

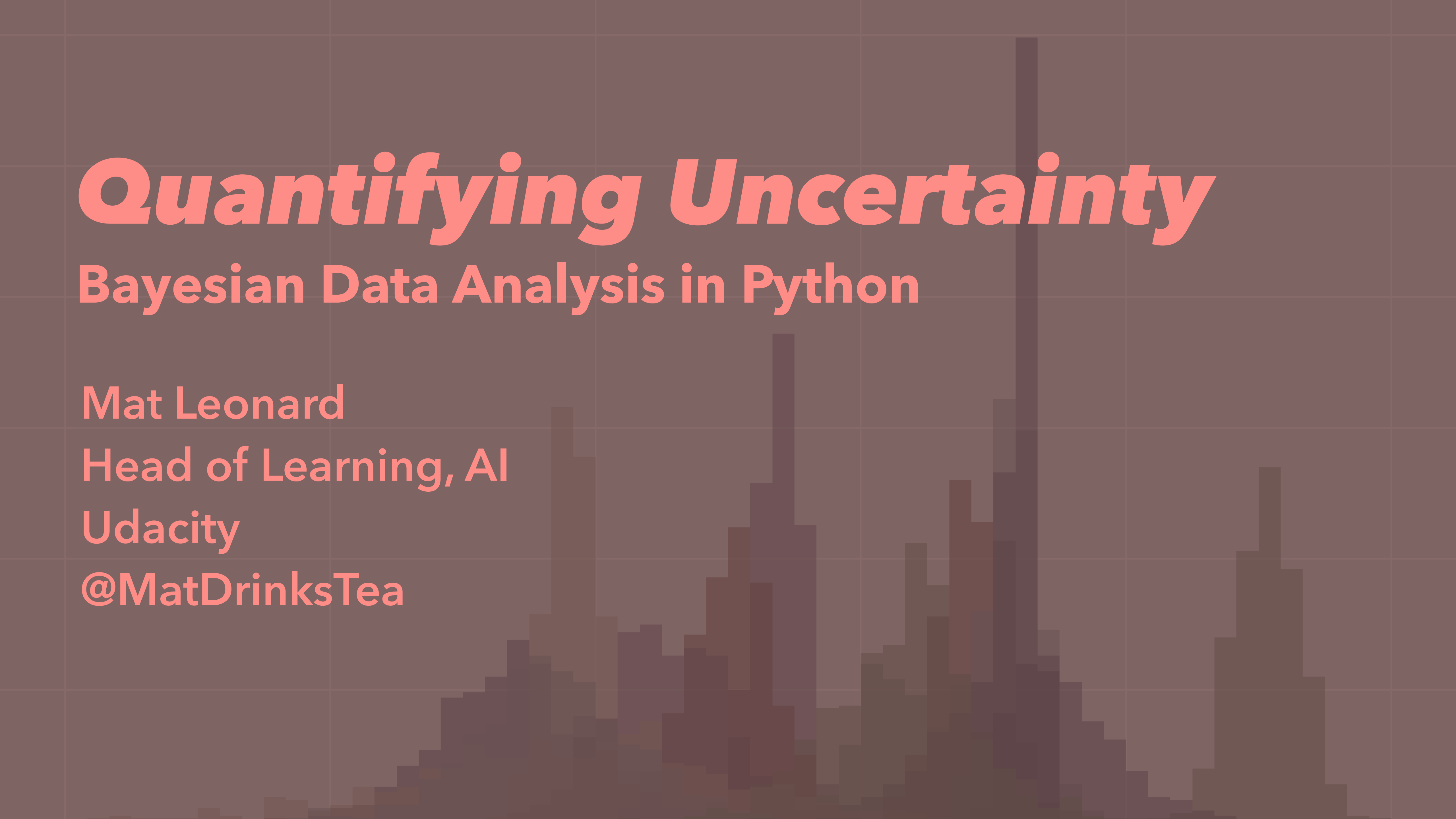
Quantifying Uncertainty

Bayesian Data Analysis in Python

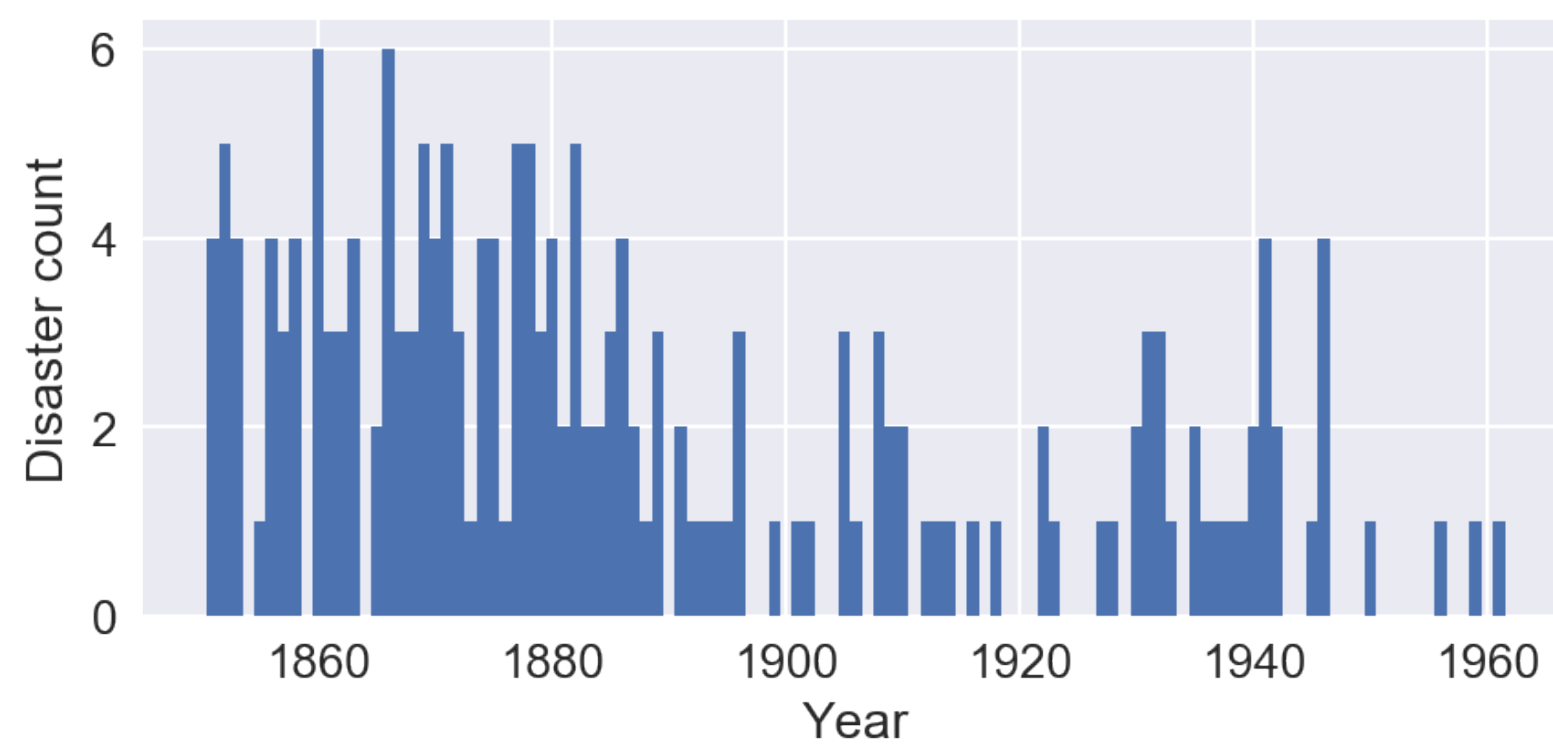
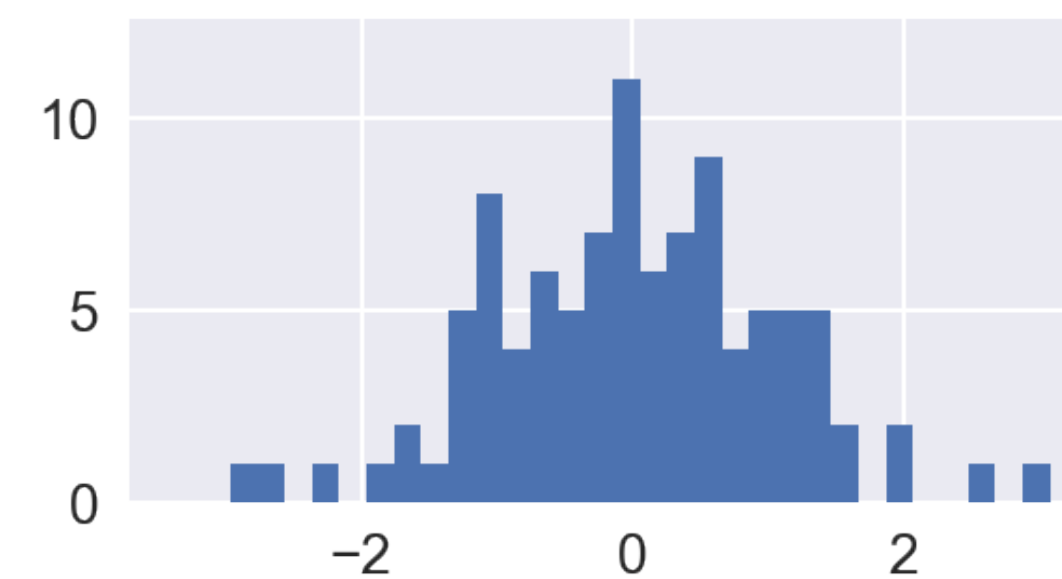
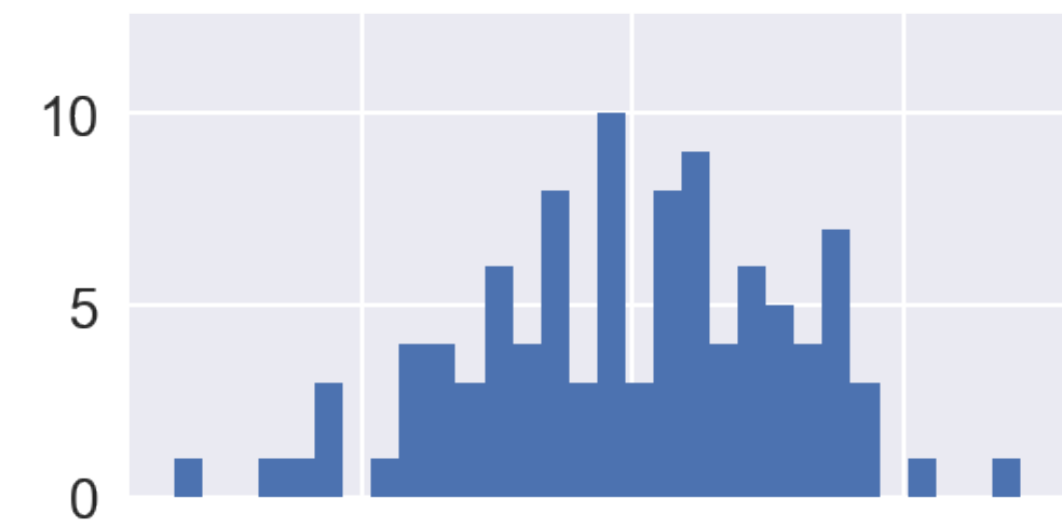
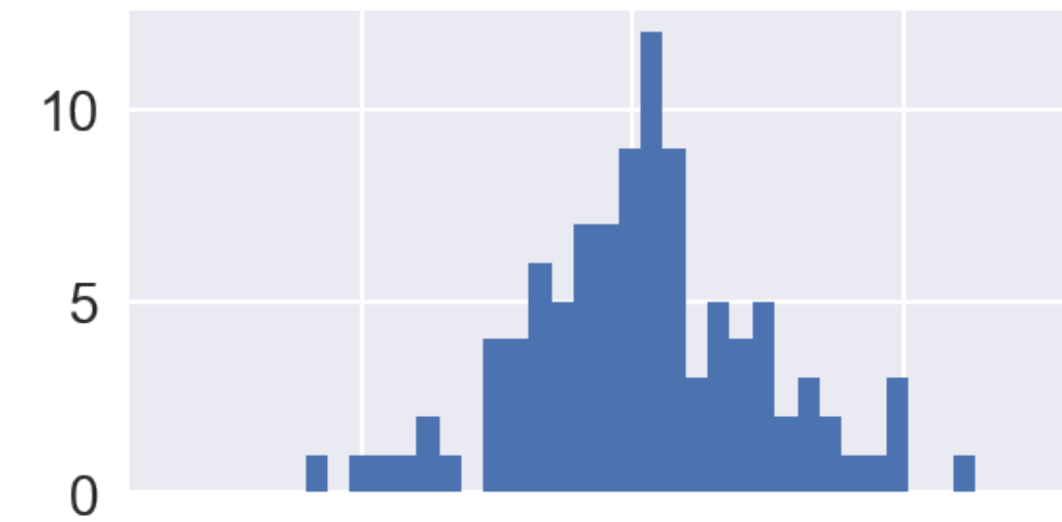
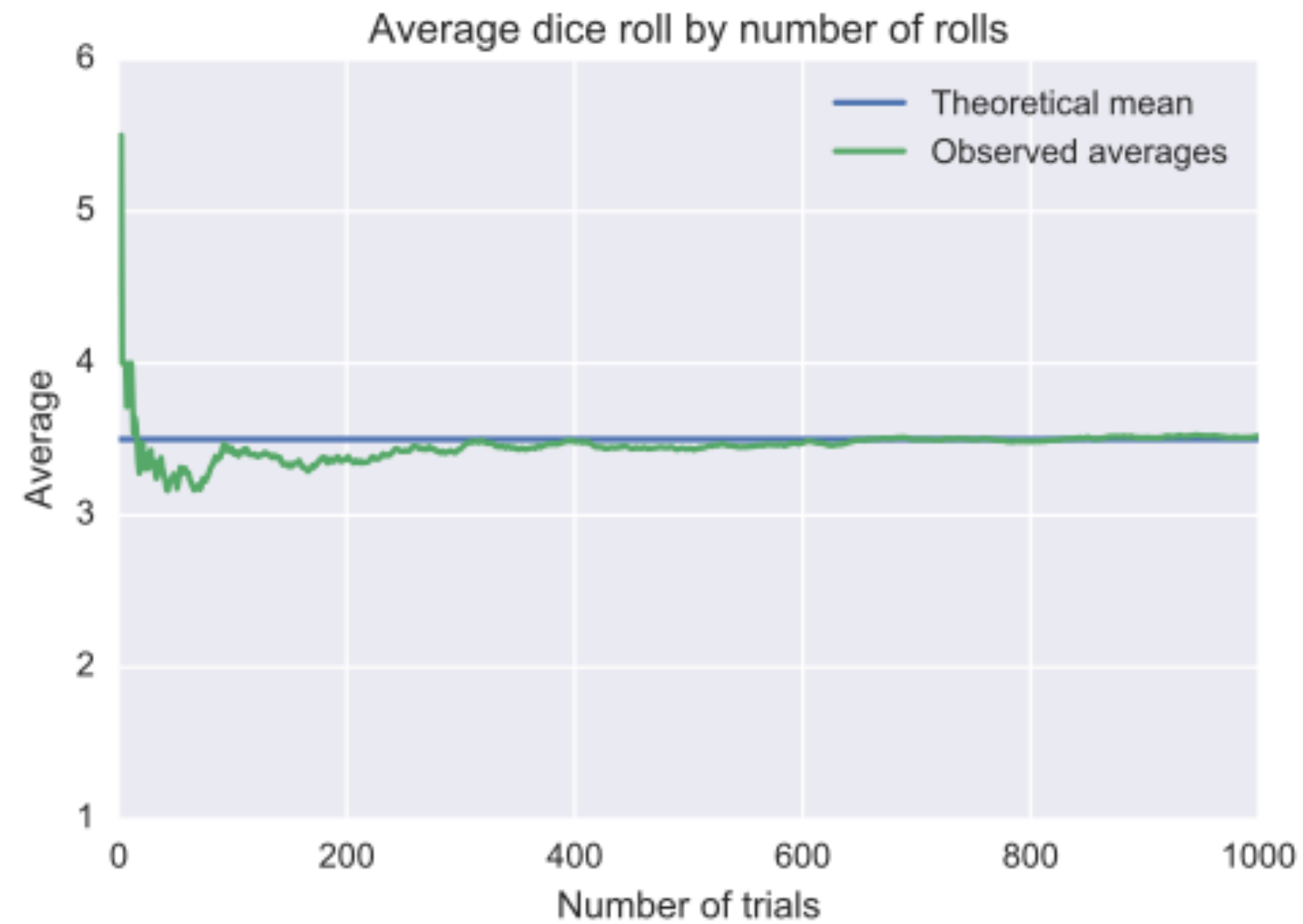
Mat Leonard

