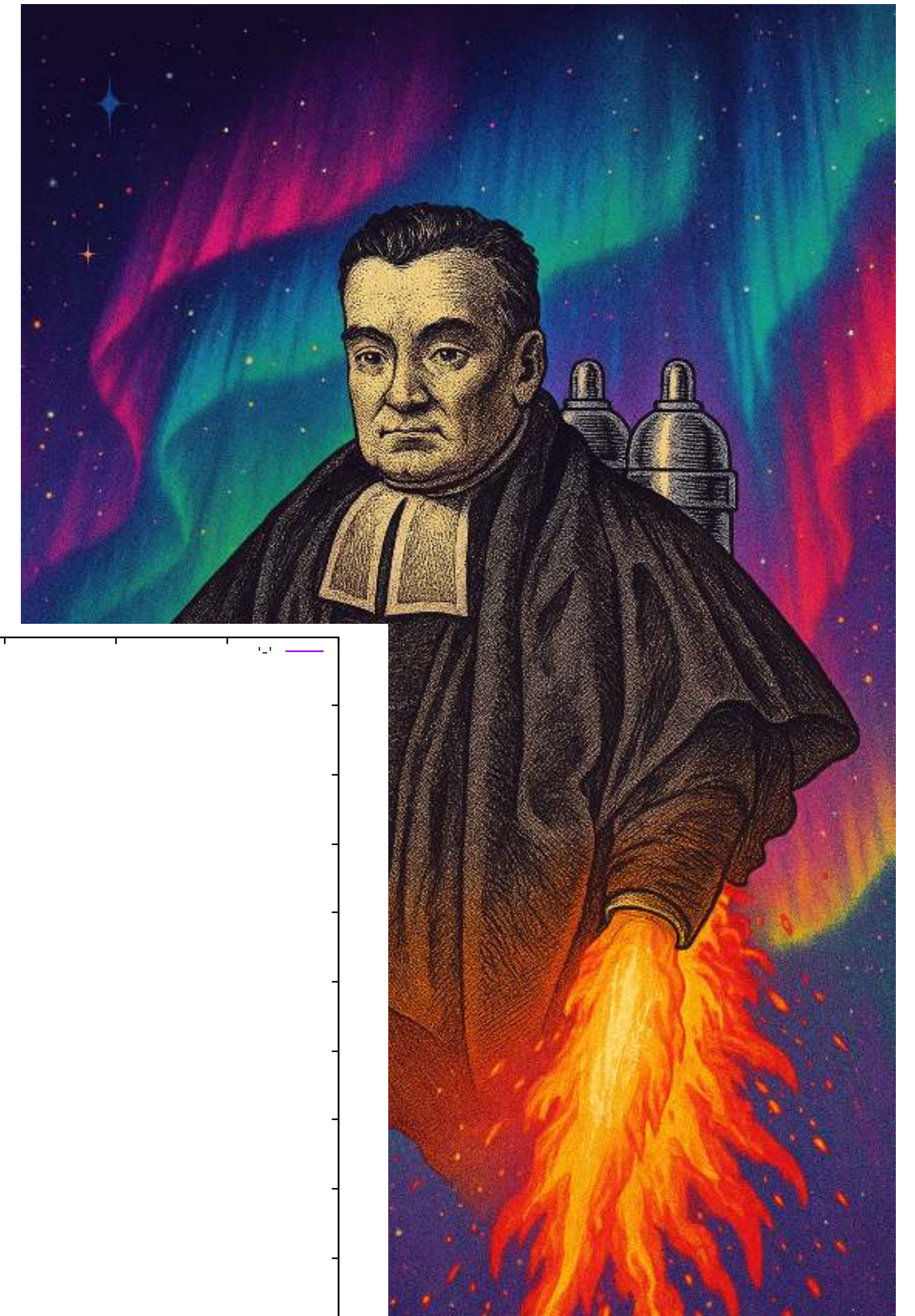
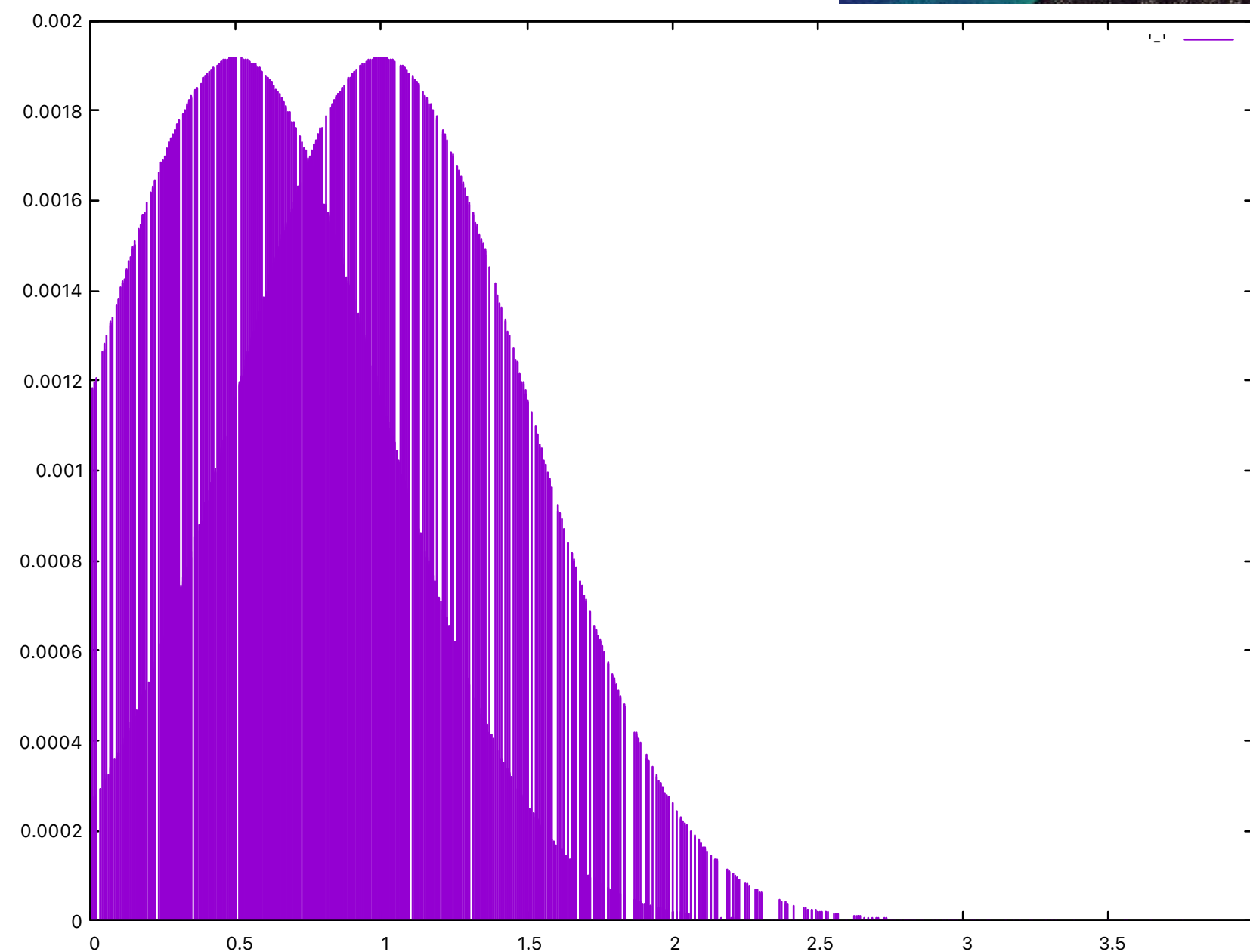
Probabilistic programming

Probabilistic constructs

- **sample**: draw a sample from a distribution
- **assume**, **factor**, **observe**: condition the model on observed data
- **infer**: compute the posterior distribution of a model given the inputs

