# Bayesian Statistics

# Agenda

Part I - Generative Models

Part II - Bayesian Analysis

# Part I - Generative Models

$$p(y = c | \mathbf{x}, \boldsymbol{\theta}) \propto p(\mathbf{x} | y = c, \boldsymbol{\theta}) p(y = c | \boldsymbol{\theta})$$

posterior          likelihood          prior

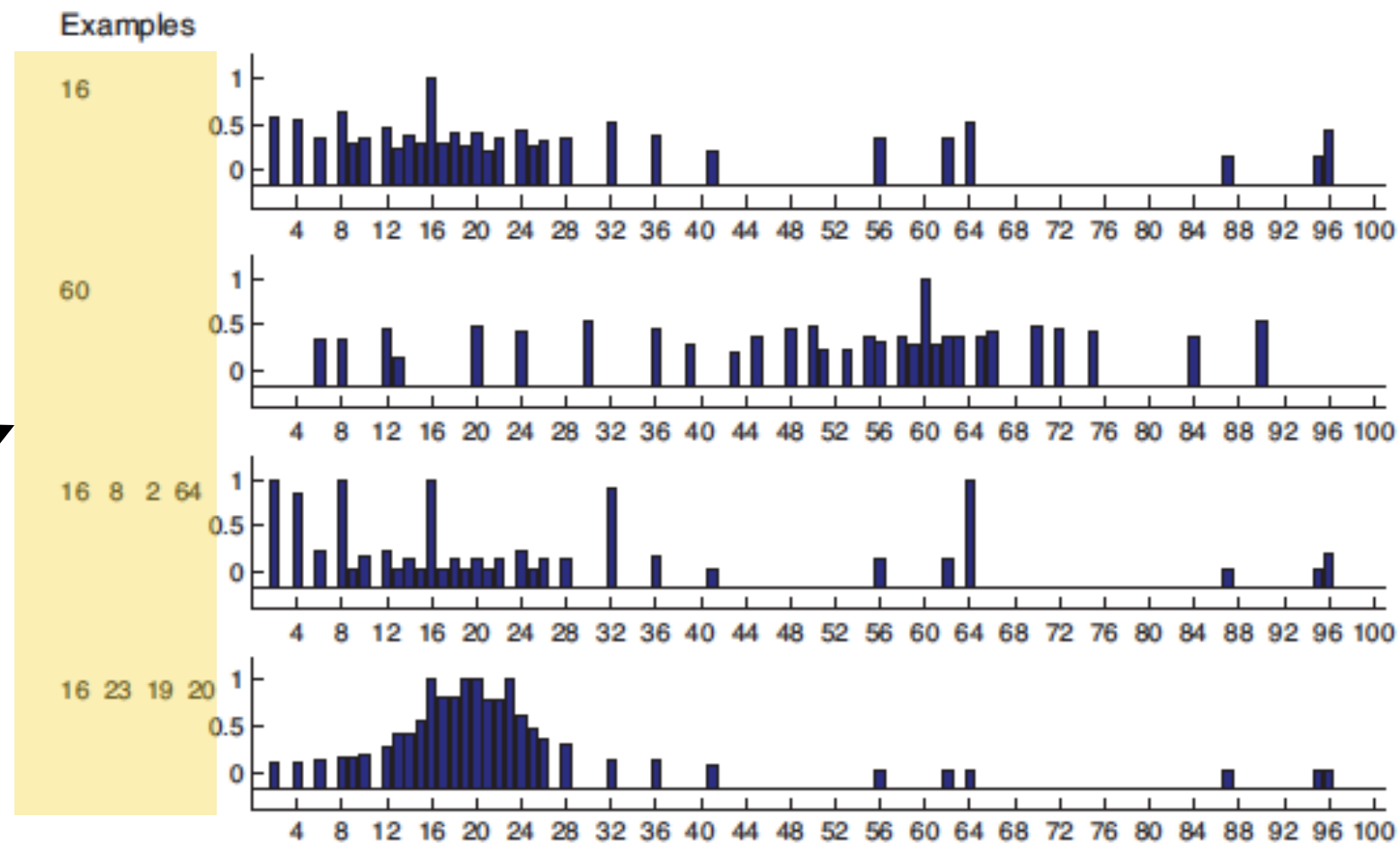# The numbers game - human estimates

The game proceeds as follows.