**Head of Learning, AI
Udacity**

@MatDrinksTea



Data is Life, Life is Uncertain



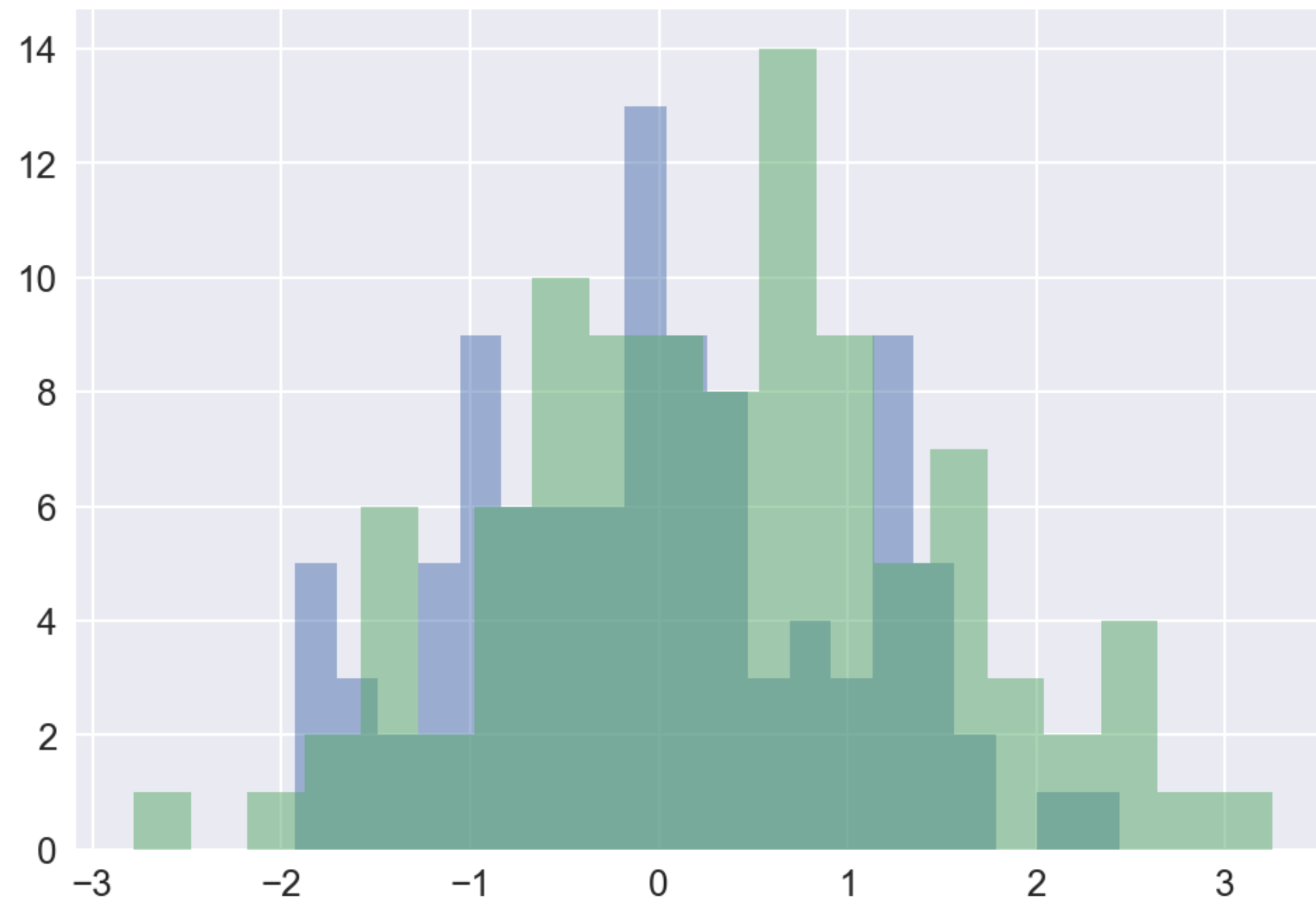
Quantify Uncertainty

with

STATISTICS

Hypothesis Testing Refresher

$d = 0.3, n=100$

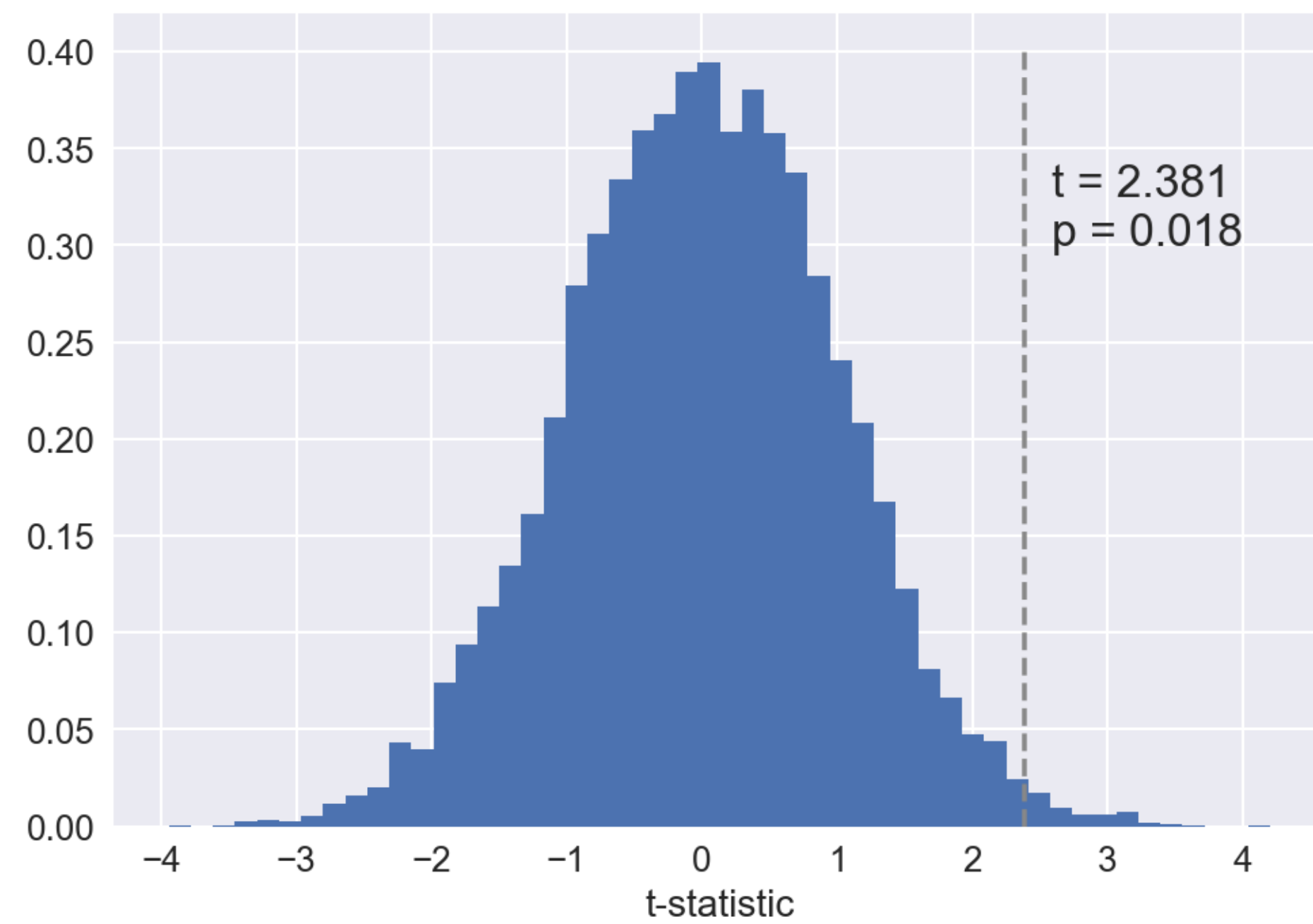


$$t = \sqrt{n} \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{s_{X_1}^2 + s_{X_2}^2}}$$

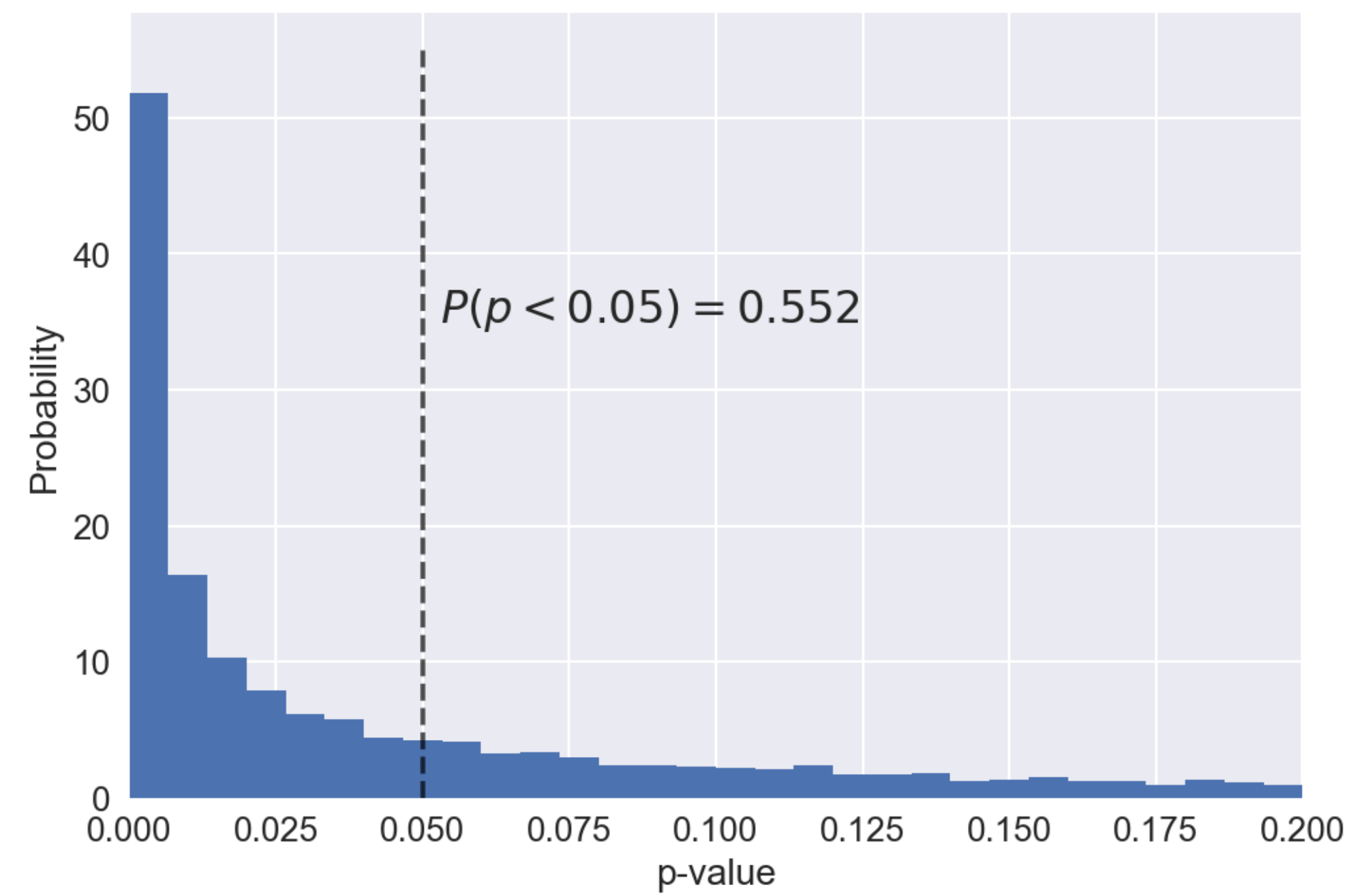
$$H_0 : t = 0$$

$$H_1 : t \neq 0$$

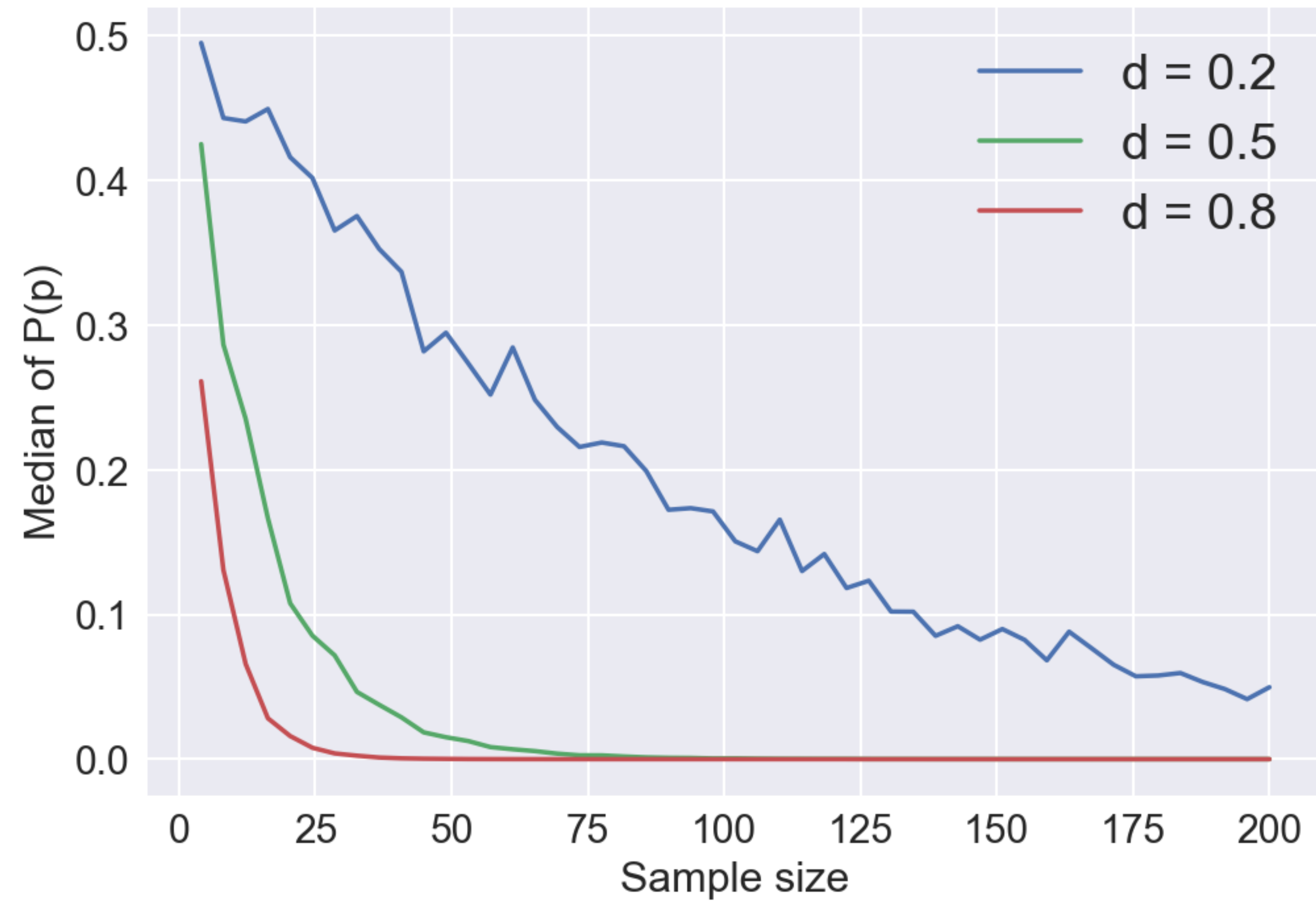
$$H_0 : t = 0$$



$$d = 0.3, n=100$$



Median of the p sampling distribution: function of the effect size (d) and the sample size



