



Computer
Science

CSC535: Probabilistic Graphical Models

Bayesian Probability and Inference

Prof. Jason Pacheco

Administrative Items

- HW1 grades / solutions posted
- For homeworks, please submit PDF *as separate file*
 - D2L makes it difficult to grade otherwise
- To generate uniform random integers in the range [1,6]

```
numpy.random.randint(1, 7)
```



Second argument is (exclusive)

This is the case for most Numpy functions that take intervals

What is Probability?

What does it mean that the probability of heads is $\frac{1}{2}$?



Two schools of thought...

Frequentist Perspective

Proportion of successes (heads) in repeated trials (coin tosses)

Bayesian Perspective

Belief of outcomes based on assumptions about nature and the physics of coin flips

Neither is better/worse, but we can compare interpretations...

Frequentist & Bayesian Modeling

We will use the following notation throughout:

θ - Unknown (e.g. coin bias)

y - Data

Frequentist

(Conditional Model)

$$p(y; \theta)$$

- θ is a non-random unknown parameter
- $p(y; \theta)$ is the *sampling / data generating distribution*

Bayesian

(Generative Model)

Prior Belief $\rightarrow p(\theta)p(y | \theta) \leftarrow$ Likelihood

- θ is a random variable (latent)
- Requires specifying $p(\theta)$ the prior belief

Bayes' Rule

Posterior represents all uncertainty after observing data...

The diagram illustrates Bayes' Rule with the following components and labels:

- prior probability**: Labeled in red, with an arrow pointing to the numerator term $p(\theta)$.
- likelihood function for the parameters**: Labeled in red, with an arrow pointing to the numerator term $p(y | \theta)$.
- posterior probability**: Labeled in red, with an arrow pointing to the entire left side of the equation $p(\theta | y)$.
- marginal likelihood or: evidence or: partition function or: normalizer**: Labeled in red, with an arrow pointing to the denominator term $p(y)$.

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

Bayes' Rule : Marginal Likelihood

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)} \propto \underbrace{p(\theta)p(y | \theta)}$$

Often hard to calculate

Often know this (the model)

Marginal likelihood integrates (marginalizes) over unknown θ :

$$p(y) = \int p(\theta)p(y | \theta) d\theta$$

Marginal likelihood is less problematic in discrete models (not always)

This integral often lacks a closed form and cannot be computed...

Aside : Proportionality

Recall PMF / PDF must sum / integrate to 1,

$$\begin{array}{cc} \text{PMF} & \text{PDF} \\ \sum_x p(x) = 1 & \int p(x) dx = 1 \end{array}$$

May only know distribution constant that does not depend on RV x ,

$$\int \tilde{p}(x) dx = \mathcal{Z} \quad \text{so} \quad p(x) \propto \tilde{p}(x)$$

Properly normalized distribution by dividing our normalization constant:

$$\int p(x) dx = \int \frac{1}{\mathcal{Z}} \tilde{p}(x) dx = \frac{1}{\int \tilde{p}(x) dx} \int \tilde{p}(x) dx = 1$$

Aside : Proportionality

Example Let X be a Bernoulli RV (coinflip) with probabilities *proportional to*:

$$\tilde{p}(X = 0) = 0.5$$

$$\tilde{p}(X = 1) = 1.5$$

Greater than 1, but
It is an *unnormalized*
probability

Compute normalization constant,

$$\mathcal{Z} = \tilde{p}(X = 0) + \tilde{p}(X = 1) = 2.0$$

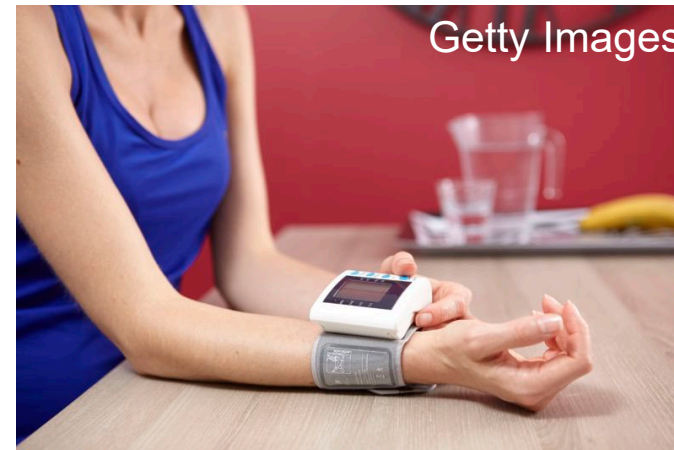
Normalize probability distribution,

$$p(X) = \frac{1}{\mathcal{Z}} \tilde{p}(X) = \begin{pmatrix} 1/4 \\ 3/4 \end{pmatrix}$$

Sums to 1

Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



A recent home test states that you have high BP. Should you start medication?