More general than classic Bayesian reasoning

```
def weird() =  
  b = sample(Bernoulli(0.5))  
  mu = 0.5 if (b = 1) else 1.0  
  theta = sample(Gaussian(mu, 1.0))  
  if theta > 0.:  
    observe (Gaussian(mu, 0.5), theta)  
    return theta  
  else:  
    return weird ()
```



Probabilistic programming

Probabilistic constructs

- **sample**: draw a sample from a distribution
- **assume**, **factor**, **observe**: condition the model on observed data
- **infer**: compute the posterior distribution of a model given the inputs

More general than classic Bayesian reasoning

But: general case....

$$p(\theta \mid x_1, \dots, x_n) = \frac{p(\theta) p(x_1, \dots, x_n \mid \theta)}{p(x_1, \dots, x_n)} \quad (\text{Bayes' theorem})$$

$$= \frac{p(\theta) p(x_1, \dots, x_n \mid \theta)}{\int p(\theta) p(x_1, \dots, x_n \mid \theta) d\theta} = \frac{\text{🎂}}{\text{💀} \text{ intractable...}}$$

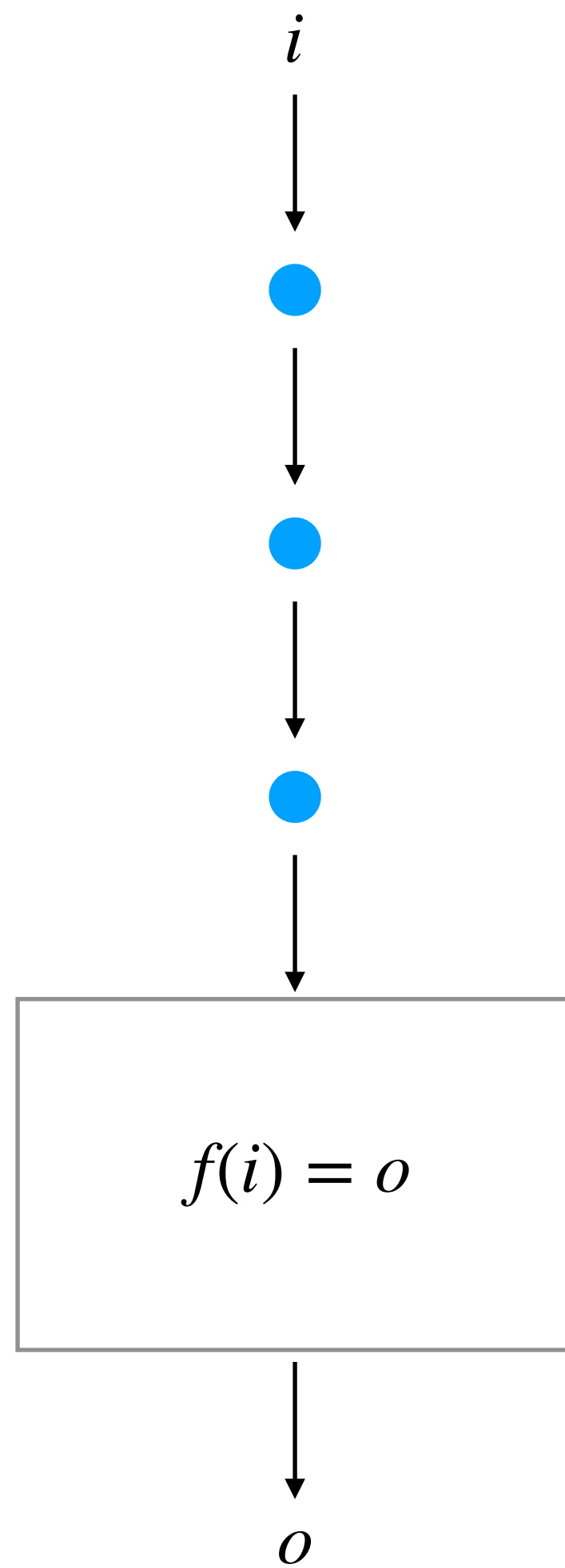


Inference

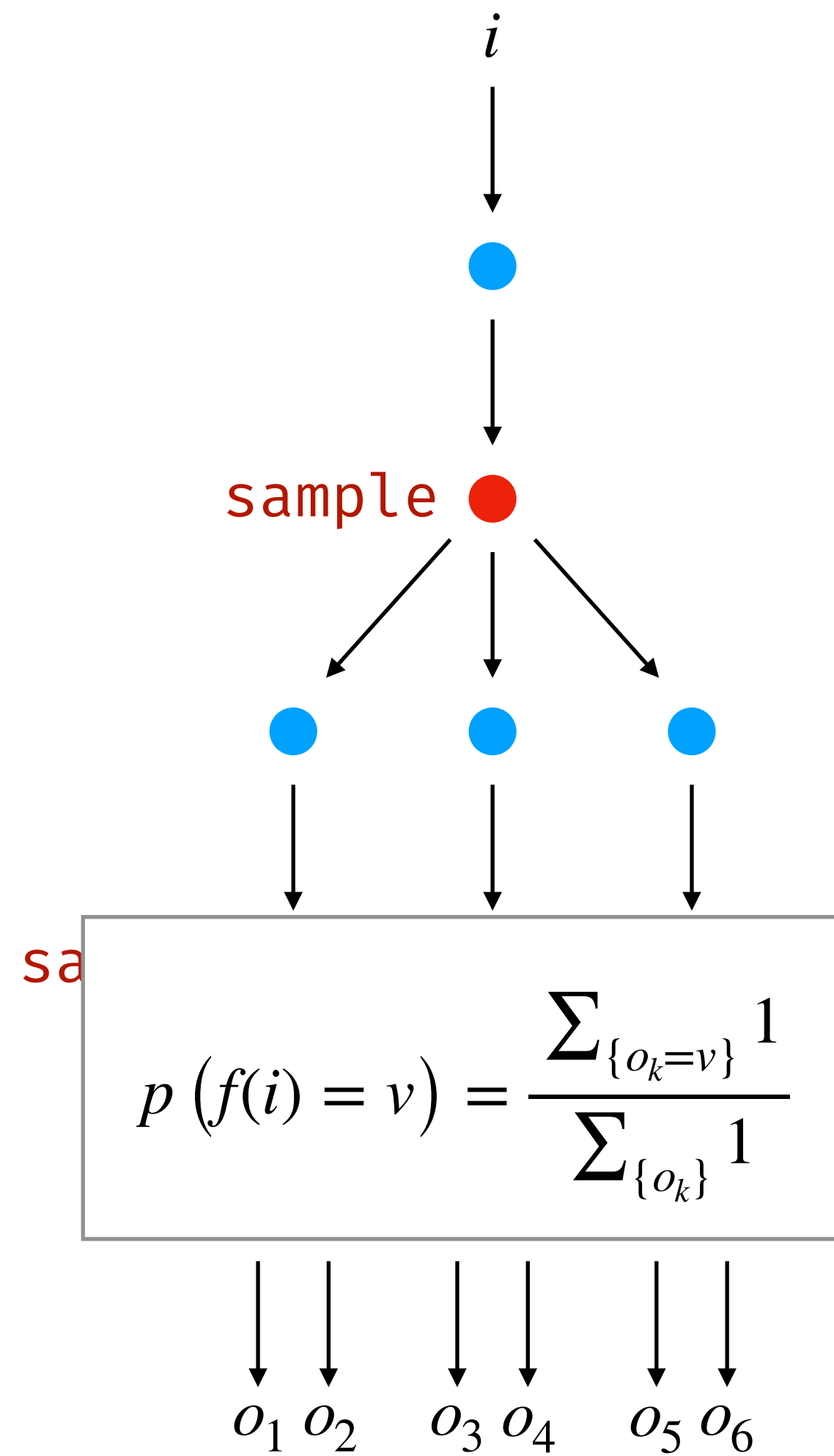
Probabilistic Programming Languages

infer : $(\alpha \rightarrow \beta) \rightarrow \alpha \rightarrow \beta$ dist

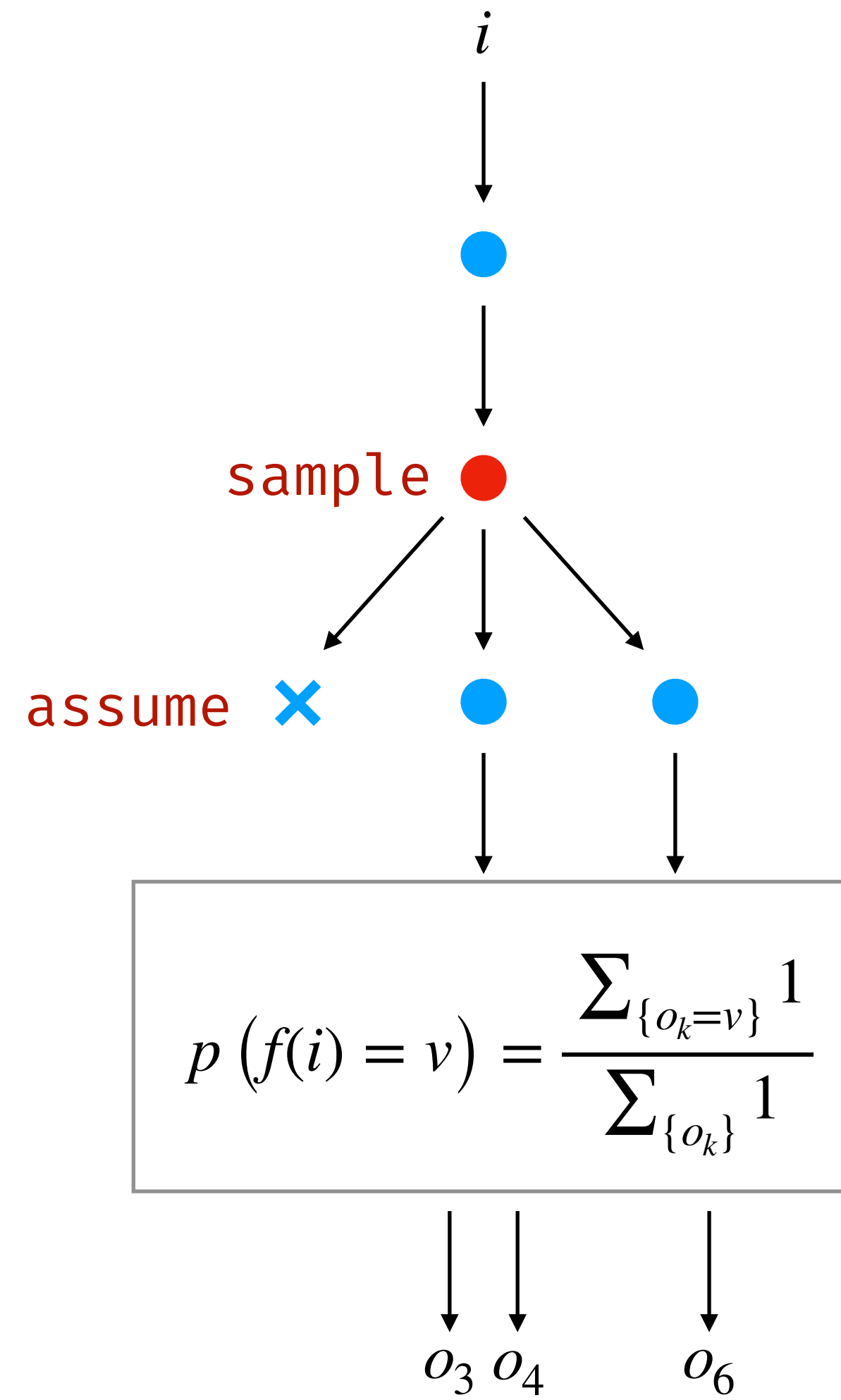
program



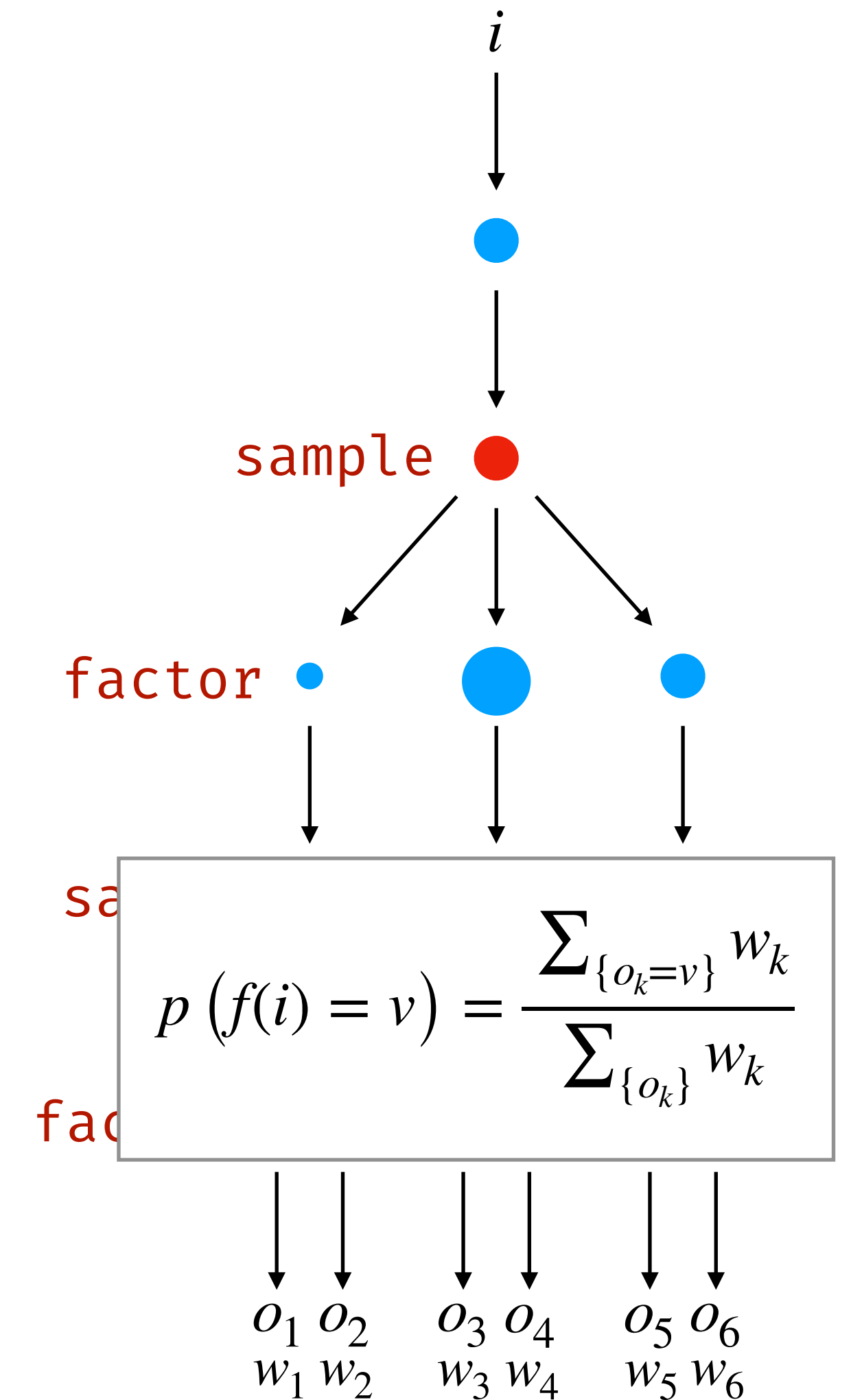
sample



assume



factor



Hard/soft conditioning

```
# Reject if [p] is not true.
def assume(self, p):
    if not p: self.score += -np.inf

# Assume [x] was sampled from [d].
def observe(self, d, x):
    v = self.sample(d)
    self.assume(v == x)
```