**p-values are *not* a
measure of uncertainty**

**Parameter estimation with
confidence intervals?**

What people think (and want)

95% probability the parameter is in the 95% confidence interval

What it actually is

If you perform your experiment an infinite number of times, calculating the 95% confidence interval for each experiment, then for 95% of those experiments the true parameter will fall within the 95% confidence interval

Bayesian Inference

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Bayes' Theorem

$$P(\theta | y) \propto \prod_i^N P(y_i | \theta) P(\theta)$$

Bayes' Theorem

$$\overbrace{P(\theta | y)}^{\text{Posterior}} \propto \prod_i^N \underbrace{P(y_i | \theta)}_{\text{Data Likelihood}} \overbrace{P(\theta)}^{\text{Prior}}$$

Flipping Coins, a Classic

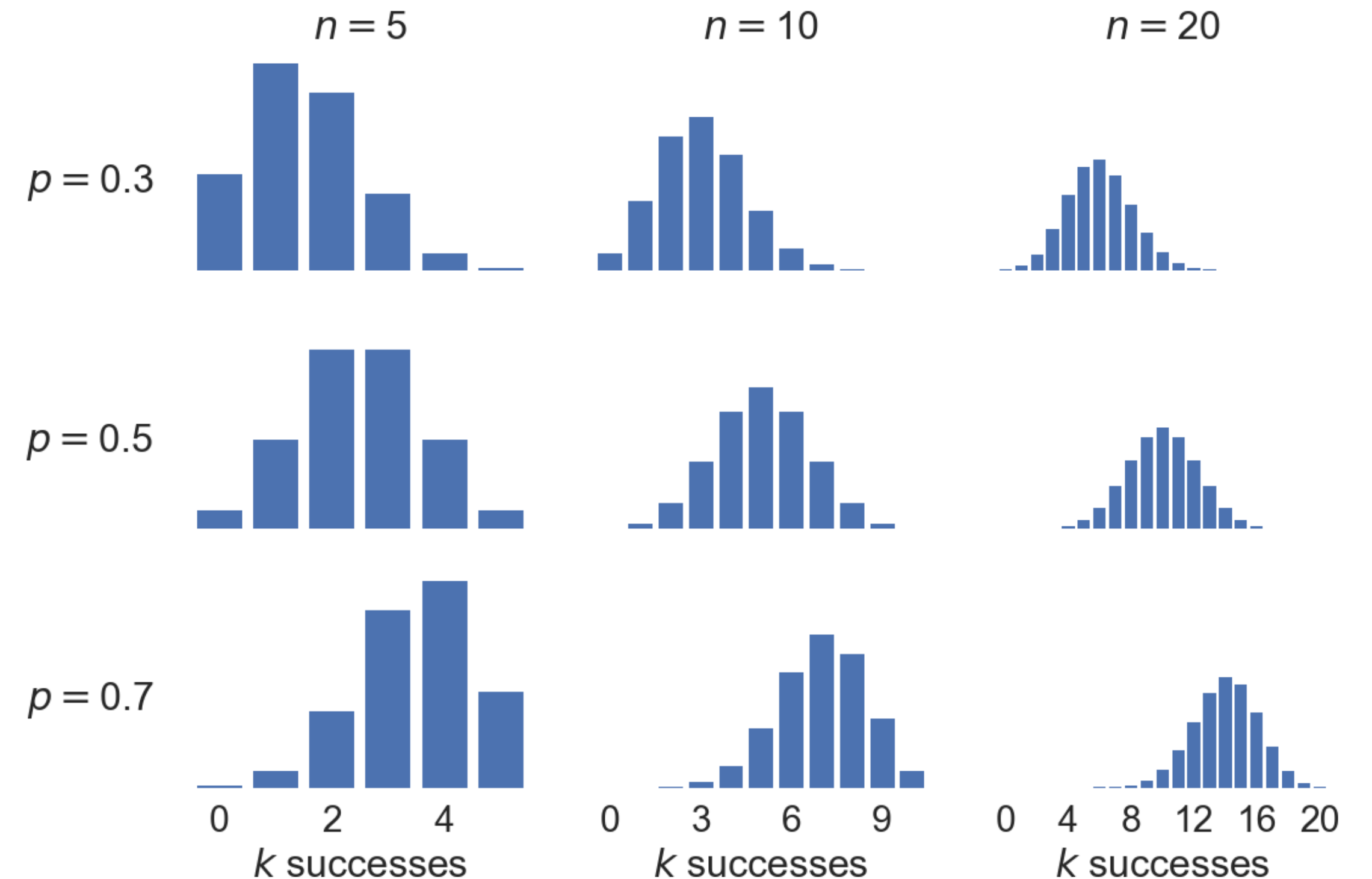
Binomial distribution

Attempts

Probability of success

$$\Pr(n, k, p) = \binom{n}{k} p^k (1-p)^{n-k}$$

Successes



Is our coin fair???