An Assessment of the Accuracy of Home Blood Pressure Monitors When Used in Device Owners

Jennifer S. Ringrose,¹ Gina Polley,¹ Donna McLean,²⁻⁴ Ann Thompson,^{1,5} Fraulein Morales,¹ and Raj Padwal^{1,4,6}

Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



- Latent quantity of interest is hypertension: $\theta \in \{true, false\}$
- Measurement of hypertension: $y \in \{true, false\}$
- Prior: $p(\theta = true) = 0.29$
- Likelihood: $p(y = true \mid \theta = false) = 0.30$
 $p(y = true \mid \theta = true) = 1.00$

Bayesian Inference Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and **no false negative error**.



Suppose we get a positive measurement, then posterior is:


$$\begin{aligned} p(\theta = true \mid y = true) &= \frac{p(\theta = true)p(y = true \mid \theta = true)}{p(y = true)} \\ &= \frac{0.29 * 1.00}{0.29 * 1.00 + 0.71 * 0.30} \approx 0.58 \end{aligned}$$

What conclusions can be drawn from this calculation?

Bayesian Updating

Suppose we plan to take another test...

Question What is our belief about blood pressure status *before* the second test?

(a) Posterior: $p(\theta = \text{true} \mid y_1 = \text{true})$ 

(b) Likelihood: $p(y_1 = \text{true} \mid \theta = \text{true})$


(c) Marginal Likelihood: $p(y_1 = \text{true})$

Bayesian Updating

Suppose we plan to take another test...

Question What is the probability that we get *true* on the second test if we have high blood pressure?

(a) Posterior: $p(\theta = \text{true} \mid y_1 = \text{true}, y_2 = \text{true})$

(b) Likelihood: $p(y_2 = \text{true} \mid \theta = \text{true})$ 

(c) Marginal Likelihood: $p(y_2 = \text{true})$

Why not: $p(y_2 = \text{true} \mid \theta = \text{true}, y_1 = \text{true})$

Bayesian Updating

Suppose we plan to take another test...

Question What is the probability that we get *true* on the second test if we have high blood pressure?

(a) Posterior: $p(\theta = \text{true} \mid y_1 = \text{true}, y_2 = \text{true})$

(b) Likelihood: $p(y_2 = \text{true} \mid \theta = \text{true})$ ←

(c) Marginal Likelihood: $p(y_2 = \text{true})$

Because $y_1 \perp y_2 \mid \theta$ *so* $p(y_2 \mid \theta, y_1) = p(y_2 \mid \theta)$

Bayesian Updating

Suppose we receive another positive test $y_2 = true \dots$

Posterior belief given *both* tests is then,

$$p(\theta = true \mid y_1 = true, y_2 = true) =$$

$$= \frac{p(\theta = true \mid y_1 = true)p(y_2 = true \mid \theta)}{p(y_2 = true \mid y_1 = true)}$$

**Probability of getting
two positive tests
regardless of BP status**

$$\propto p(\theta = true \mid y_1 = true)p(y_2 = true \mid \theta = true)$$

Inference from first test

Likelihood of positive test

Bayesian Updating

Consider two *conditionally independent* observations X_1 and X_2 , their joint distribution is:

$$p(\theta, X_1, X_2) = p(\theta)p(X_1 | \theta)p(X_2 | \theta) \stackrel{\text{Probability chain rule}}{=} p(\theta | X_1)p(X_1)p(X_2 | \theta)$$

So, conditioned on X_1 :

$$p(\theta, X_2 | X_1) = p(\theta | X_1)p(X_2 | \theta) \quad \xrightarrow{\text{Update prior belief after seeing } X_1}$$