I choose some simple arithmetical concept C , such as "prime number" or "a number between 1 and 10".

I then give you a series of randomly chosen positive examples D = {x1, . . . , xN} drawn from C , and ask you whether some new test case ˜x  belongs to C

i.e., I ask you to classify ˜x

# The numbers game - human estimates

Say D = { 16, 8, 2, 64 }

Am I taking C = "powers of 2", or C = "even numbers"?

(both are consistent with the data)

Say D = { 16, 8, 2, 64 }

Am I taking C = "powers of 2", or C = "even numbers"?

(both are consistent with the data)

**Occam's razor**

$$p(\mathcal{D}|h) = \left[\frac{1}{\text{size}(h)}\right]^N = \left[\frac{1}{|h|}\right]^N$$
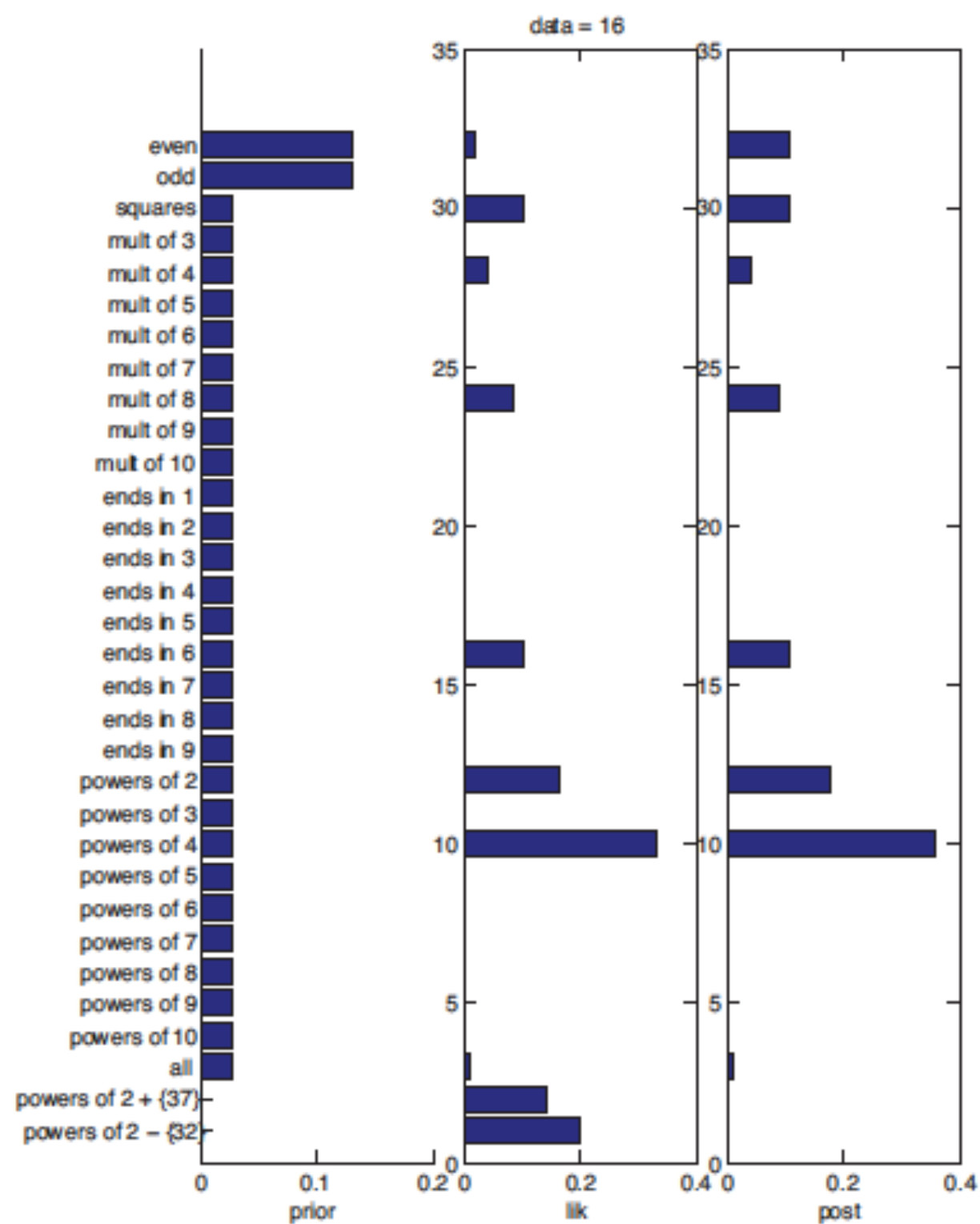
**likelihood**