Bayesian Model

Likelihood

$$\text{Binomial}(k, n \mid p)$$

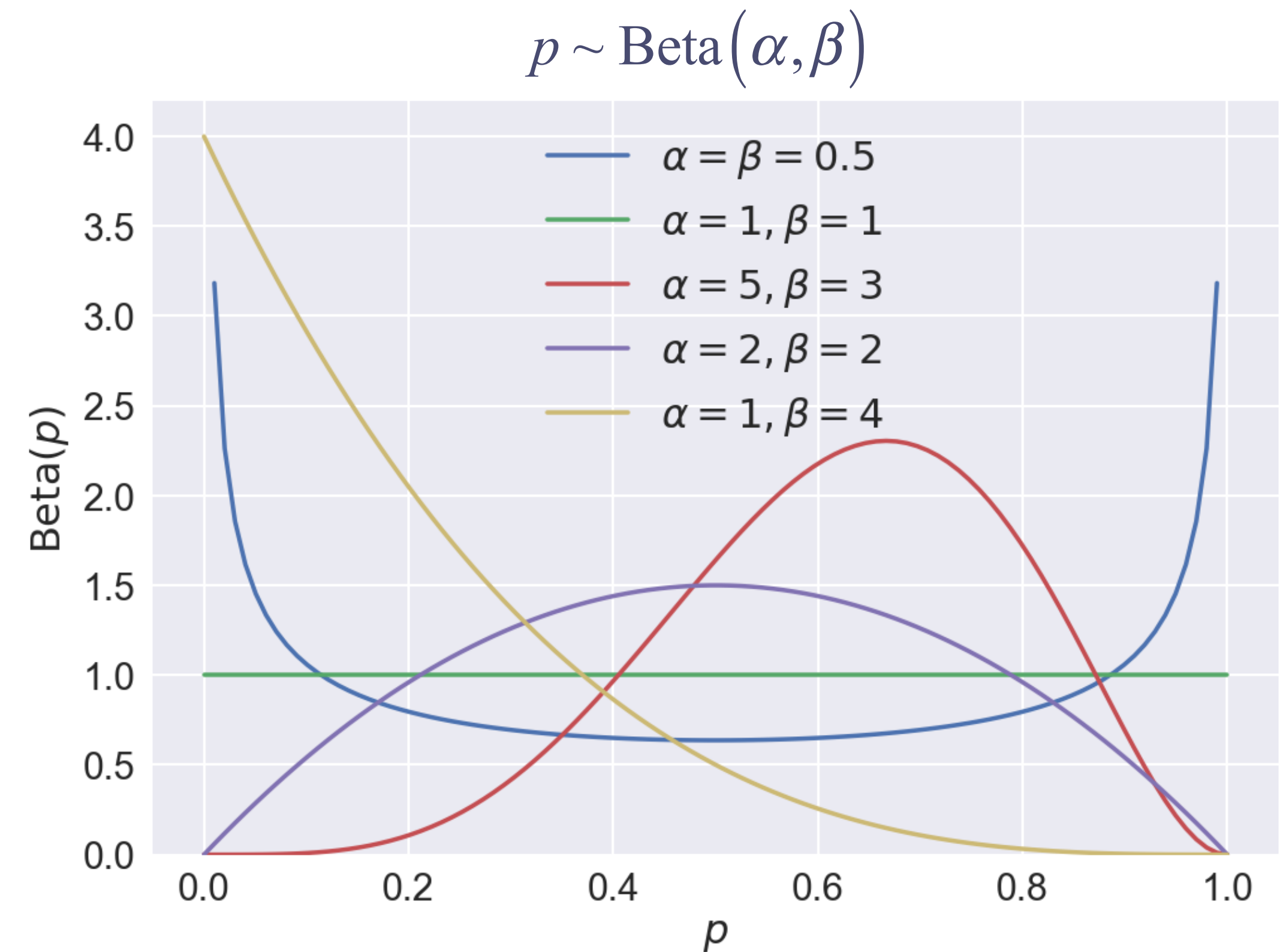
Prior

$$p \sim \text{Beta}(\alpha, \beta)$$

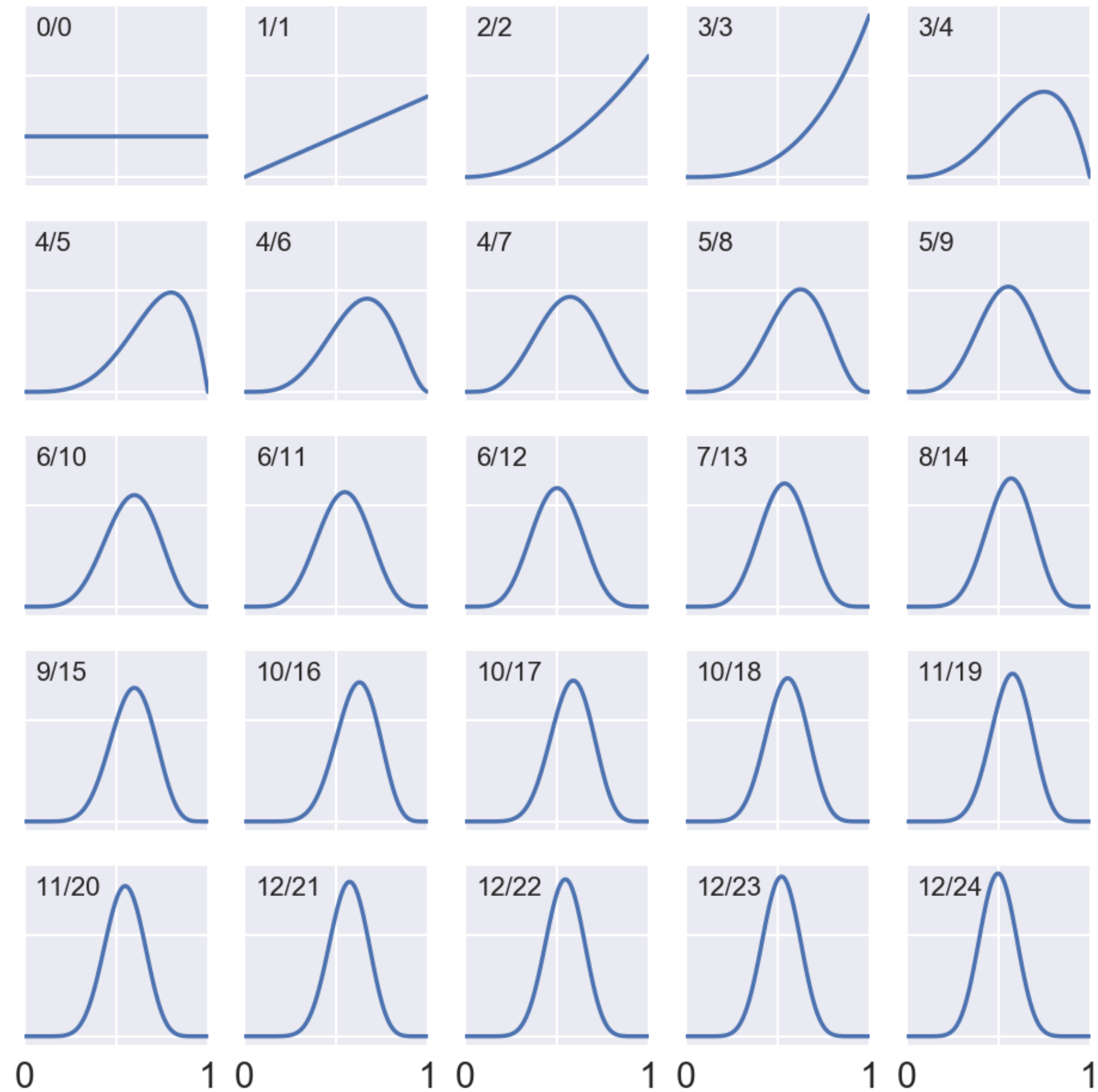
Posterior

$$P(p \mid k, n) \propto P(k, n \mid p) P(p)$$

$$P(p \mid k, n) = \text{Beta}(\alpha + k, \beta + n - k)$$



Posterior distributions, $p = 0.5$



Now with Python

```
from scipy import stats

# Our data, events, array of 0s and 1s
successes = events.sum()
failures = len(events) - successes

# Prior parameters
 $\alpha$ ,  $\beta$  = 1, 1

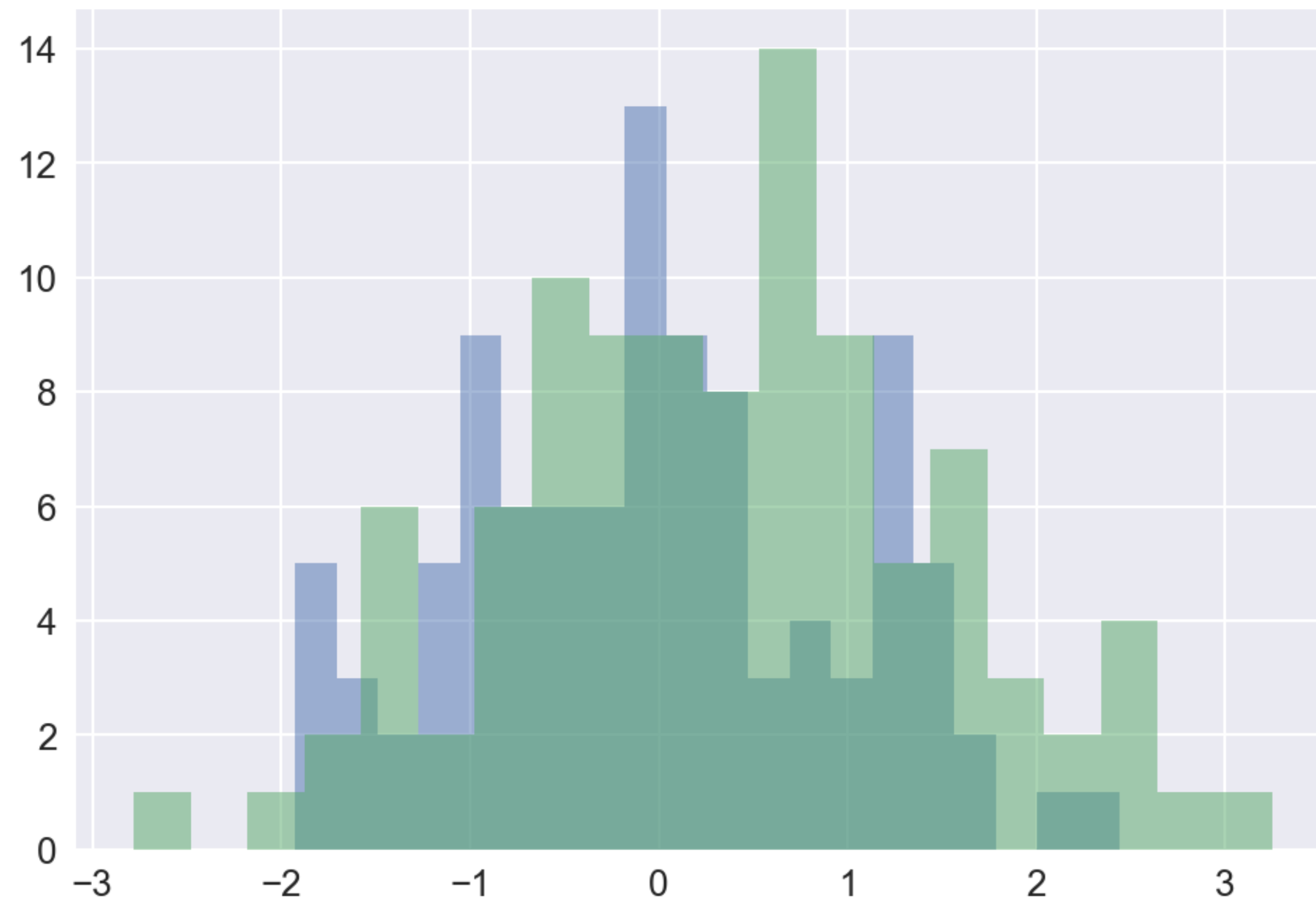
# Calculate posterior distribution
posterior = stats.beta( $\alpha$  + successes,  $\beta$  + failures)

# Mean and 95% credible interval
mean = posterior.mean()
CR = posterior.interval(0.95)

# Posterior distribution for plotting
xs = np.linspace(0, 1, num=100)
pdf = posterior.pdf(xs)
```

Back to the T-test

$d = 0.3, n=100$



What do we want?