Hard conditioning

```
# Update the (log)score.
def factor(self, s):
    self.score += s

# Assume [x] was sampled from [d].
def observe(self, d, x):
    self.factor(d.logprob(x))
```

Soft conditioning

Importance sampling

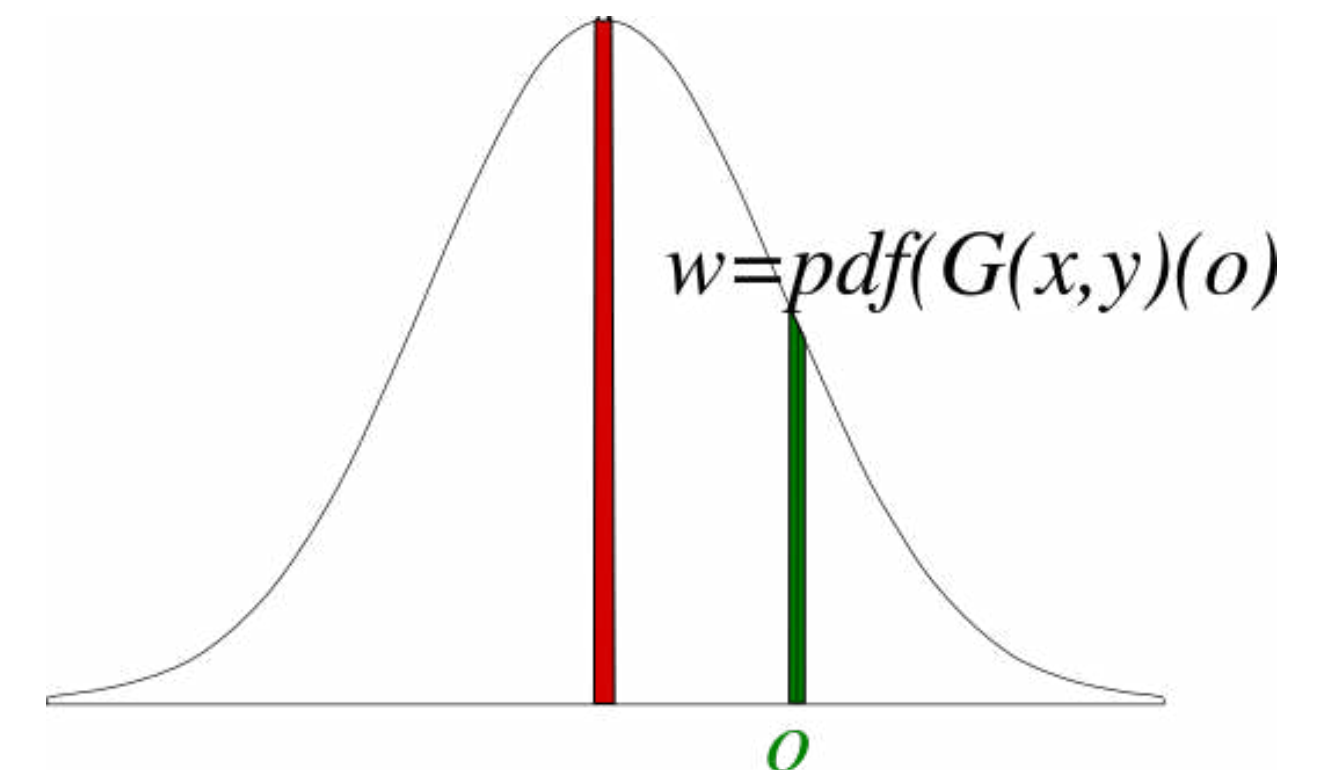
Inference algorithm

- Run a set of n independent executions
- **sample**: draw a sample from a distribution
- **factor**: associate a score to the current execution
- **infer**: gather output values v_i and score W_i to approximate the posterior distribution

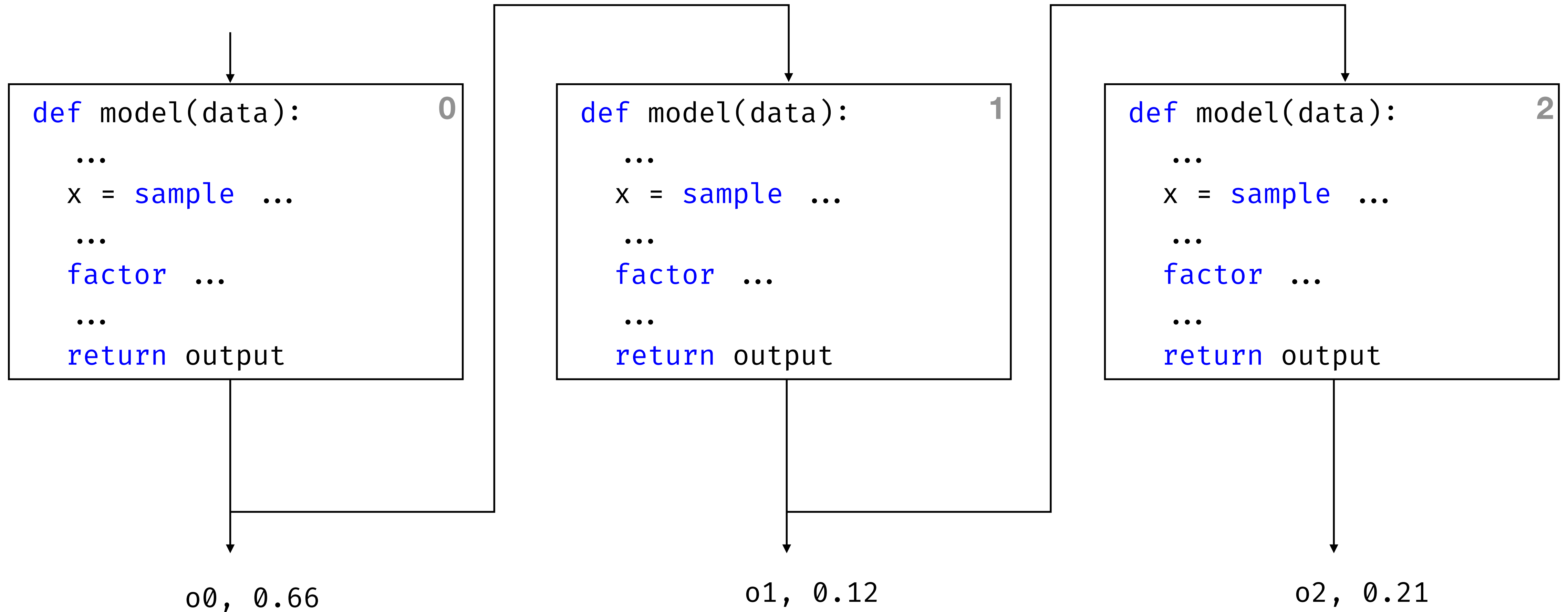
$$p_i = \frac{W_i}{\sum_{1 \leq i \leq n} W_i}$$

Conditioning

- **assume**(p): sets the score to 0 if p is false
- **observe**(d, v): multiply the score by the likelihood of d at v (density function of d)