# prior

h = "powers of 2"        h' = "powers of 2, except 32"

What now?

# prior

h = "powers of 2"     h' = "powers of 2, except 32"

What now?

the prior p(h) is the **subjective** part of Bayesian reasoning

$$p(h \mid D) = p(D \mid h) \cdot p(h) / \text{normalisation}$$

posterior     likelihood     prior

data = 16

# Recall Naïve Bayes

We want p(y=class | D=words)

What is likelihood p(D | y=c) ?

What is prior p(y = c)?

e.g.  Mary had a little lamb, little lamb, little lamb,
Mary had a little lamb, its fleece as white as snow

# Beta-Binomial

## Counting dice head, unknown bias

likelihood $\quad \text{Bin}(k|n,\theta) \triangleq \binom{n}{k} \theta^k (1-\theta)^{n-k}$

prior $\quad \text{Beta}(\theta|a,b) \propto \theta^{a-1}(1-\theta)^{b-1}$

posterior

$$p(\theta|\mathcal{D}) \propto \text{Bin}(N_1|\theta, N_0+N_1)\text{Beta}(\theta|a,b)$$

$$\propto \text{Beta}(\theta|N_1+a, N_0+b)$$

# Bayesian inference for online learning

$$p(\theta|\mathcal{D}_a, \mathcal{D}_b) \quad \propto \quad p(\mathcal{D}_b|\theta)p(\theta|\mathcal{D}_a)$$

$$\propto \quad \text{Bin}(N_1^b|\theta, N_1^b + N_0^b)\text{Beta}(\theta|N_1^a + a, N_0^a + b)$$

$$\propto \quad \text{Beta}(\theta| N_1^a + N_1^b + a, N_0^a + N_0^b + b)$$

# Part II - Bayesian Analysis

Frequentist approach:

$$\hat{\theta}_{MLE} = \arg \max_{\theta} p(D|\theta)$$
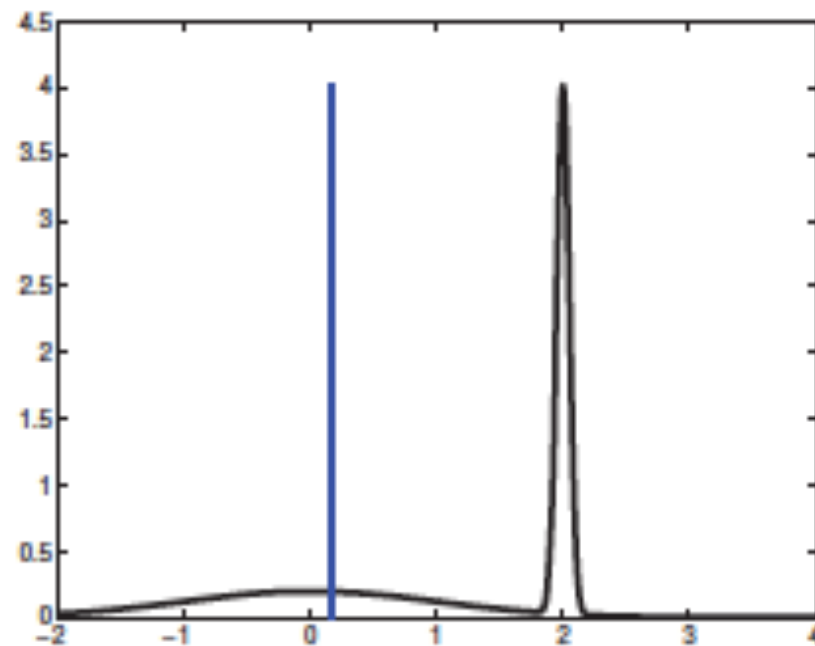
Bayesian approach:

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(D|\theta)p(\theta)$$
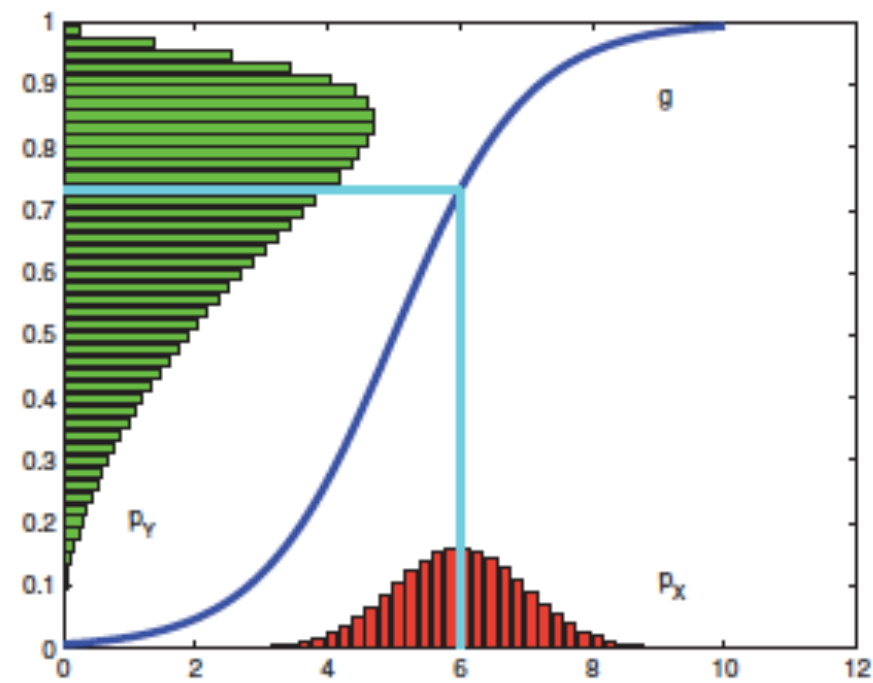
we care about the prior

# problems with a point estimate



# the mode is an untypical point

**Figure 5.2** Example of the transformation of a density under a nonlinear transform. Note how the mode of the transformed distribution is not the transform of the original mode. Based on Exercise 1.4 of (Bishop 2006b). Figure generated by `bayesChangeOfVar`.
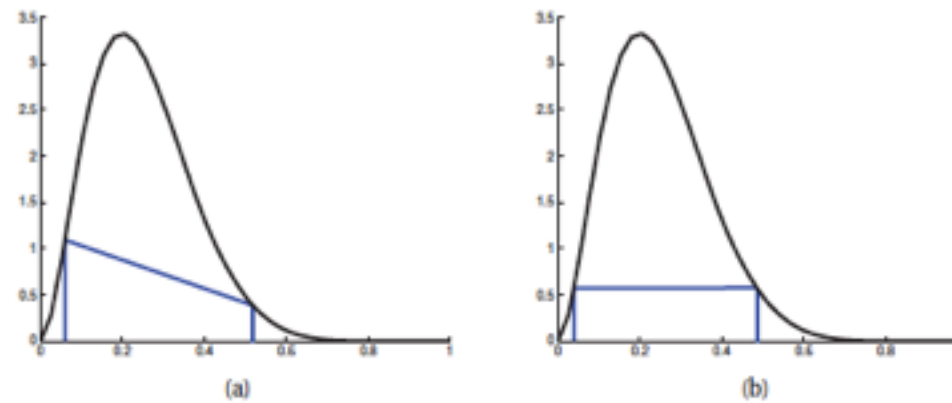
# credible interval (CI) vs HPD