Means: μ_1, μ_2

Standard deviations: σ_1, σ_2

Estimates: $\mu_2 - \mu_1$ and $\sigma_2 - \sigma_1$

Assess the *uncertainty* in our estimates

Forget the T-test

Likelihood

$$y_1 \sim \text{Normal}(\mu_1, \sigma_1)$$

$$y_2 \sim \text{Normal}(\mu_2, \sigma_2)$$

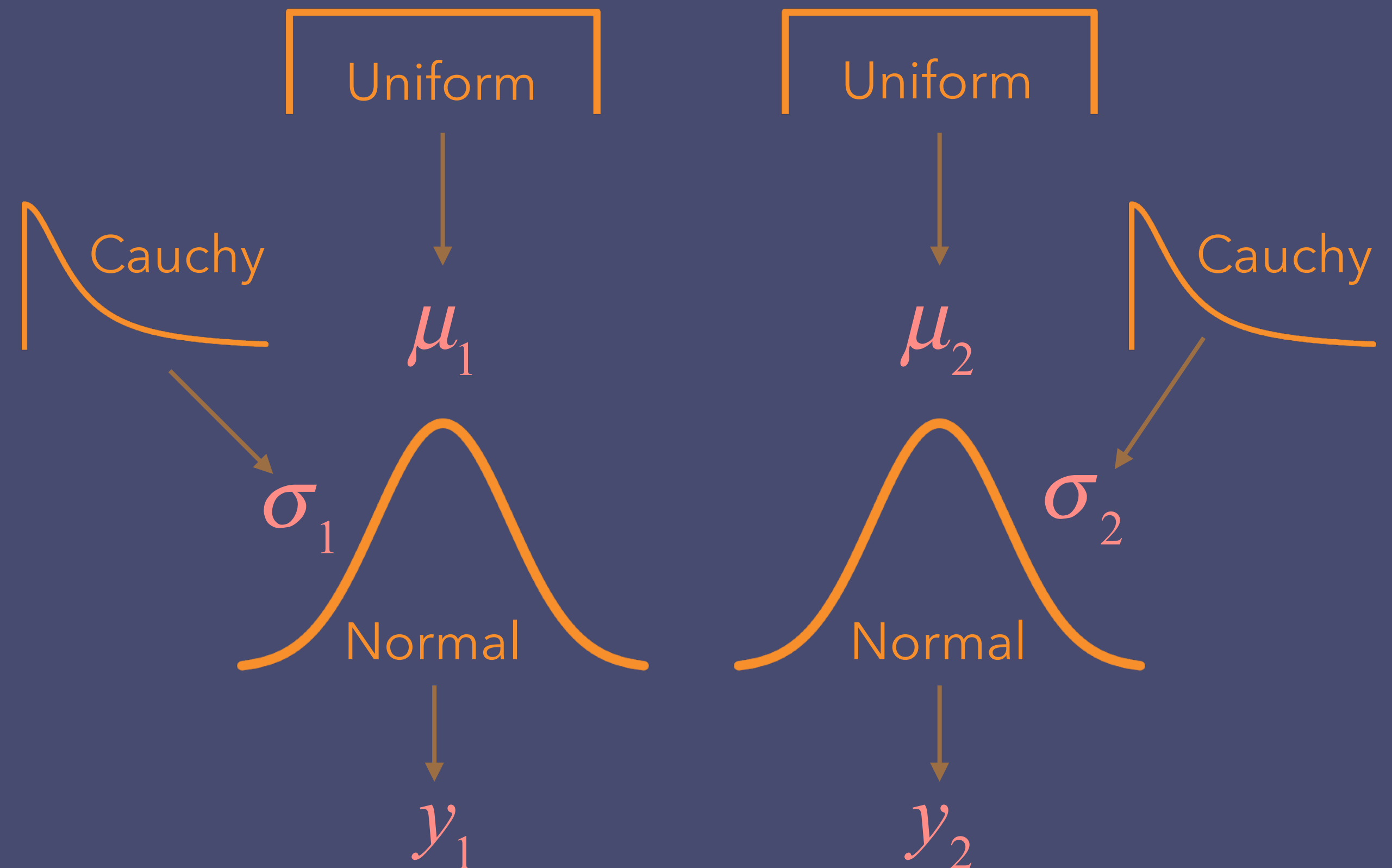
Priors

$$\mu_1 \sim \text{Uniform}(-5, 5)$$

$$\sigma_1 \sim \text{Half Cauchy}(\gamma = 5)$$

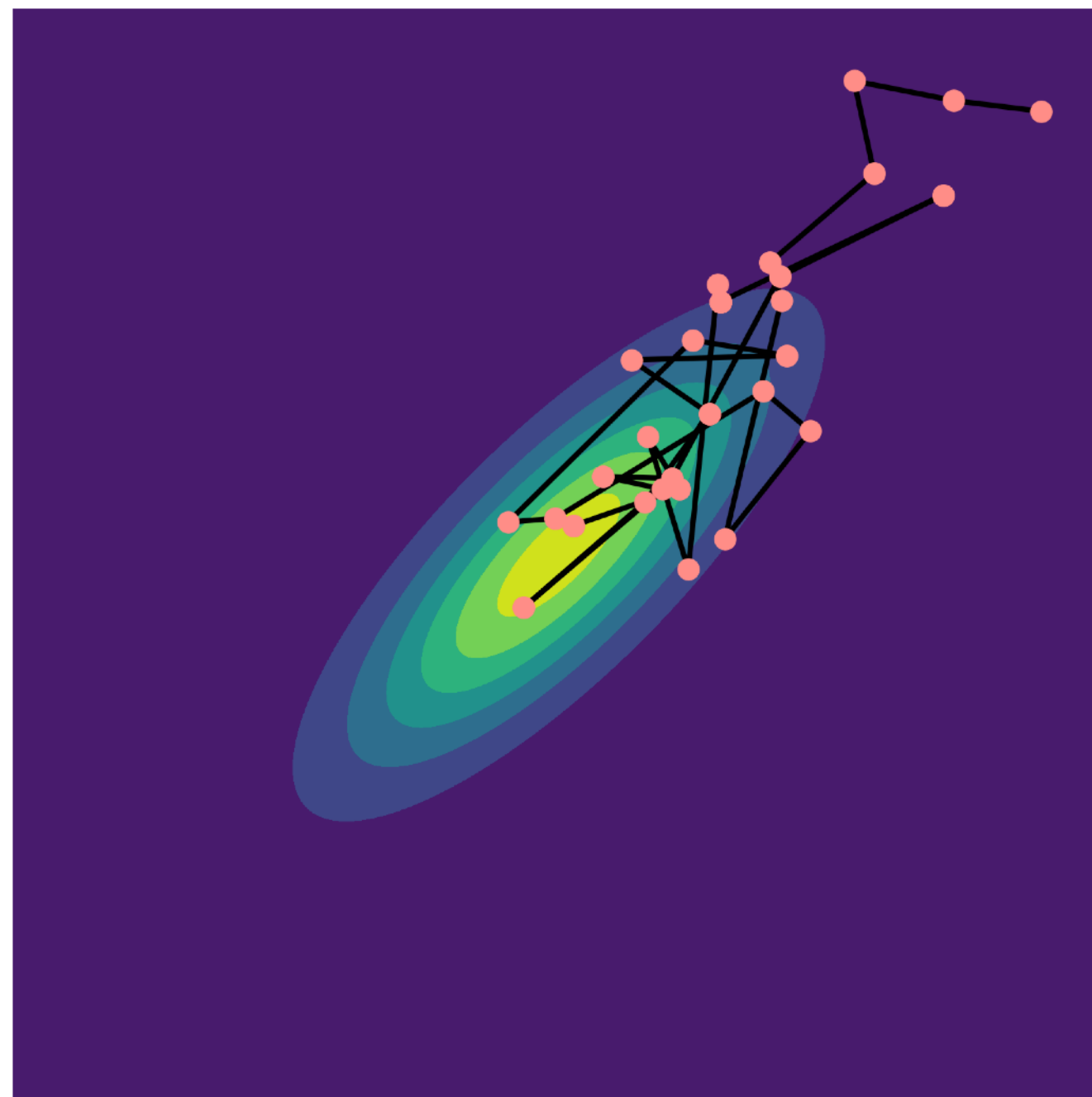
$$\mu_2 \sim \text{Uniform}(-5, 5)$$

$$\sigma_2 \sim \text{Half Cauchy}(\gamma = 5)$$

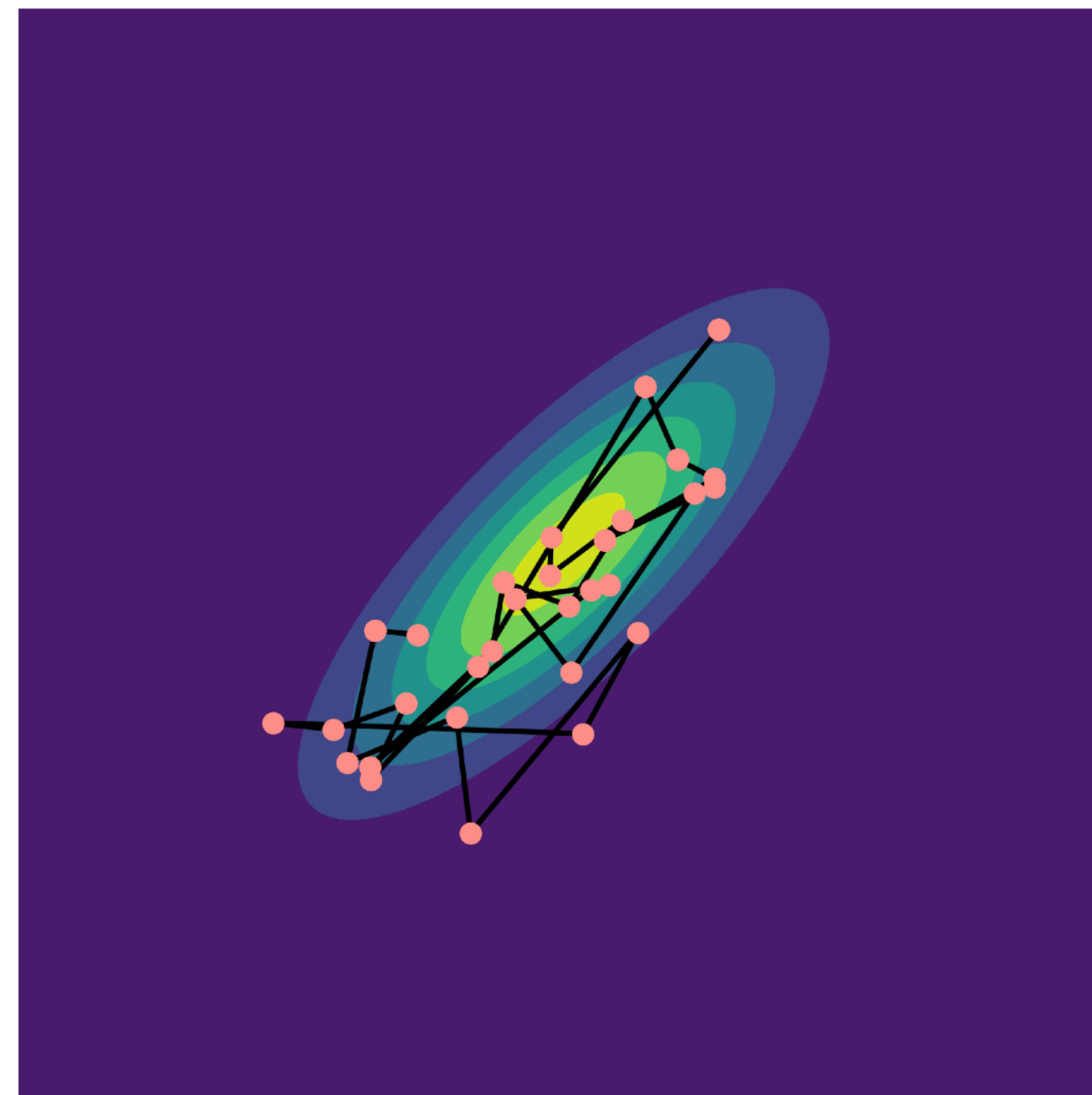


Markov Chain Monte Carlo

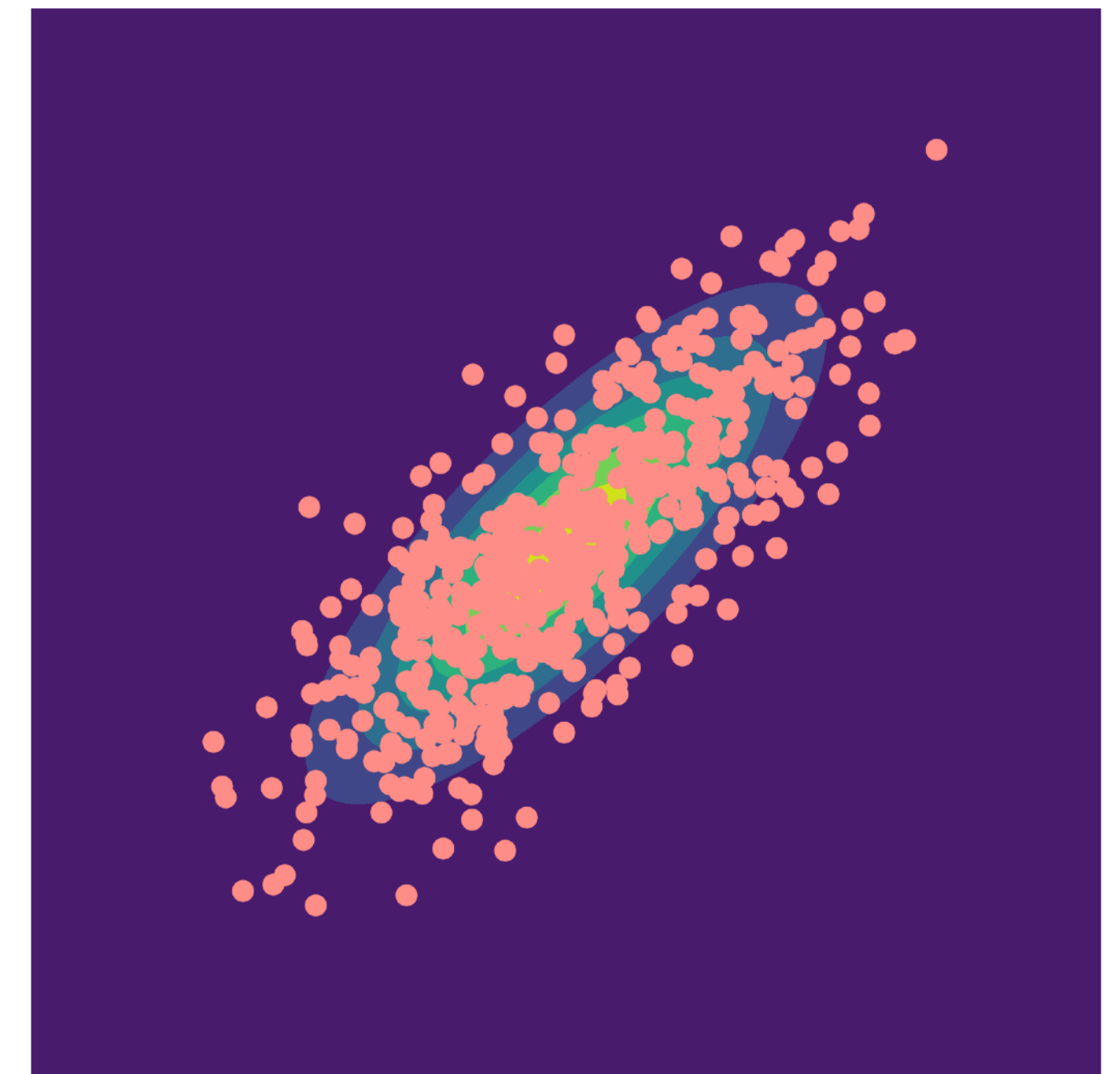
Starting out



Converged

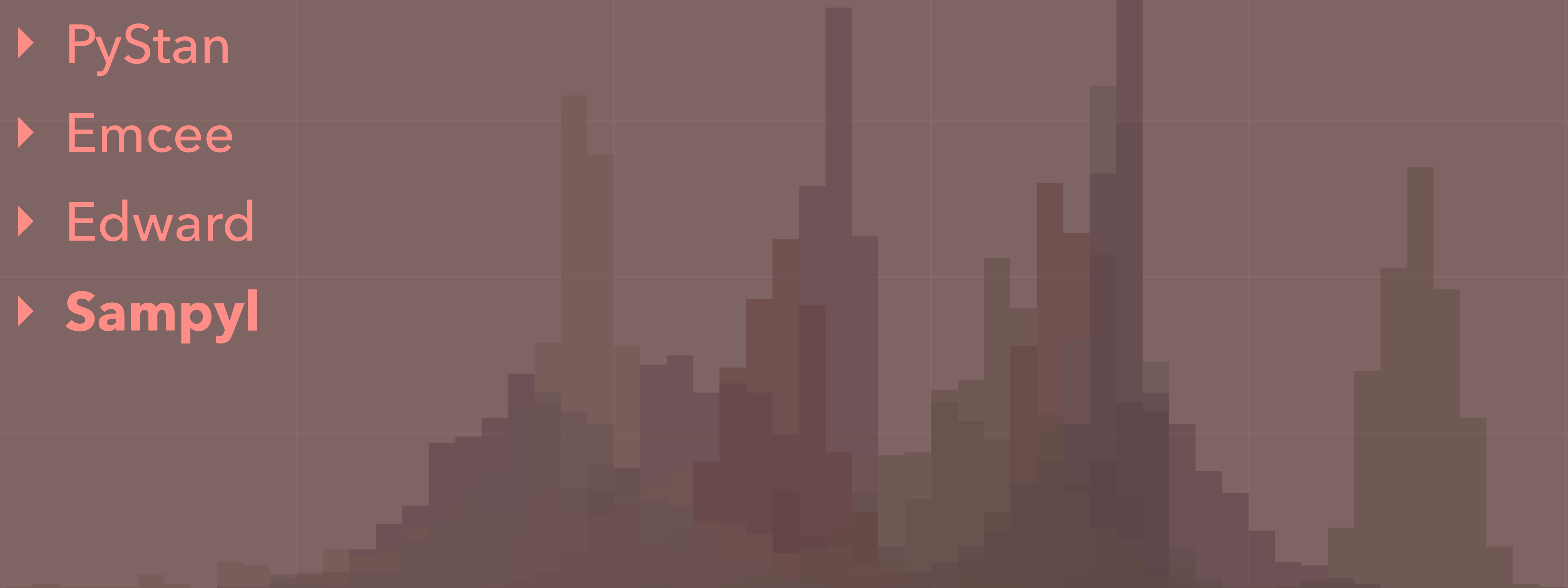


Sampled



Bayesian Data Analysis in Python

- ▶ PyMC
- ▶ PyStan
- ▶ Emcee
- ▶ Edward
- ▶ **Sampyl**



Independent Samples Model

Define the model in SamPyI

```
import sampyl as smp  
from sampyl import np
```

Independent Samples Model

Define the model in SamPy

```
def logp( $\mu_1$ ,  $\sigma_1$ ,  $\mu_2$ ,  $\sigma_2$ ):  
  
    model = smp.Model()  
  
    # Priors for means  
    model.add(smp.uniform( $\mu_1$ , -5, 5),  
              smp.uniform( $\mu_2$ , -5, 5))  
  
    # Priors for standard devs  
    model.add(smp.half_cauchy( $\sigma_1$ , beta=5),  
              smp.half_cauchy( $\sigma_2$ , beta=5))  
  
    # Data Likelihood  
    model.add(smp.normal(group1, mu= $\mu_1$ , sig= $\sigma_1$ ),  
              smp.normal(group2, mu= $\mu_2$ , sig= $\sigma_2$ ))  
  
    return model()
```

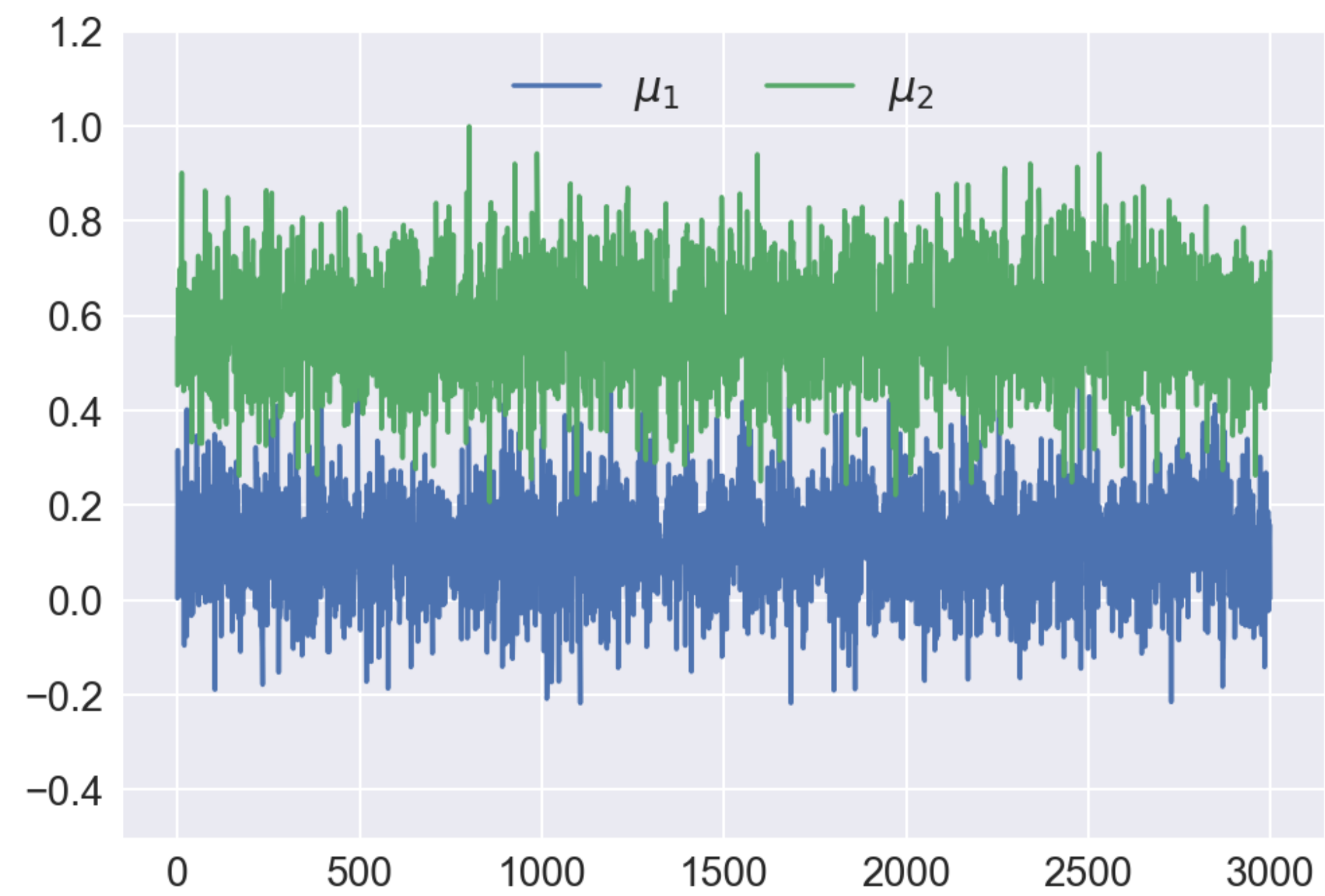

Independent Samples Model

Sample from the posterior

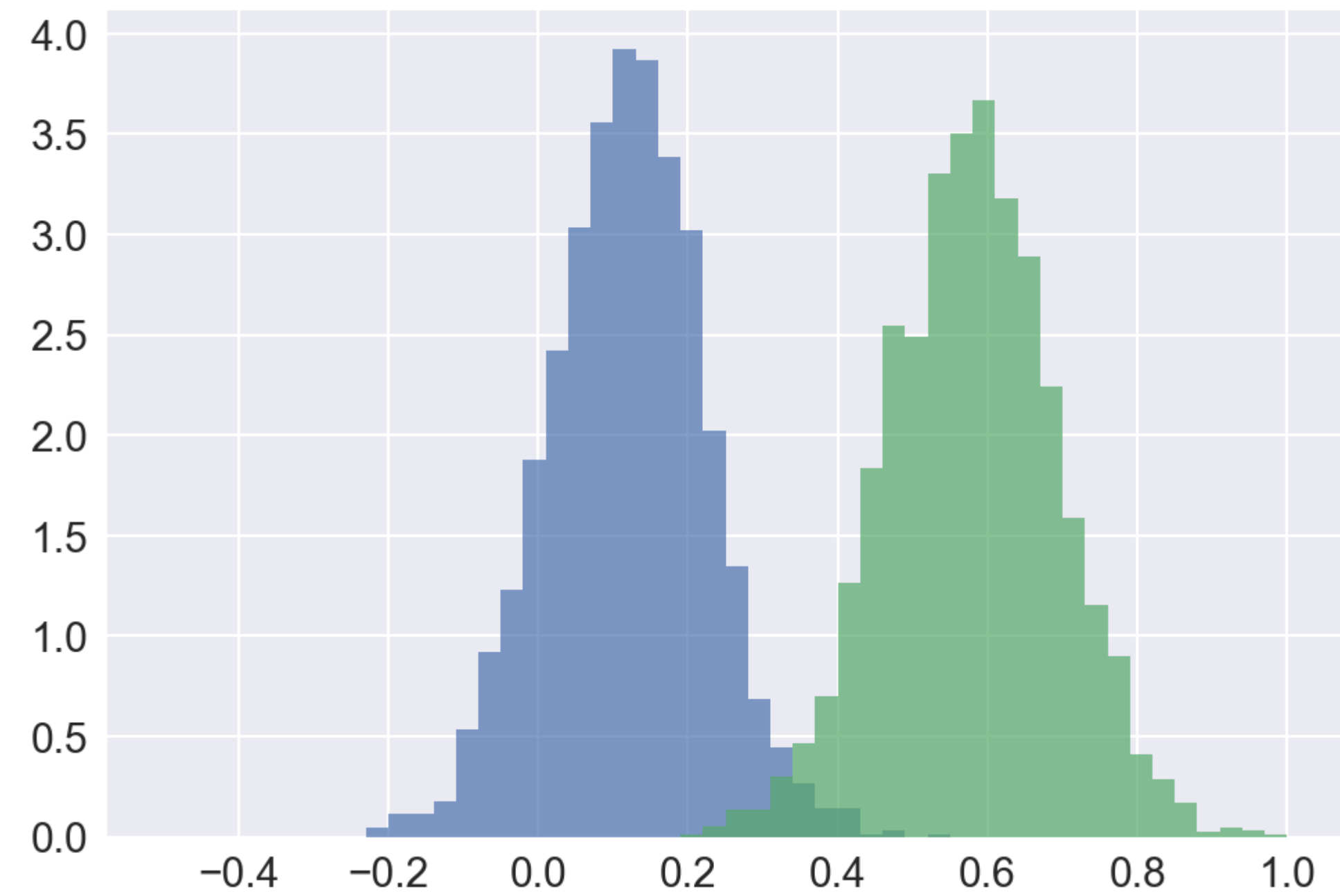
```
start = {'μ_1': 0., 'σ_1': 1., 'μ_2': 0., 'σ_2': 1.}
sampler = smp.NUTS(logp, start)
chain = sampler(6100, burn=100, thin=2)
```