This is proportional to the **full posterior** by Bayes' rule:

$$p(\theta | X_1, X_2) \propto p(\theta | X_1)p(X_2 | \theta) \quad \text{Normalizer is } p(X_2 | X_1)$$

Step 1: Do inference after seeing X_1

Step 2: Update posterior by multiplying likelihood of X_2

Bayesian Updating

Given conditionally independent X_1, \dots, X_N posterior belief is,

$$p(\theta \mid X_1, \dots, X_N)$$

Receive $N+1^{\text{th}}$ observation X_{N+1} and update posterior,

$$p(\theta \mid X_1, \dots, X_{N+1}) \propto p(\theta \mid X_1, \dots, X_N) p(X_{N+1} \mid \theta)$$

**Belief after seeing
N+1th observation**



**Belief before seeing
N+1th observation**



**Belief about
N+1th observation**



Updates are more complicated if observations are dependent...

Frequentist vs. Bayesian Inference

We have data X_1, \dots, X_N and want to infer unknown parameter θ

Frequentist Inference

The data *uniquely determines* θ , e.g. by the likelihood:

Not a distribution on parameter $p(X_1, \dots, X_N; \theta)$ **How well it explains the data**

Bayesian Inference

The data *updates our belief* about θ , which is random:

$$p(\theta \mid X_1, \dots, X_N) \propto p(\theta \mid X_1, \dots, X_{N-1})p(X_N \mid \theta)$$

Our belief changes with more data

Minimum Mean Squared Error (MMSE)

Posterior mean minimizes squared error,

$$\hat{\theta}^{\text{MMSE}} = \arg \min \mathbb{E}[(\hat{\theta} - \theta)^2 \mid y] = E[\theta \mid y]$$

- Minimizes error conditioned on observed data
- MMSE is an **unbiased estimator**
- MMSE is **asymptotically unbiased** and **asymptotically normal**,

$$\sqrt{N}(\hat{\theta}^{\text{MMSE}} - \theta) \rightarrow \mathcal{N}(0, \sigma^2)$$

Example: Beta-Bernoulli MMSE

Let $Y_1, \dots, Y_N \sim \text{Bernoulli}(\pi)$ and $\pi \sim \text{Beta}(\alpha, \beta)$.

- Beta is a distribution on probabilities $\pi \in [0, 1]$
- *Shape* parameters α and β with mean,

$$\mathbf{E}[\theta] = \frac{\alpha}{\alpha + \beta}$$

- Beta-Bernoulli has Beta posterior distribution,

$$p(\pi | X_1^N) = \text{Beta}(\alpha + \text{number of heads}, \beta + \text{number of tails})$$

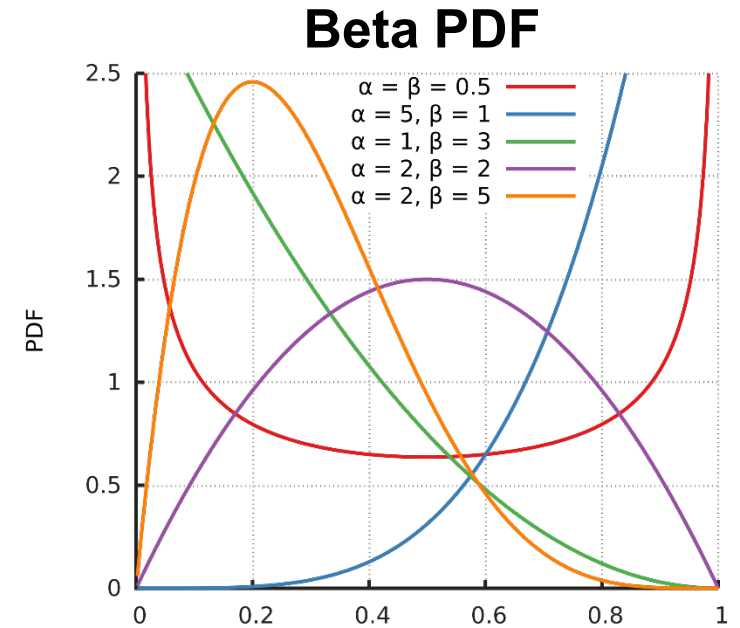
MMSE given by posterior mean,

Q What happens to MMSE when we have limited data?

$$\hat{\pi}^{\text{MMSE}} = \frac{\alpha + \text{number of heads}}{\alpha + \beta + N}$$

Prior belief (pseudo-heads)

Q What happens to MMSE when we have a lot of data?