Importance sampling

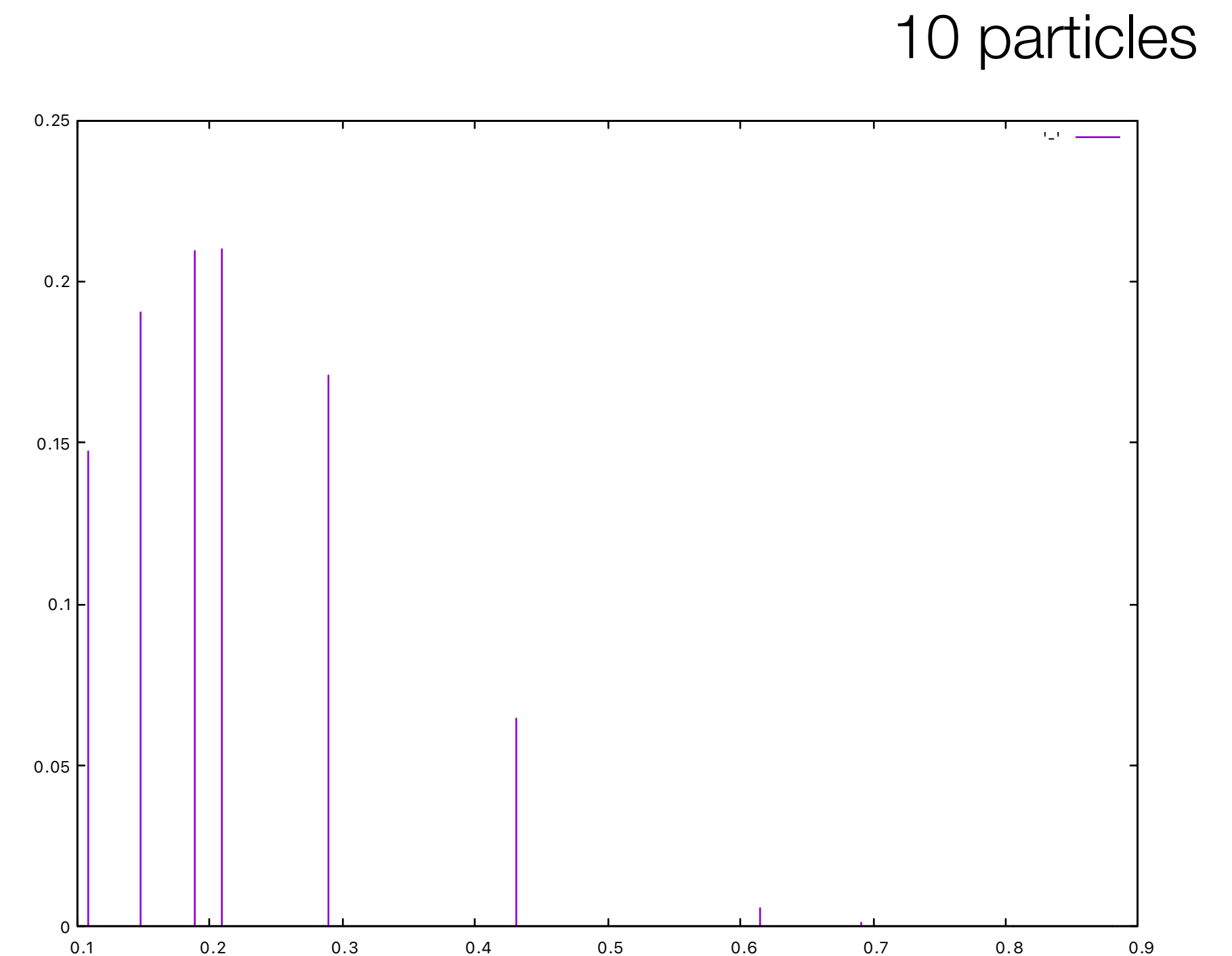


Example: coin



```
def coin(obs: list[int]) → float:
    p = sample(Uniform(0, 1))
    for o in obs:
        observe(Bernoulli(p), o)
    return p

with ImportanceSampling(num_particles=10):
    dist = infer(coin, [0, 0, 0, 0, 0, 0, 0, 0, 1, 1])
    viz(dist)
```



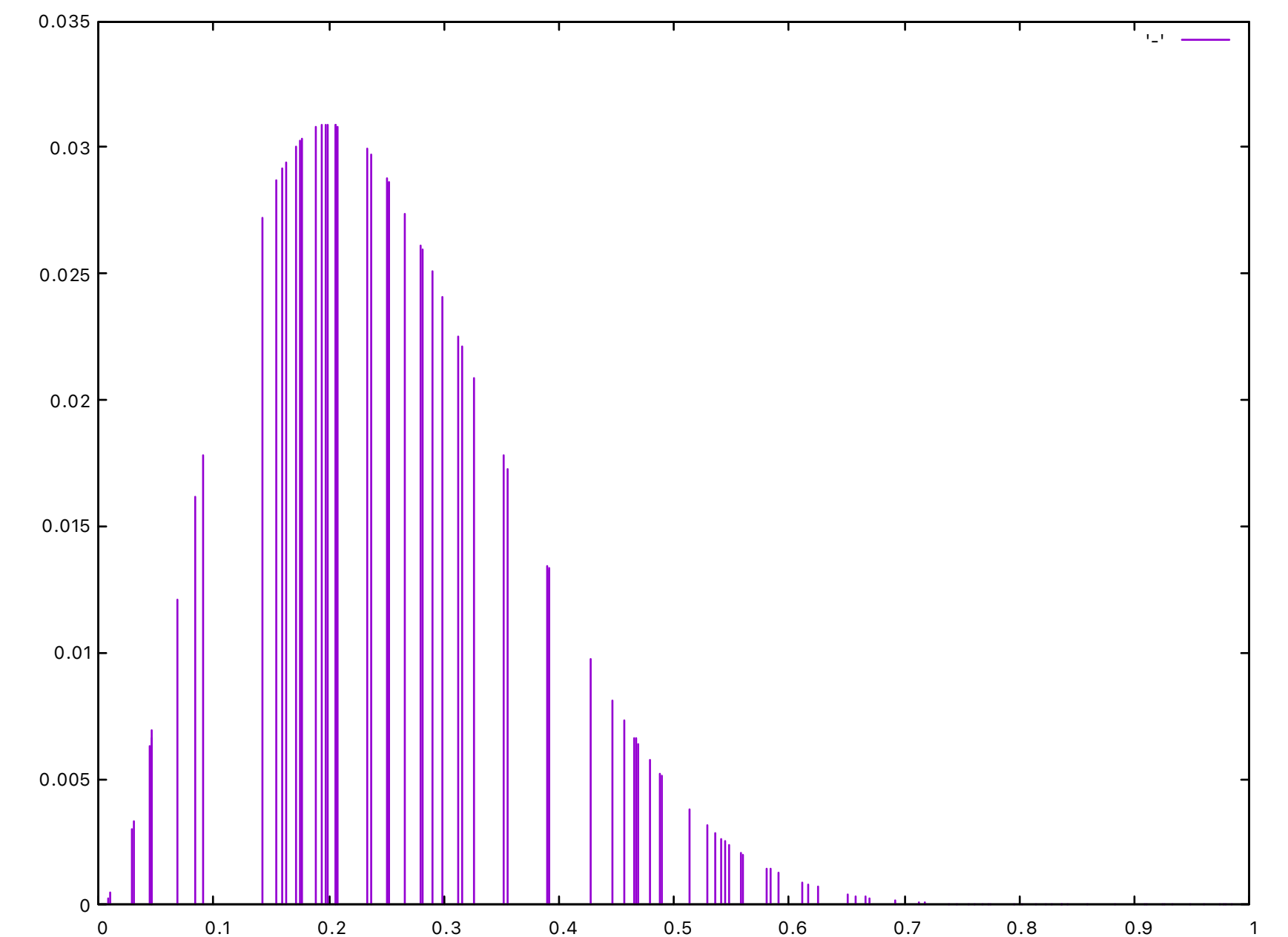
Example: coin



```
def coin(obs: list[int]) → float:  
    p = sample(Uniform(0, 1))  
    for o in obs:  
        observe(Bernoulli(p), o)  
    return p
```

```
with ImportanceSampling(num_particles=100):  
    dist = infer(coin, [0, 0, 0, 0, 0, 0, 0, 0, 1, 1])  
    viz(dist)
```