Progress: [#####] 6100 of 6100 samples

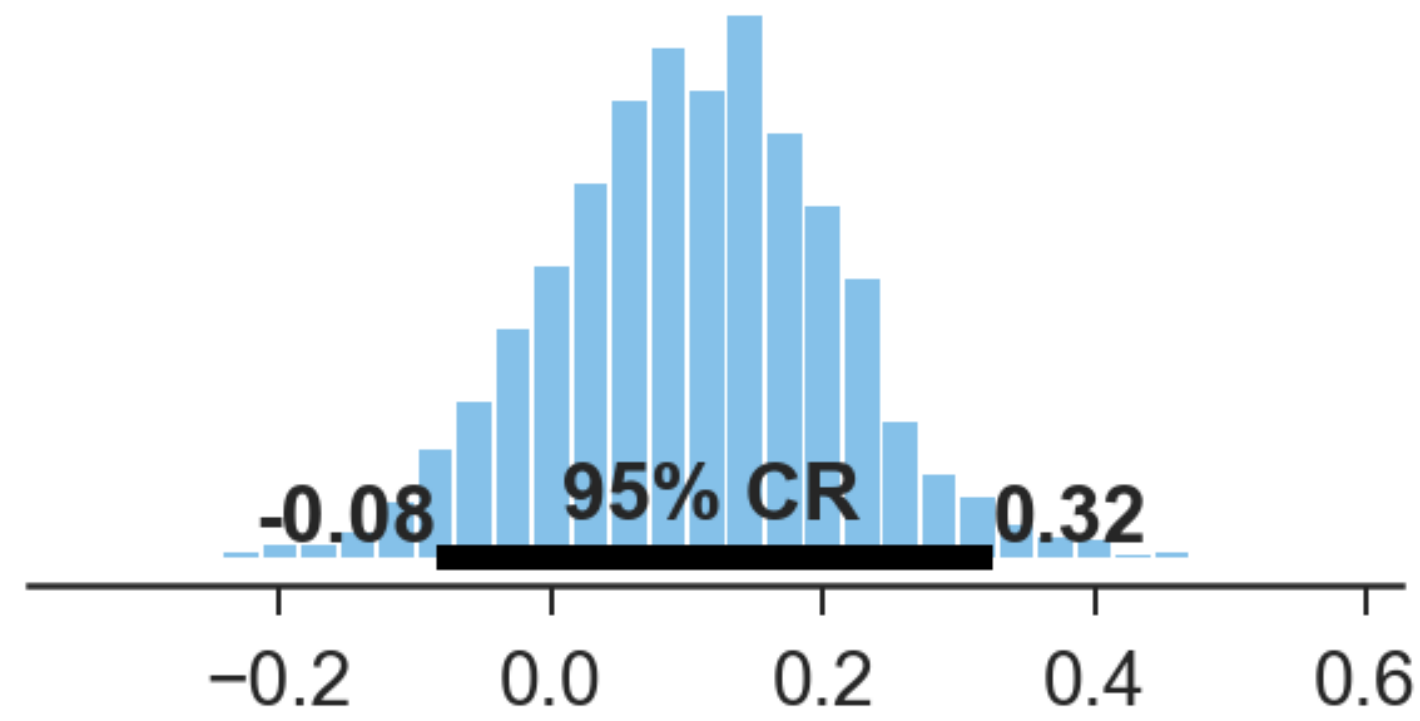
Posterior chains



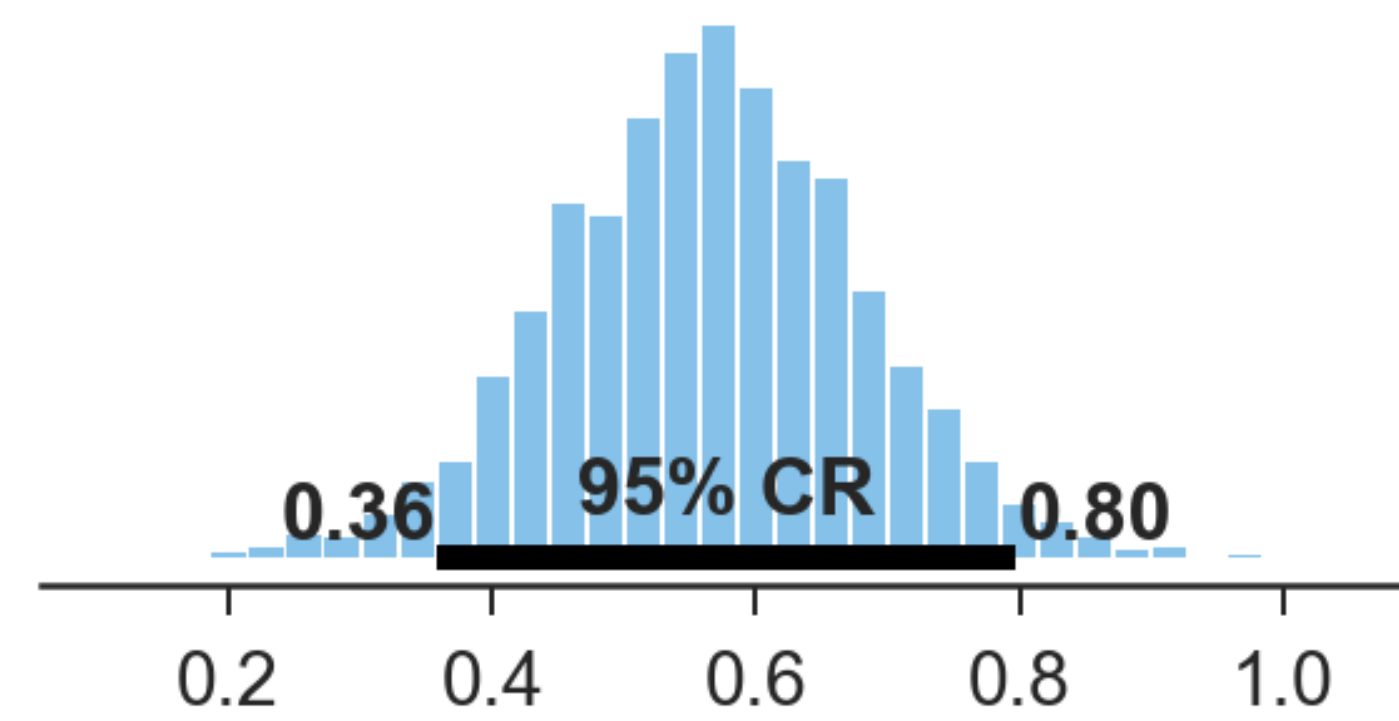
Posterior distributions



Group 1 mean
mean = 0.12

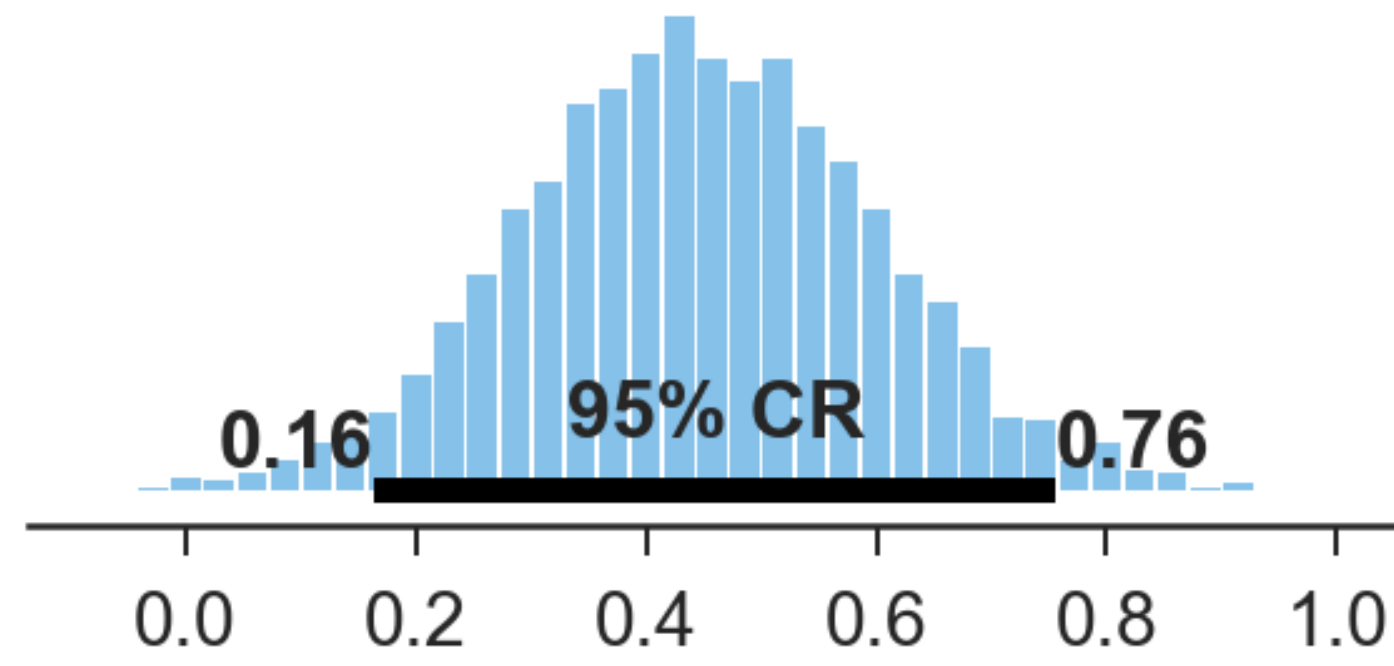


Group 2 mean
mean = 0.58

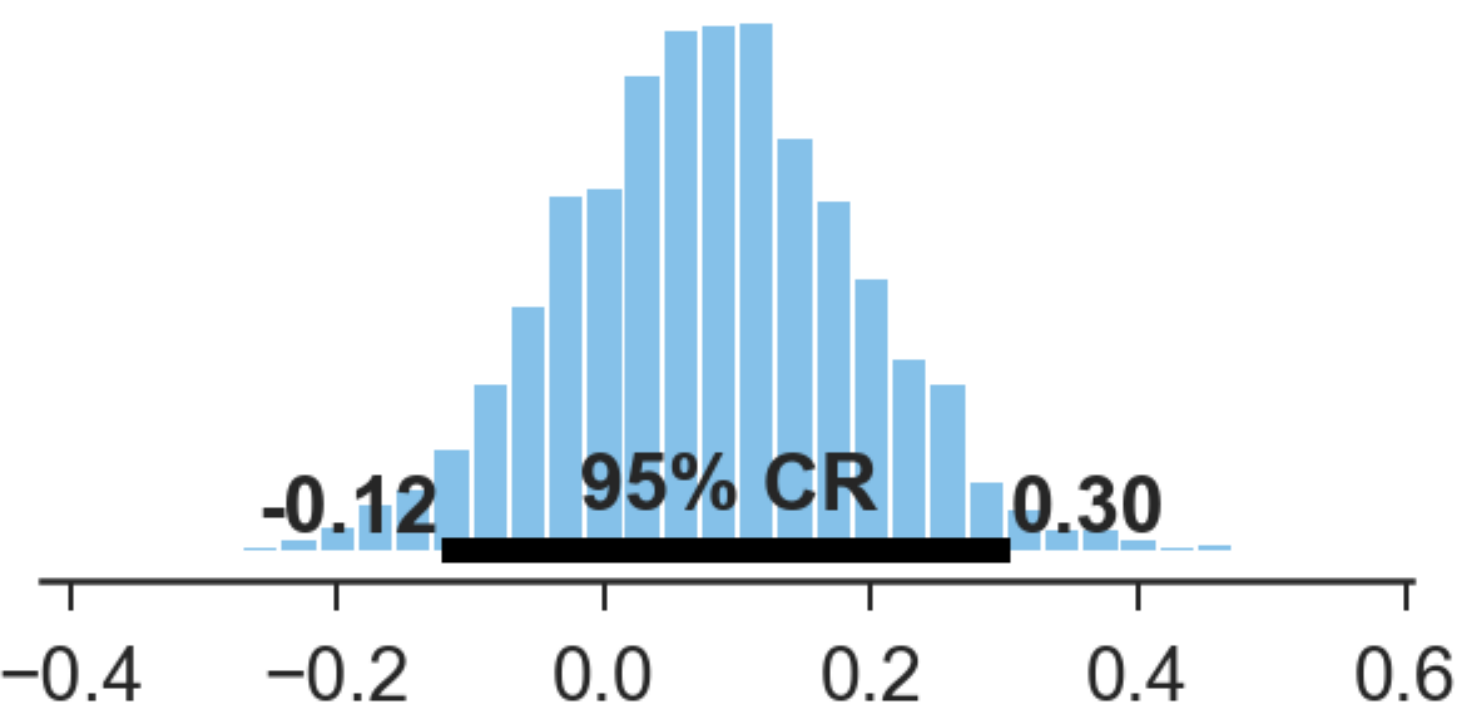


Difference of means
mean = 0.46

0.03% < 0 < 99.97%



Difference of std devs
mean = 0.09



Boston Housing Prices

Linear regression, Bayesian style

From attributes predict median home value



	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	black	lstat	medv
ID														
1	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	4.98	24.0
2	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	9.14	21.6
4	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4
5	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	5.33	36.2
7	0.08829	12.5	7.87	0	0.524	6.012	66.6	5.5605	5	311	15.2	395.60	12.43	22.9
11	0.22489	12.5	7.87	0	0.524	6.377	94.3	6.3467	5	311	15.2	392.52	20.45	15.0
12	0.11747	12.5	7.87	0	0.524	6.009	82.9	6.2267	5	311	15.2	396.90	13.27	18.9
13	0.09378	12.5	7.87	0	0.524	5.889	39.0	5.4509	5	311	15.2	390.50	15.71	21.7
14	0.62976	0.0	8.14	0	0.538	5.949	61.8	4.7075	4	307	21.0	396.90	8.26	20.4
15	0.63796	0.0	8.14	0	0.538	6.096	84.5	4.4619	4	307	21.0	380.02	10.26	18.2

Linear Model

Define the model

Likelihood

$$\hat{y} = \sum_i^m \beta_i x_i + \alpha$$
$$y \sim \text{Normal}(\hat{y}, \epsilon)$$

Priors

$$\beta \sim \text{Uniform}(-10, 10)$$
$$\alpha \sim \text{Uniform}(0, 100)$$
$$\epsilon \sim \text{Half Cauchy}(\gamma = 5)$$