Bayes Estimators

Minimizes expected loss function,

$$\hat{\theta} = \arg \min_{\hat{\theta}} \mathbf{E} \left[L(\theta, \hat{\theta}) \mid y \right]$$

Expected loss referred to as *Bayes risk*.

MMSE minimizes squared-error loss $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$

Minimum absolute error (MAE) is posterior *median*,

$$\arg \min \mathbf{E}[|\hat{\theta} - \theta| \mid y] = \text{median}(\theta \mid y)$$

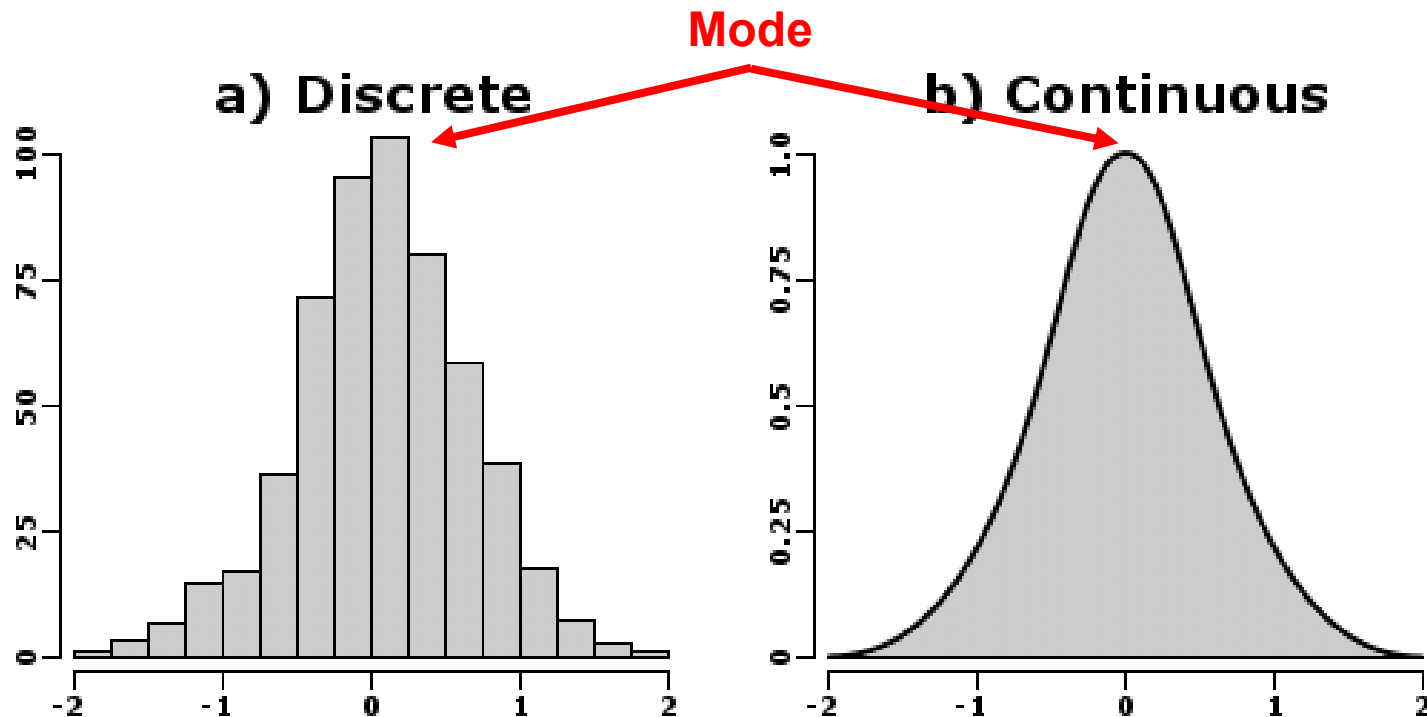
Note: Same answer for linear function: $L(\theta, \hat{\theta}) = c|\hat{\theta} - \theta|$

Maximum a Posteriori (MAP)