100 particles



Example: coin

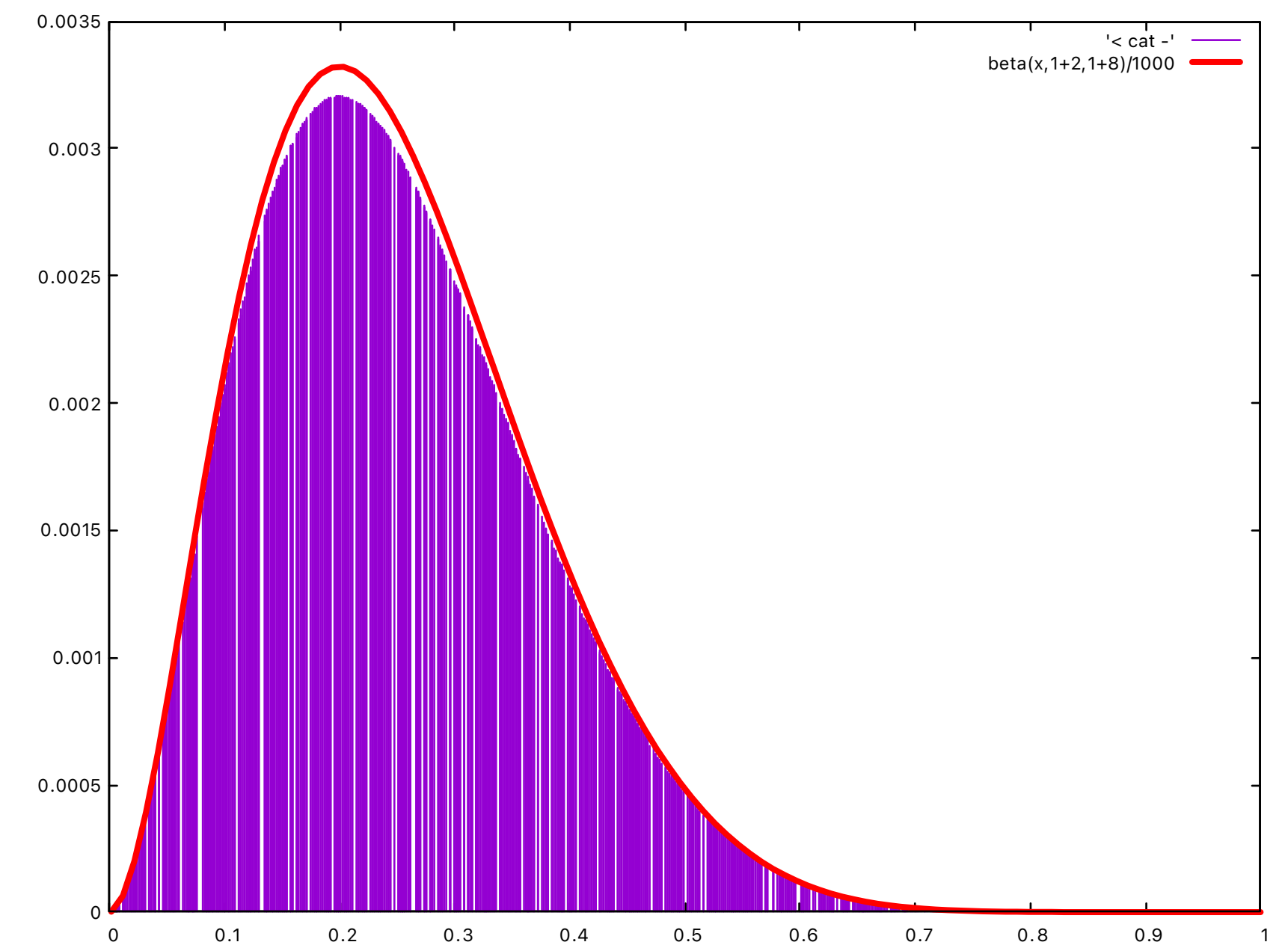


```
def coin(obs: list[int]) → float:  
    p = sample(Uniform(0, 1))  
    for o in obs:  
        observe(Bernoulli(p), o)  
    return p
```

```
with ImportanceSampling(num_particles=1000):  
    dist = infer(coin, [0, 0, 0, 0, 0, 0, 0, 0, 1, 1])  
    viz(dist)
```

Exact solution: $\text{Beta}(\text{\#heads} + 1, \text{\#tails} + 1)$

1000 particles



Examples: $\mu\text{-ppl}$

Reminders: Measure Theory

Probabilistic Programming Languages

Measurable Space

A σ -algebra on a set X is a collection of subsets

- containing \emptyset
- closed under complement
- closed under countable union

A measurable space is a pair (X, Σ_X)

- X : set
- Σ_X : measurable sets

Examples

- $\emptyset \dots$
- Singleton: $\{\emptyset, \{\text{unit}\}\}$
- Booleans: $\{\emptyset, \{\text{true}\}, \{\text{false}\}, \{\text{true}, \text{false}\}\}$
- Borel sets (intervals on \mathbb{R})
- Product: $\Sigma_{A \otimes B}$ generated by the rectangles

Measure

A measure maps a measurable set to a positive score: $\mu : \Sigma_X \rightarrow [0, \infty)$

- $\Sigma_X, \mu(U) \geq 0$ for all set U in Σ_X
- $\mu(\emptyset) = 0$
- $\mu(\bigcup_{i \in I} U_i) = \sum_{i \in I} \mu(U_i)$ for all countable collections $\{U_i\}_{i \in I}$ of pairwise disjoint sets in Σ_X

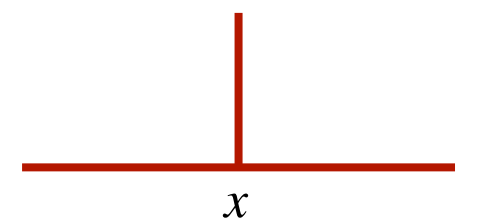
Probability distributions are normalized measures, i.e., $\mu(X) = 1$

Examples

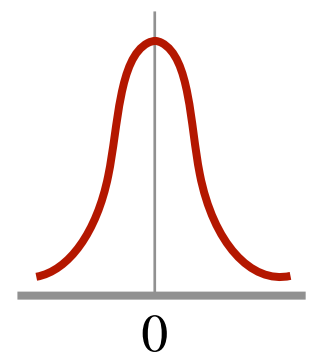
- Lebesgue measure
 $\lambda(a, b) = b - a$
- Bernoulli distribution (discret



- Dirac delta measure
$$\delta_x(U) = \begin{cases} 1 & \text{if } x \in U \\ 0 & \text{otherwise} \end{cases}$$



- Normal distribution
 $\mathcal{N}(0, 1)([0, 1]) \approx 0.34$
 $\mathcal{N}(0, 1)((-\infty, 0]) = 0.5$



Measurable Function

$f : X \rightarrow Y$ is measurable if $f^{-1}(U) \in \Sigma_X$ for all $U \in \Sigma_Y$

Pushforward: transfer a measure from a measurable space to another one

- $f : X \rightarrow Y$ measurable
- $\mu : \Sigma_X \rightarrow [0, \infty)$
- $f_*(\mu) : \Sigma_Y \rightarrow [0, \infty)$

Examples

- $\mu : \Sigma_{A \times B} \rightarrow [0, \infty)$
- $\pi_{1*}(\mu) : \Sigma_A \rightarrow [0, \infty)$
- $\pi_{2*}(\mu) : \Sigma_B \rightarrow [0, \infty)$

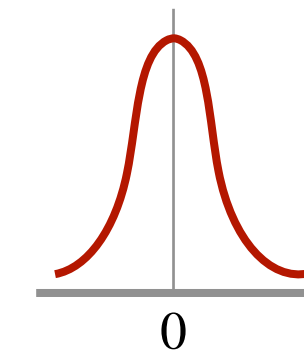
Integration

Given $f : X \rightarrow [0, \infty)$ a measurable function and a measure $\mu : X \rightarrow [0, \infty)$