Linear Model

Define the model in SamPy

```
# Here,  $\beta$  is an array of coefficients
def logp( $\beta$ ,  $\alpha$ ,  $\epsilon$ ):

    model = smp.Model()

    model.add(smp.half_cauchy( $\epsilon$ , 5))

    model.add(smp.uniform( $\beta$ , lower=-10, upper=10),
              smp.uniform( $\alpha$ , lower=0, upper=100))

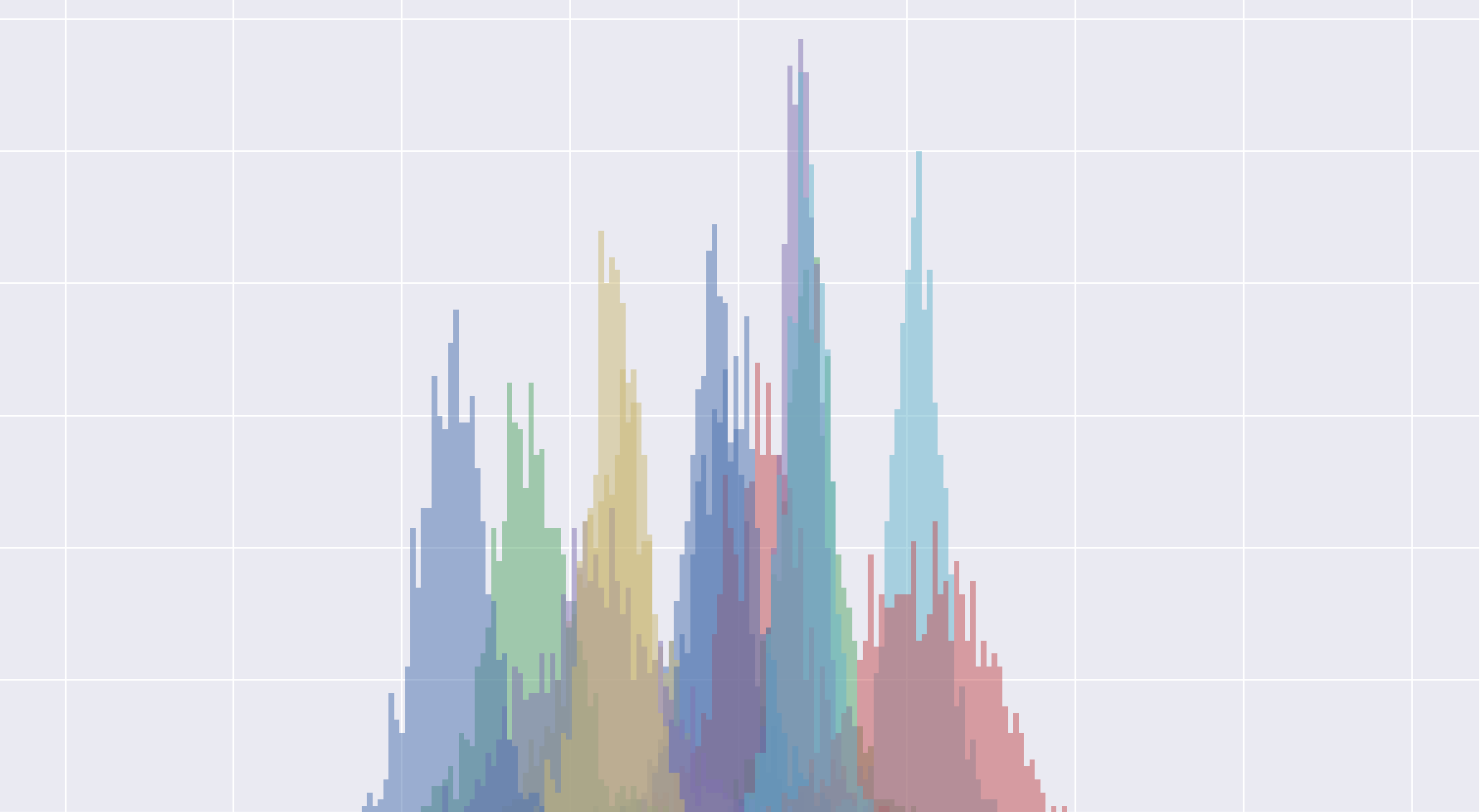
    estimate = np.dot(X,  $\beta$ ) +  $\alpha$ 
    model.add(smp.normal(Y, mu=estimate, sig= $\epsilon$ ))

    return model()
```

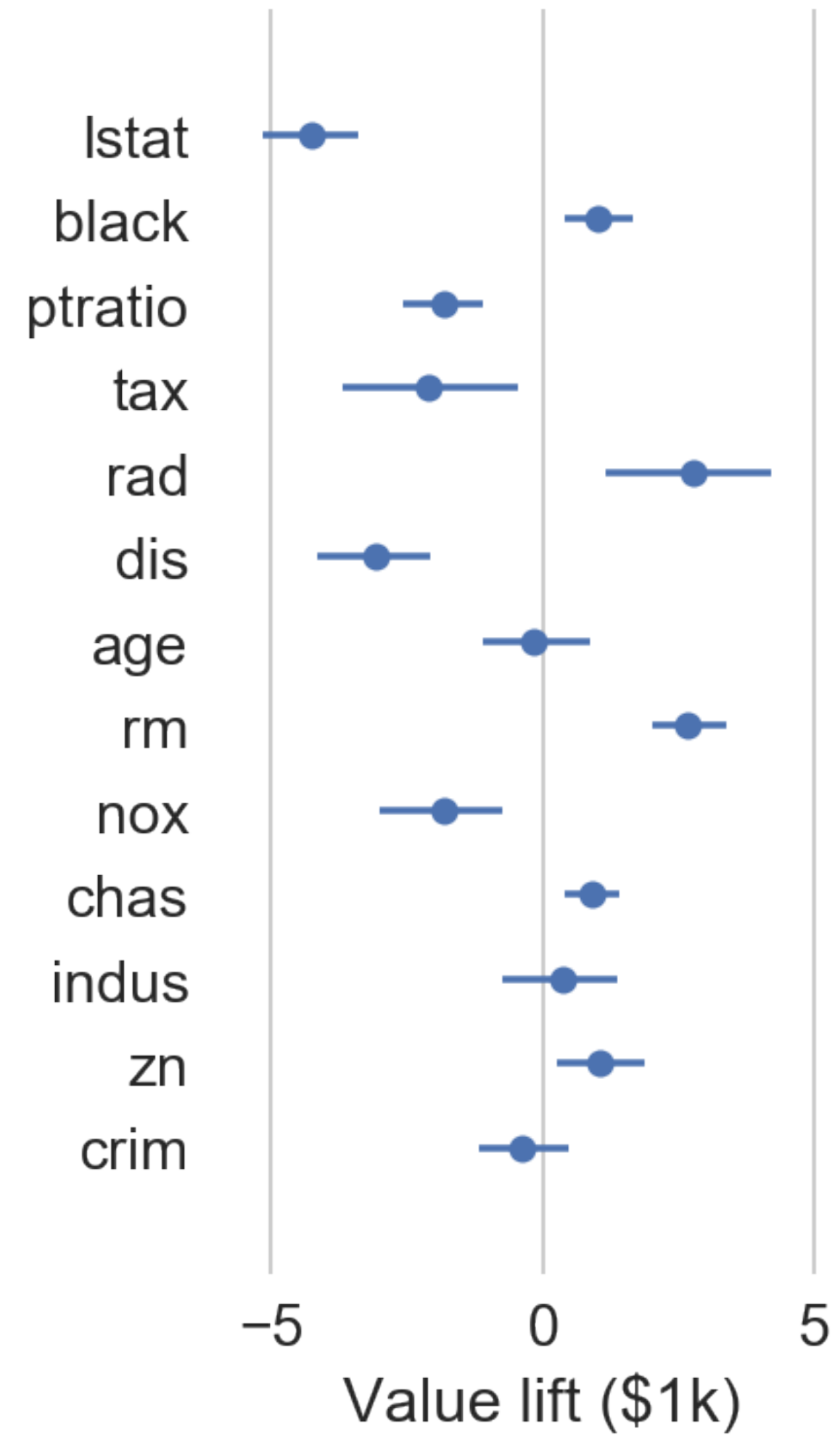
Linear Model

Sample from the posterior

```
start = {' $\beta$ ': np.zeros(X.shape[1]), ' $\alpha$ ': 0., ' $\epsilon$ ': 1.}  
sampler = smp.Slice(logp, start)  
chain = sampler(10000, burn=1000, thin=3)
```

Model Coefficients



Linear Model

Making Predictions

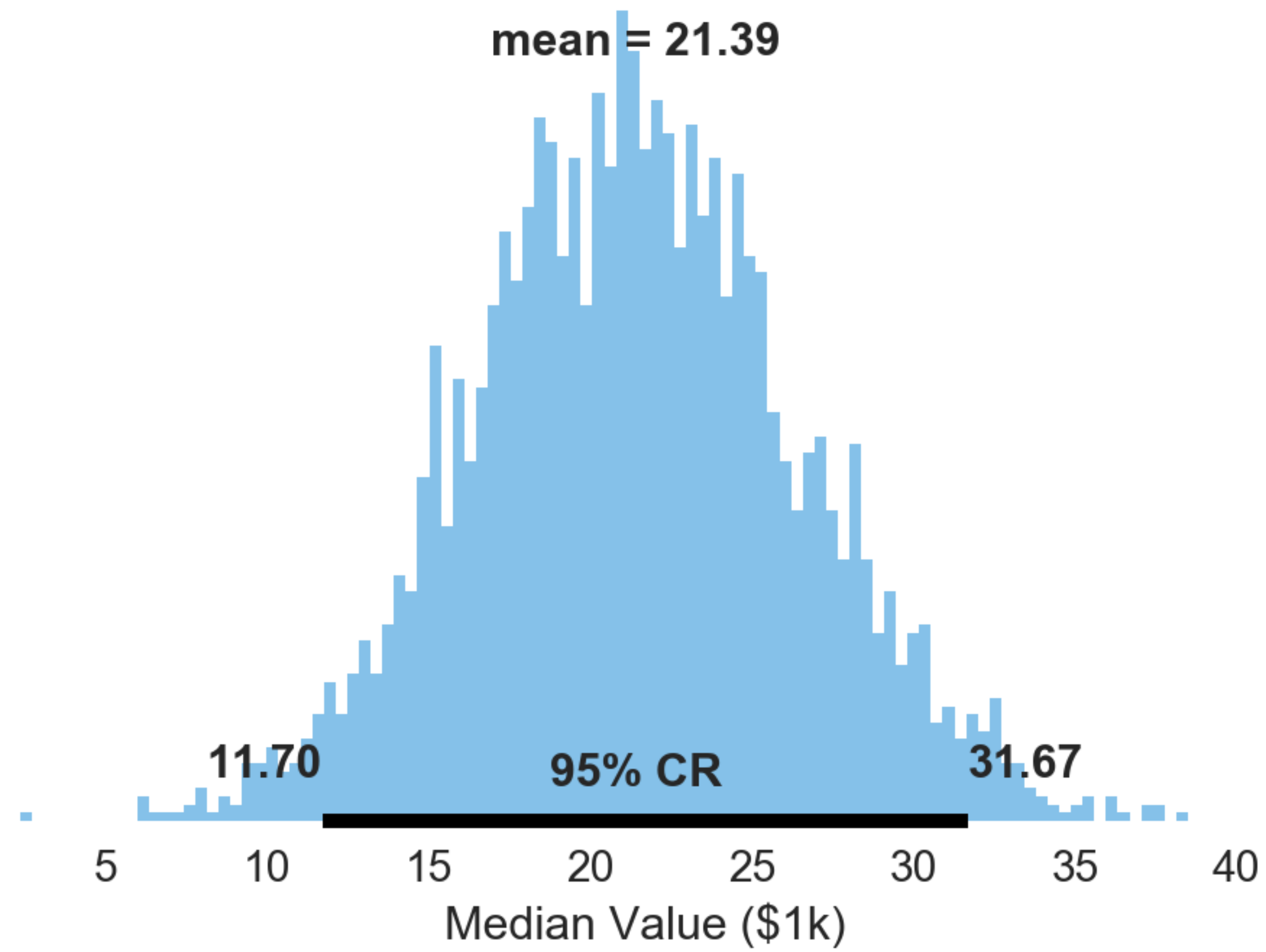
```
# Get one random input
predict_data = np.random.randn(X.shape[1])

idx = np.random.choice(np.arange(chain.size))
post_sample = chain[idx]

estimate = np.dot(predict_data, post_sample.β) + post_sample.α
single_predict = stats.norm(loc=estimate, scale=post_sample.ε).rvs()

# Now make posterior predictive distribution to get uncertainty
estimate = np.dot(predict_data, chain.β.T) + chain.α
post_predict = stats.norm(loc=estimate, scale=chain.ε).rvs()
```