**Figure 5.3** (a) Central interval and (b) HPD region for a Beta(3,9) posterior. The CI is (0.06, 0.52) and the HPD is (0.04, 0.48). Based on Figure 3.6 of (Hoff 2009). Figure generated by `betaHPD`.
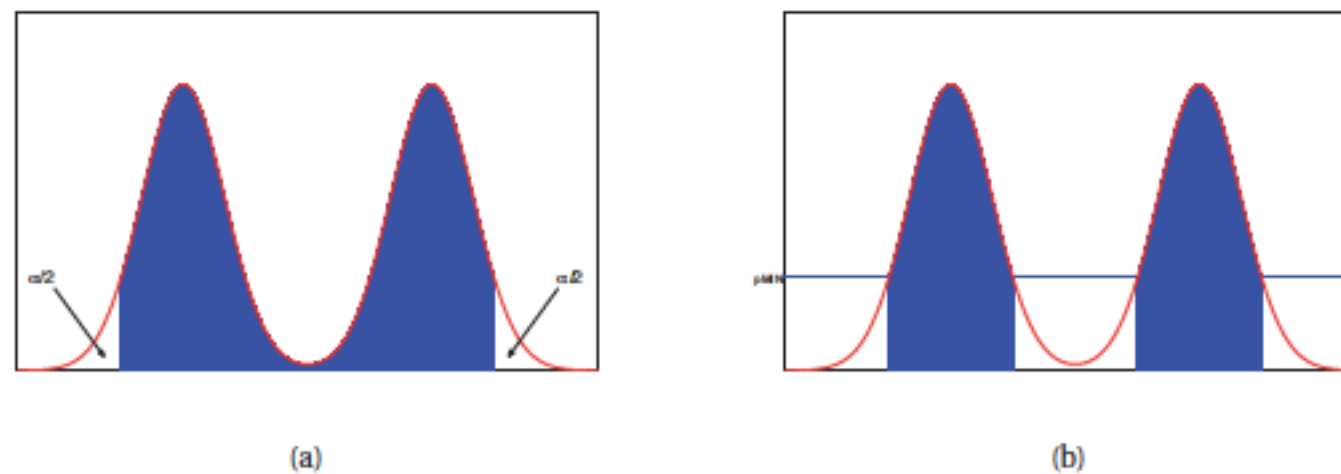


**Figure 5.4** (a) Central interval and (b) HPD region for a hypothetical multimodal posterior. Based on Figure 2.2 of (Gelman et al. 2004). Figure generated by `postDensityIntervals`.

# Going further than MAP

| Method | Definition |
|---|---|
| Maximum likelihood | $\hat{\theta} = \mathrm{argmax}_{\theta}\, p(\mathcal{D}|\theta)$ |
| MAP estimation | $\hat{\theta} = \mathrm{argmax}_{\theta}\, p(\mathcal{D}|\theta)p(\theta|\eta)$ |
| ML-II (Empirical Bayes) | $\hat{\eta} = \mathrm{argmax}_{\eta} \int p(\mathcal{D}|\theta)p(\theta|\eta)d\theta = \mathrm{argmax}_{\eta}\, p(\mathcal{D}|\eta)$ |
| MAP-II | $\hat{\eta} = \mathrm{argmax}_{\eta} \int p(\mathcal{D}|\theta)p(\theta|\eta)p(\eta)d\theta = \mathrm{argmax}_{\eta}\, p(\mathcal{D}|\eta)p(\eta)$ |
| Full Bayes | $p(\theta, \eta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta|\eta)p(\eta)$ |