Very common to produce maximum probability estimates,

$$\hat{\theta}^{\text{MAP}} = \arg \max p(\theta | y)$$

*MAP is the **mode** (highest probability outcome) of the posterior*



Maximum a Posteriori (MAP)

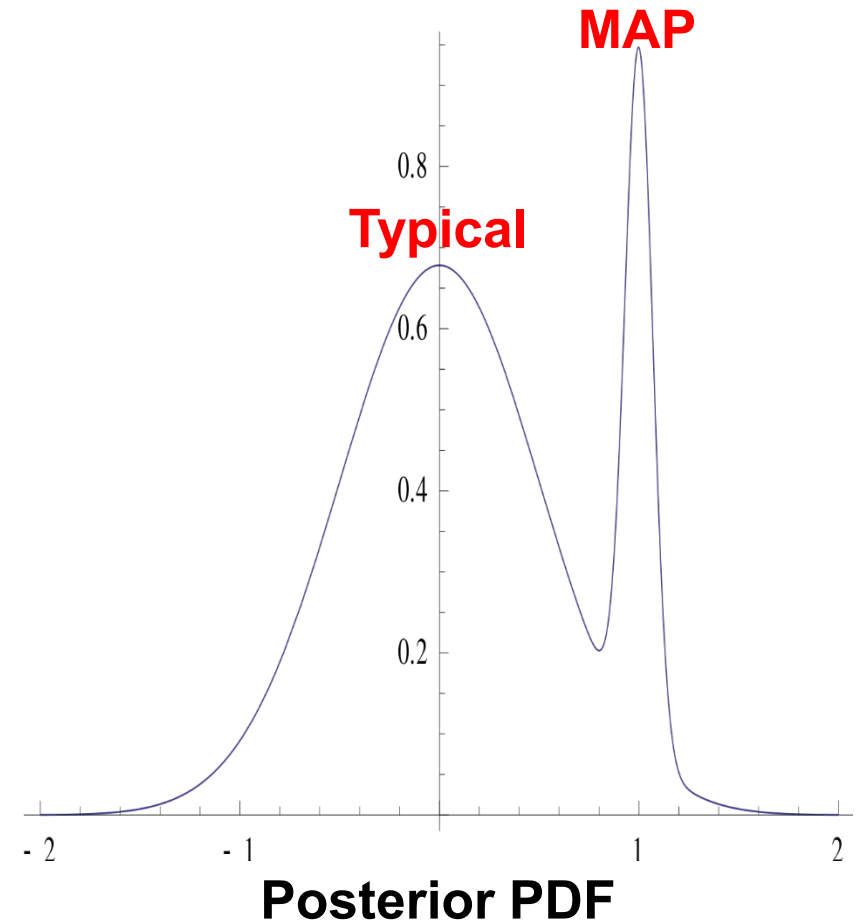
MAP (mode) may not be representative of typical outcomes

Also, not a Bayes estimator (unless discrete),

$$\lim_{c \rightarrow 0} L(\theta, \hat{\theta}) = \begin{cases} 0, & \text{if } |\hat{\theta} - \theta| < c \\ 1, & \text{otherwise} \end{cases}$$

Degenerate loss function

Despite its issues, MAP is frequently used in “Bayesian” inference and estimation



Example: Beta-Bernoulli MAP

Let $X_1, \dots, X_N \sim \text{Bernoulli}(\pi)$ and $\pi \sim \text{Beta}(\alpha, \beta)$ then posterior is,

$$p(\pi \mid X_1^N) = \text{Beta}(\alpha + \underbrace{\text{number of heads}}_{N_H}, \beta + \text{number of tails})$$