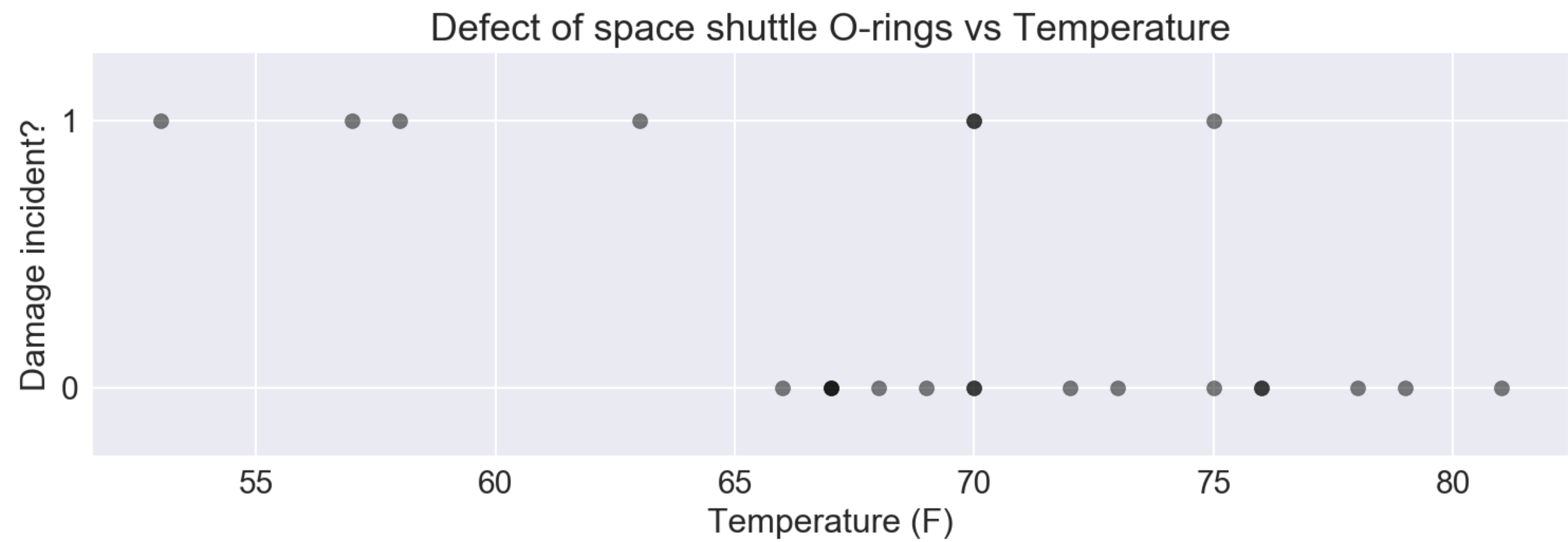
Posterior predictive distribution



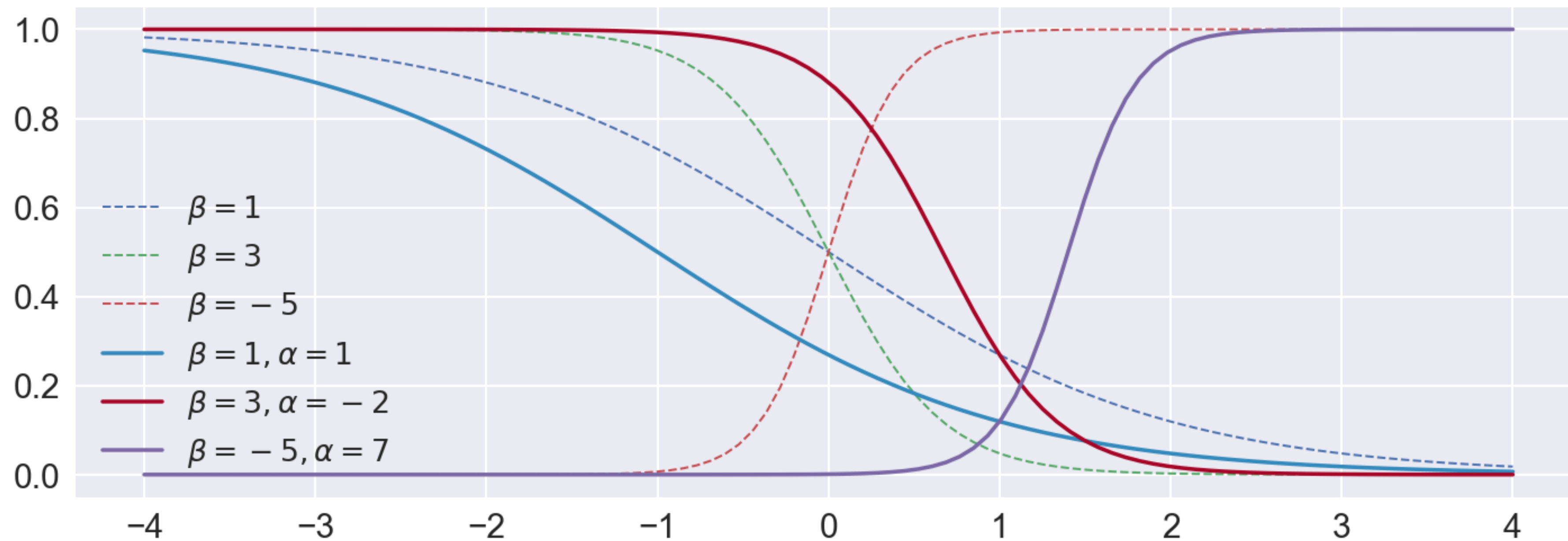
Challenger Disaster

Logistic regression, Bayesian style





The logistic function $\sigma = \frac{1}{1 + e^{\beta t + \alpha}}$



Logistic Model

Define the model

Likelihood

$$\sigma = \frac{1}{1 + e^x}$$

$$p = \sigma(\beta t + \alpha)$$

$$y \sim \text{Bernoulli}(p)$$

Priors

$$\alpha \sim \text{Normal}(0, 1000)$$

$$\beta \sim \text{Normal}(0, 1000)$$

Logistic Model

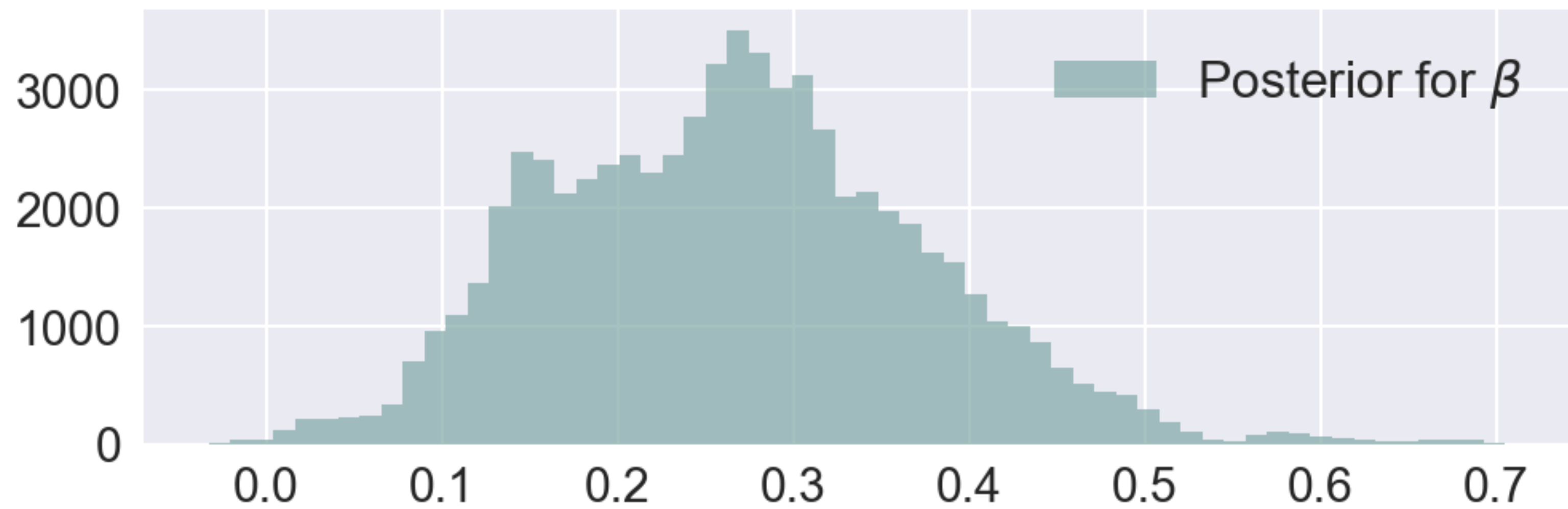
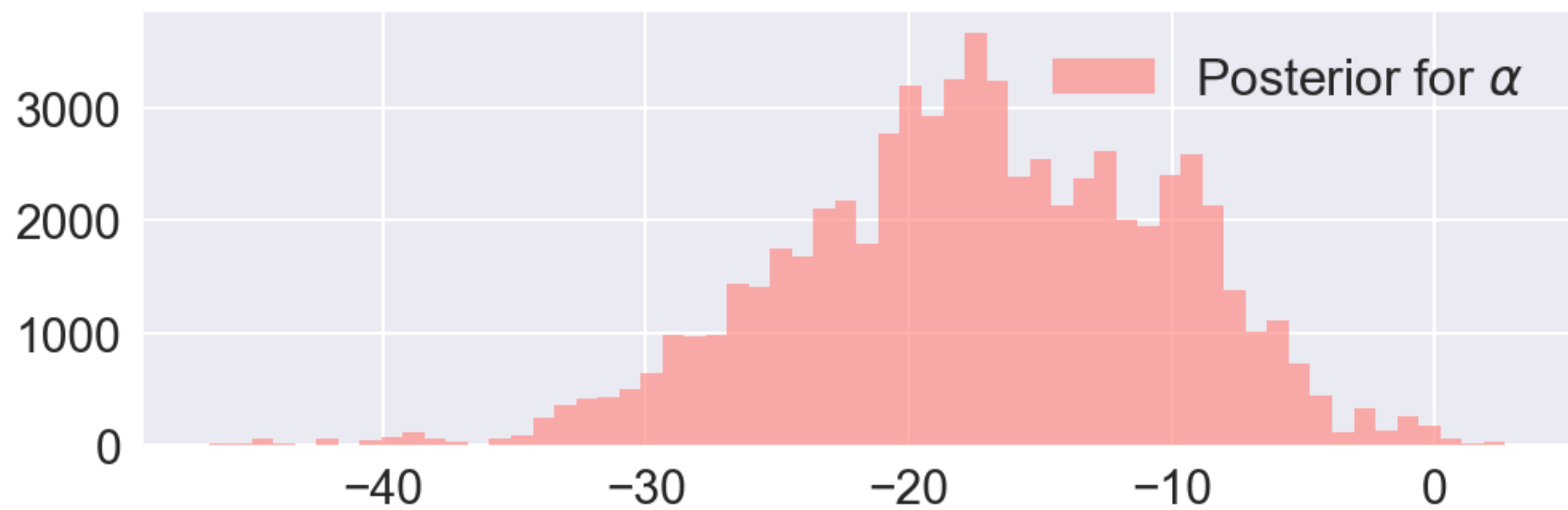
Define the model in SamPy

```
def logistic(x):  
    return 1/(1 + np.exp(x))  
  
def logp( $\beta$ ,  $\alpha$ ):  
  
    model = smp.Model()  
  
    # Wide normal priors on coefficients  
    model.add(smp.normal( $\alpha$ , sig=1000),  
              smp.normal( $\beta$ , sig=1000))  
  
    # Add log-likelihood  
    model.add(smp.bernoulli(damage, logistic(temperature* $\beta$  +  $\alpha$ )))  
  
    return model()
```

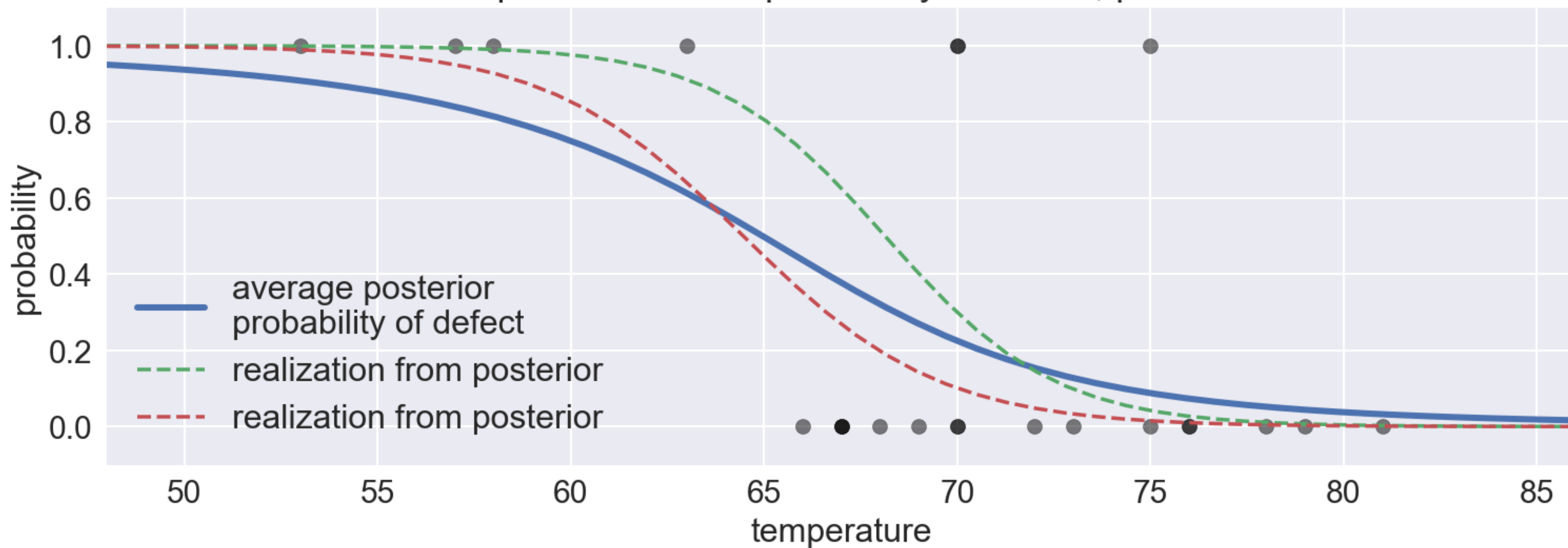
Logistic Model

Sample from the posterior

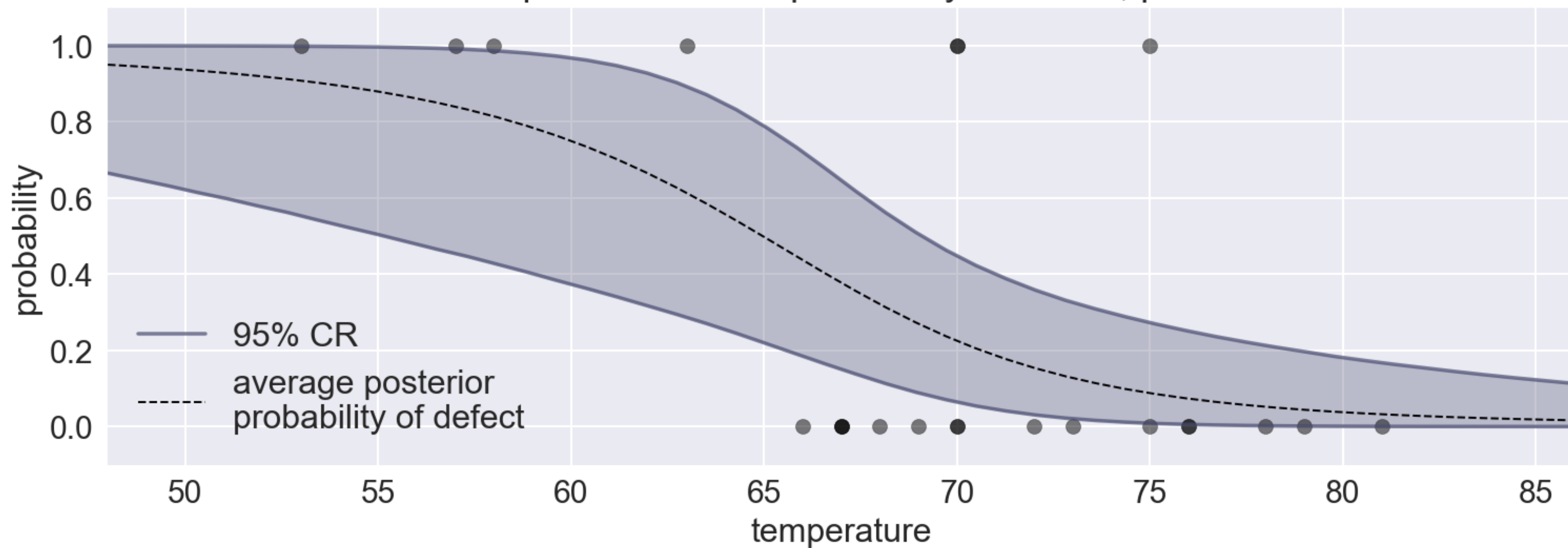
```
start = {' $\beta$ ': 0., ' $\alpha$ ': 0.}
sampler = smp.Slice(logp, start)
chain = sampler(220000, burn=20000, thin=3)
```



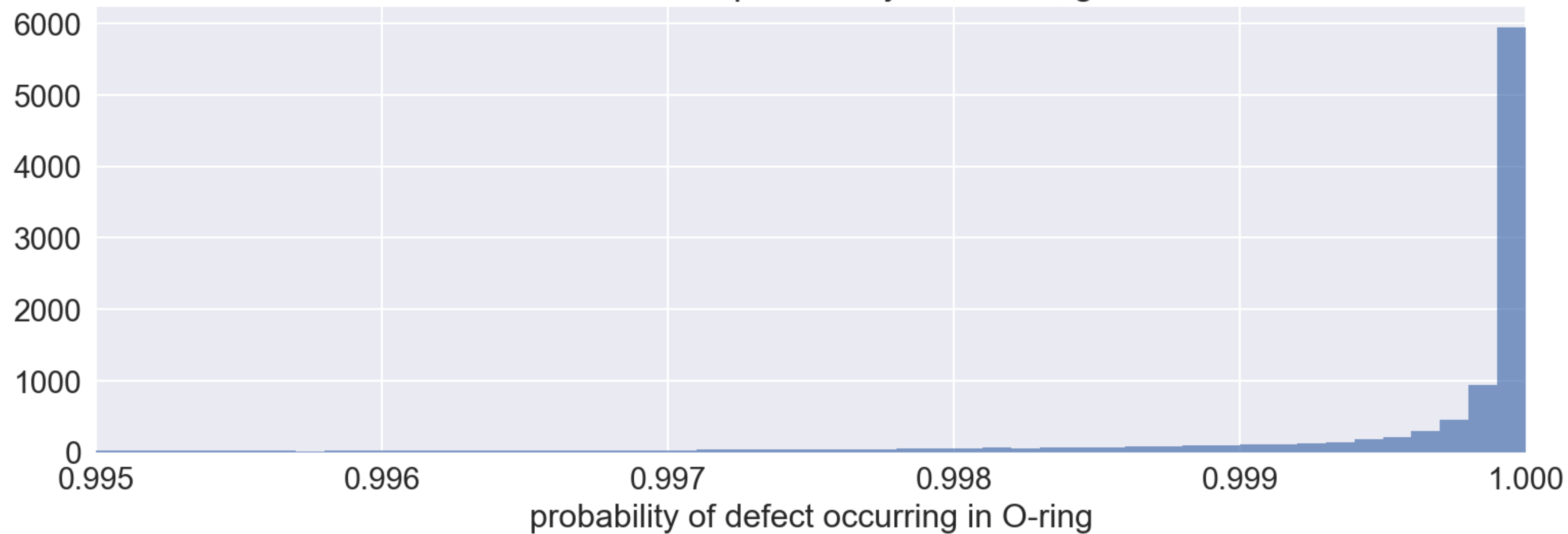
Posterior expected value of probability of defect; plus realizations



Posterior expected value of probability of defect; plus 95% CR



Posterior distribution of probability of defect, given $t = 31$



Thank you!

Questions?