- $\int_X f(x) \mu(dx) \in [0, \infty)$
- $\int_X \delta_x(U) \mu(dx) = \int_U \mu(dx) = \mu(U)$
- $\int_X f(x) r \mu(dx) = r \int_X f(x) \mu(dx)$
- $\int_X (f(x) + g(x)) \mu(dx) = \int_X f(x) \mu(dx) + \int_X g(x) \mu(dx)$

We can also define new measures from integration

- $f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$
- $\int_X f(x) \lambda(dx) = \int_X f(x) dx = \mathcal{N}(0, 1)(X)$ using the Lebesgue measure λ
- $f(x)$ is the density w.r.t. the Lebesgue measure λ



Expected value

Expected value

- If $X \sim F$ with density f w.r.t. the Lebesgue measure
- If g is a measurable function

$$\mathbb{E}(g(X)) = \int xF(dx) = \int xf(x)dx$$

Kernel

A kernel $k : X \times \Sigma_Y \rightarrow [0, \infty)$ is a function such that:

- $k(x, _) : \Sigma_Y \rightarrow [0, \infty)$ for all $x \in X$ is a measure
- $k(_, U) : X \rightarrow [0, \infty)$ for all $U \in \Sigma_Y$ is measurable

A probability kernel is such that $k(x, Y) = 1$ for all $x \in X$

Law of large numbers

Sample average

- If X_1, X_2, \dots, X_n are iid.
- Sample average: $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$

Strong law of large numbers

- If g is a measurable function and X, X_1, \dots, X_n are iid
- If the mean of $g(X)$ is finite
- Then the sample average approximates the mean

$$\overline{g(X)}_n \xrightarrow{a.s.} \mathbb{E}(g(X)) \text{ when } n \rightarrow \infty$$

Applications

- $g = id$ to compute the mean
- $g = \mathbb{1}_{[x>0]}$ to compute $P(X > 0)$
- $g = \pi_1$ to compute the first marginal of (X_1, X_2)

Correctness of importance sampling

Model

- Samples X_i from the priors with density $\pi(x)$
- Observed data: $Y = y$
- Likelihood (product of all observe statements): $\ell(x; y)$
- Posterior: $p(x|y) = \frac{\pi(x)\ell(x; y)}{Z}$ where $Z = \int \pi(x)\ell(x; y)dx$ (evidence)

Importance sampling estimation with n particles:
$$\text{IS}(g(X)|y)_n = \frac{\sum_{i=1}^n W_i g(X_i)}{\sum_{i=1}^n W_i} = \frac{\sum_{i=1}^n \ell(X_i; y)g(X_i)}{\sum_{i=1}^n \ell(X_i; y)} = \frac{\overline{(\ell(X; y)g(X))}_n}{\overline{\ell(X; y)}_n}$$

Correctness with SLLN

- $\overline{(\ell(X; y)g(X))}_n \xrightarrow{a.s.} \mathbb{E}(\ell(X; y)g(X)) = \int g(x)\ell(x; y)\pi(x)dx$
- $\overline{\ell(X; y)}_n \xrightarrow{a.s.} \mathbb{E}(\ell(X; y)) = \int \ell(x; y)\pi(x)dx = Z$

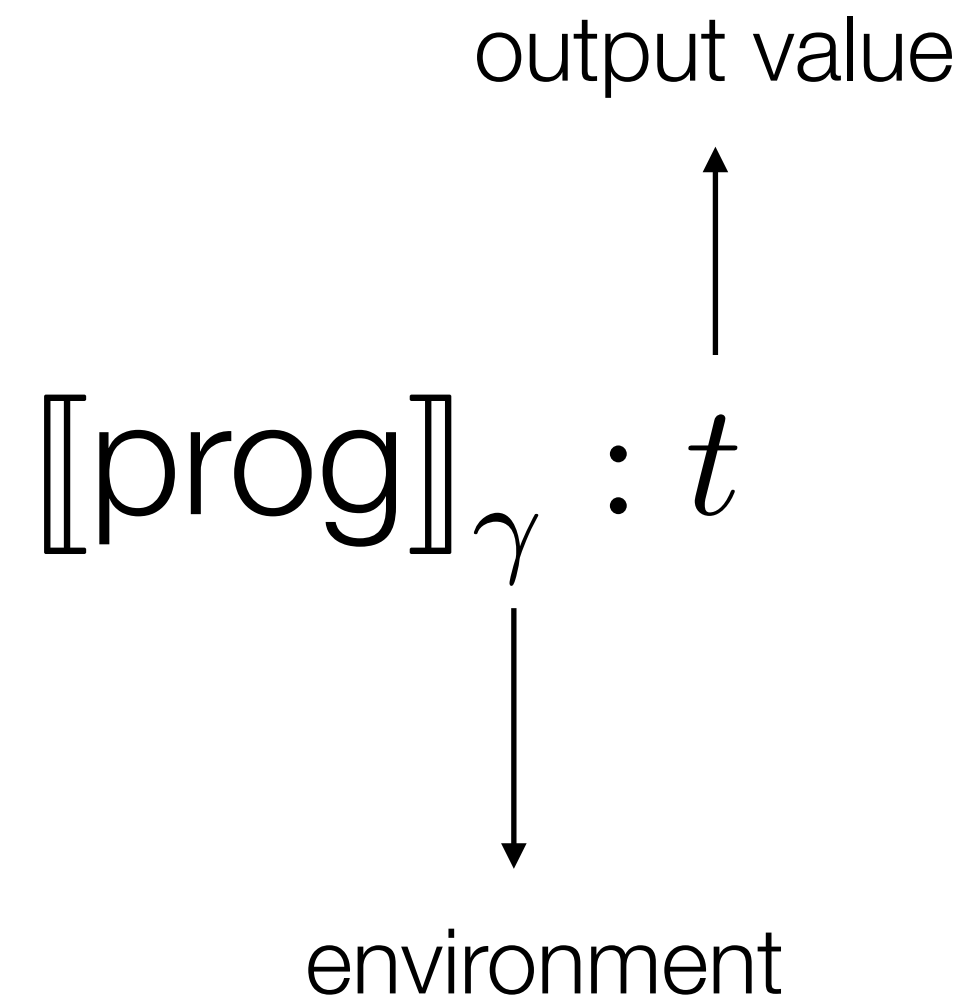
$$\text{IS}(g(X)|y)_n \xrightarrow{a.s.} \int g(x) \frac{\ell(x; y)\pi(x)}{Z} dx = \int g(x)p(x|y)dx = \mathbb{E}(g(X)|y)$$

Careful with 0, and ∞ ...

Semantics

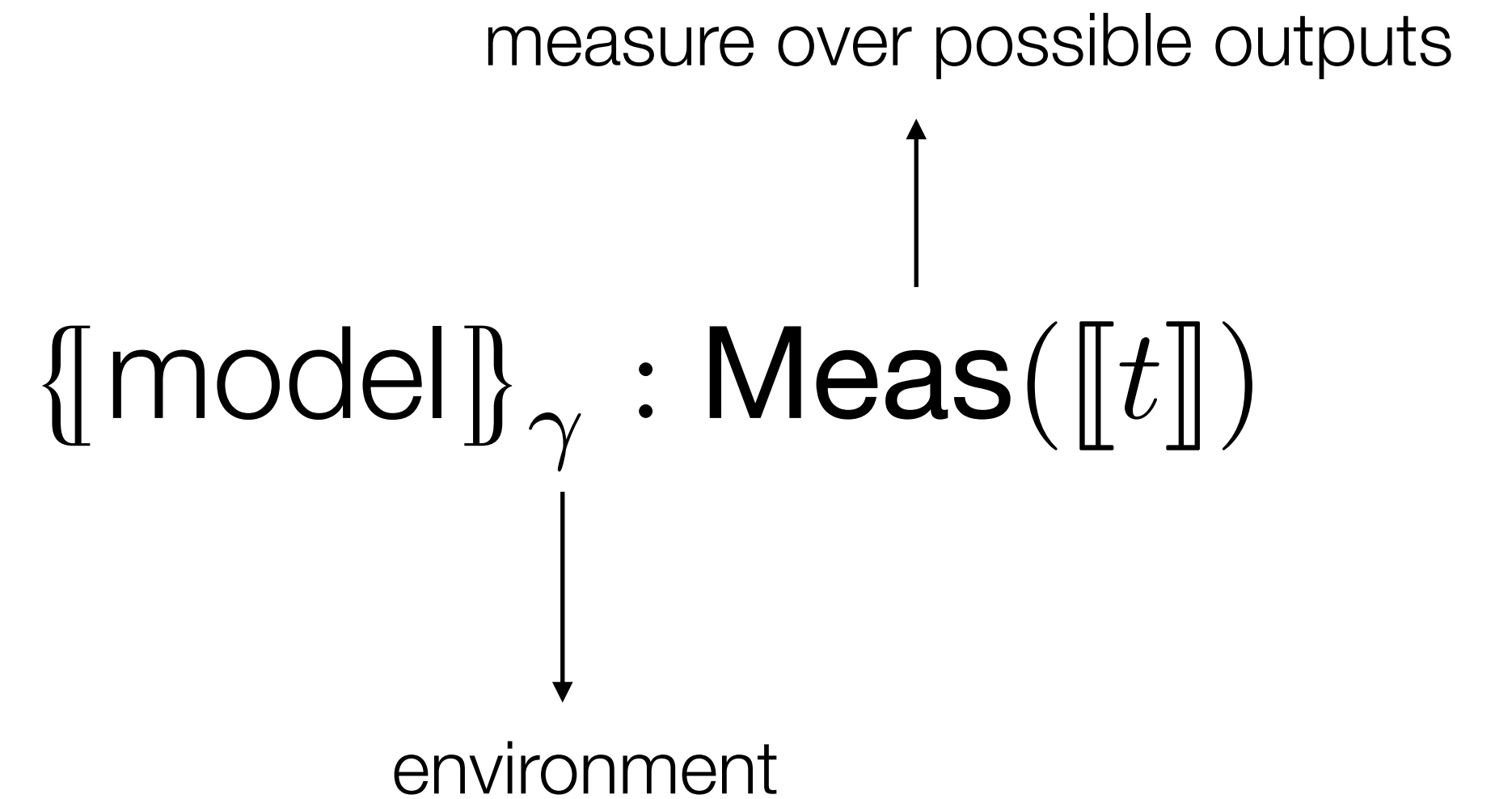
Probabilistic Programming Languages

Deterministic vs. probabilistic semantics



$$\llbracket 2 + 5 \rrbracket_\square = 7$$

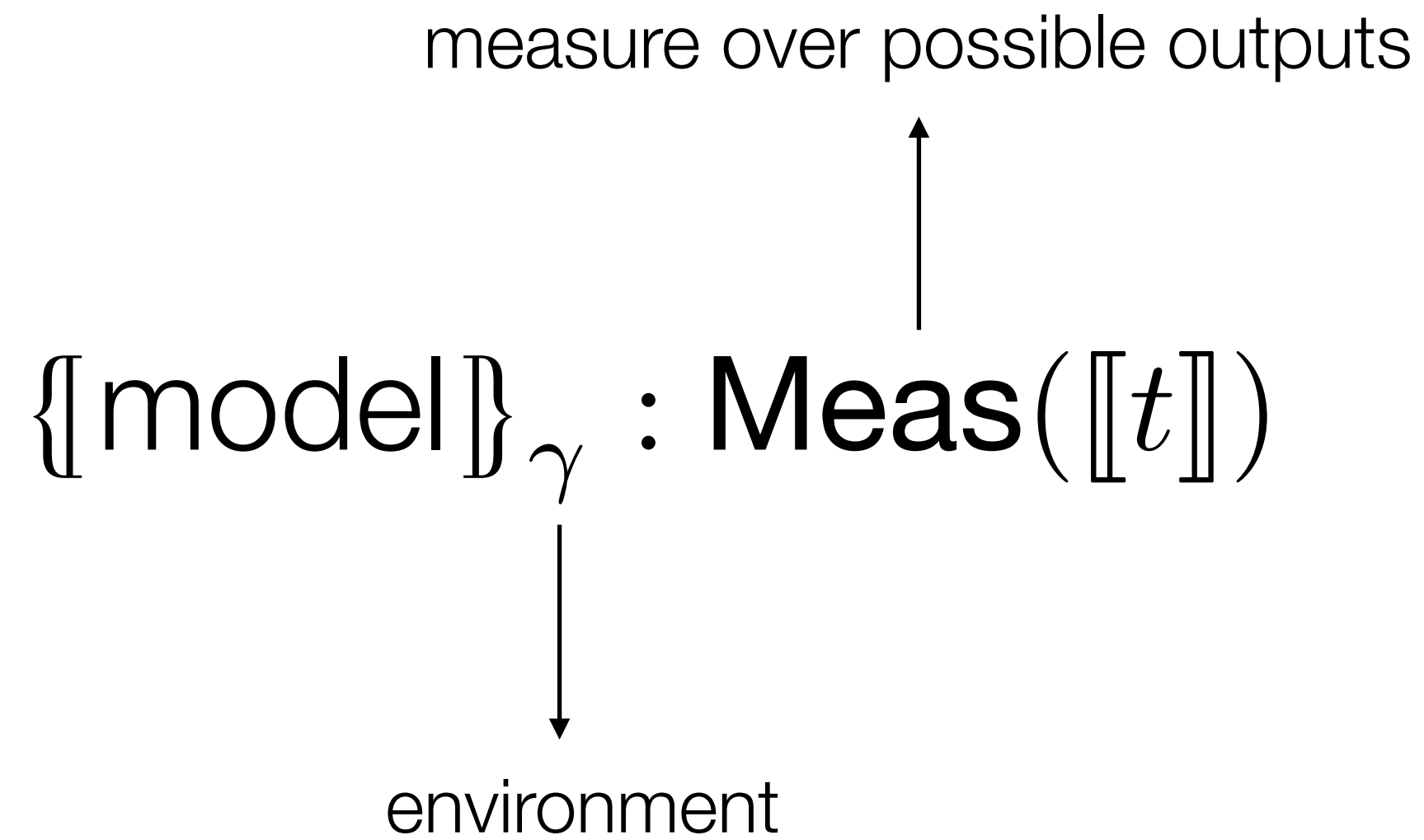
$$\llbracket x + y \rrbracket_{[x \leftarrow 2, y \leftarrow 5]} = 7$$



$$\{2 + 5\}_\square(U) = \delta(7)(U)$$

$$\{\text{sample}(\text{normal}(x, y))\}_{[x \leftarrow 1, y \leftarrow 0]}(U) = \mathcal{N}(0, 1)(U)$$

(Un)normalized measures



Unnormalized measure

$$\llbracket \text{infer}(\text{model}) \rrbracket_\gamma = \frac{\llbracket \text{model} \rrbracket_\gamma}{\llbracket \text{model} \rrbracket_\gamma (\llbracket \text{typeOf}(\text{model}) \rrbracket)}$$

Distribution



normalize over all possible values

Takeaway

I - Probabilistic programming

- A program describes a distribution
- Sample: draw from distribution
- Bayesian reasoning: condition on data
- Inference computes the distribution: Intractable problem in general

II - Approximate inference

- Importance sampling: weighted sampler returns pairs (value, score)
- Correctness with the law of large numbers

III - Semantics

- Deterministic semantics: expression to value
- Probabilistic semantics: expression to measure of possible values



<https://github.com/mpri-probprog/probprog-25-26>