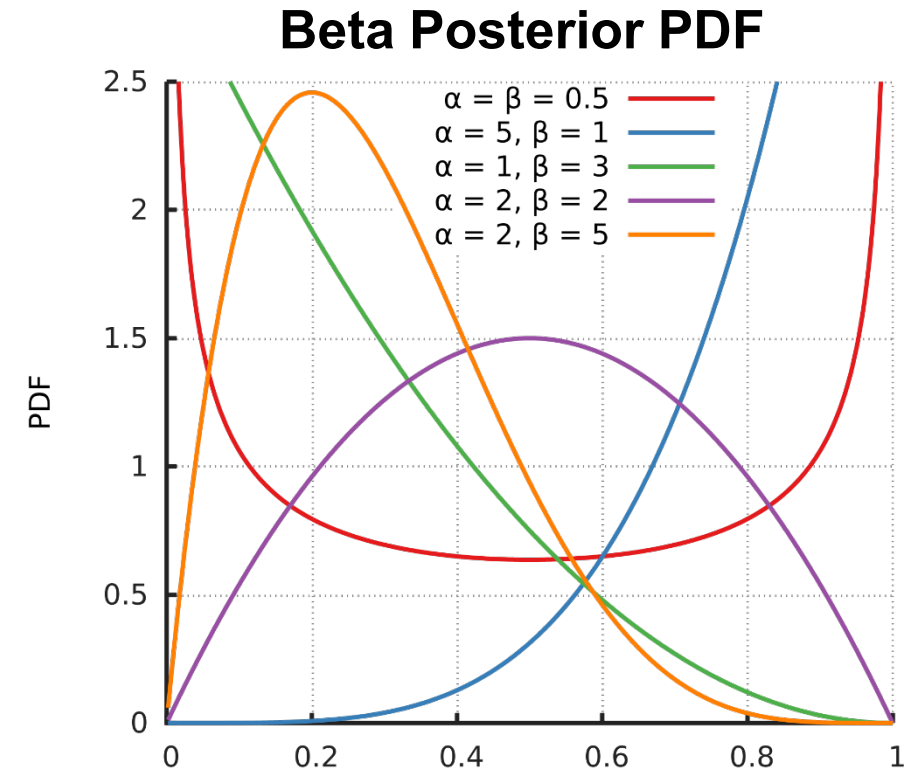
Highest probability (mode) of Beta given by,

Take derivative,
set to zero, solve.

$$\hat{\pi}^{\text{MAP}} = \frac{\alpha + N_H - 1}{\alpha + \beta + N - 2}$$

Beta distribution is not always convex!

- MAP is any value for $\alpha = \beta = 1$
- Two modes (bimodal) for $\alpha, \beta < 1$



Maximum a Posteriori (MAP)

Equivalent to maximizing joint probability,

$$\arg \max_{\theta} p(\theta | y) = \arg \max_{\theta} \frac{p(\theta, y)}{p(y)} = \arg \max_{\theta} p(\theta, y)$$

Constant

For iid y_1, \dots, y_N solve in log-domain (like *maximum likelihood est.*),

$$\hat{\theta}^{\text{MAP}} = \arg \max_{\theta} \log p(\theta, y_1, \dots, y_N) = \underbrace{\sum_i \log p(y_i | \theta)}_{\substack{\text{Log-Likelihood} \\ \text{(how well it fits data)}}} + \underbrace{\log p(\theta)}_{\substack{\text{Log-Prior} \\ \text{(how well it} \\ \text{agrees with prior)}}$$

Intuition MAP is like MLE but with a “penalty” term (log-prior)

Bayes' Rule : Reminder

prior probability

likelihood function for the parameters

$$p(\theta | D) = \frac{p(\theta)p(D | \theta)}{p(D)}$$

posterior probability

normalizer, often is not of interest

The diagram shows the equation for Bayes' Rule. The term $p(\theta)$ is labeled as 'prior probability'. The term $p(D | \theta)$ is labeled as 'likelihood function for the parameters'. The term $p(D)$ is labeled as 'normalizer, often is not of interest'. The entire fraction is labeled as 'posterior probability'.

Bayes' Rule : Reminder

prior probability

likelihood function
for the parameters

$$p(\theta | D) \propto p(\theta)p(D | \theta)$$

posterior probability

Posterior is **proportional**
to the joint

Bayes' Rule : Reminder

prior probability

likelihood function
for the parameters

$$q(\theta | D) \propto p(\theta)\omega(D | \theta)$$

posterior probability

In general, distributions are
different functions

Bayes' Rule : Reminder

Prior and **likelihood** chosen for model

$$q(\theta | D) \propto p(\theta)\omega(D | \theta)$$

Posterior determined
by algebra

In general, distributions are
different functions

Conjugate Pairs

For some special models the posterior takes a simple form

$$p(\theta | D) \propto p(\theta) \omega(D | \theta)$$

**Prior and posterior are the same
distribution (with different parameters)**

We have already seen one example, the Beta-Bernoulli conjugate pair:

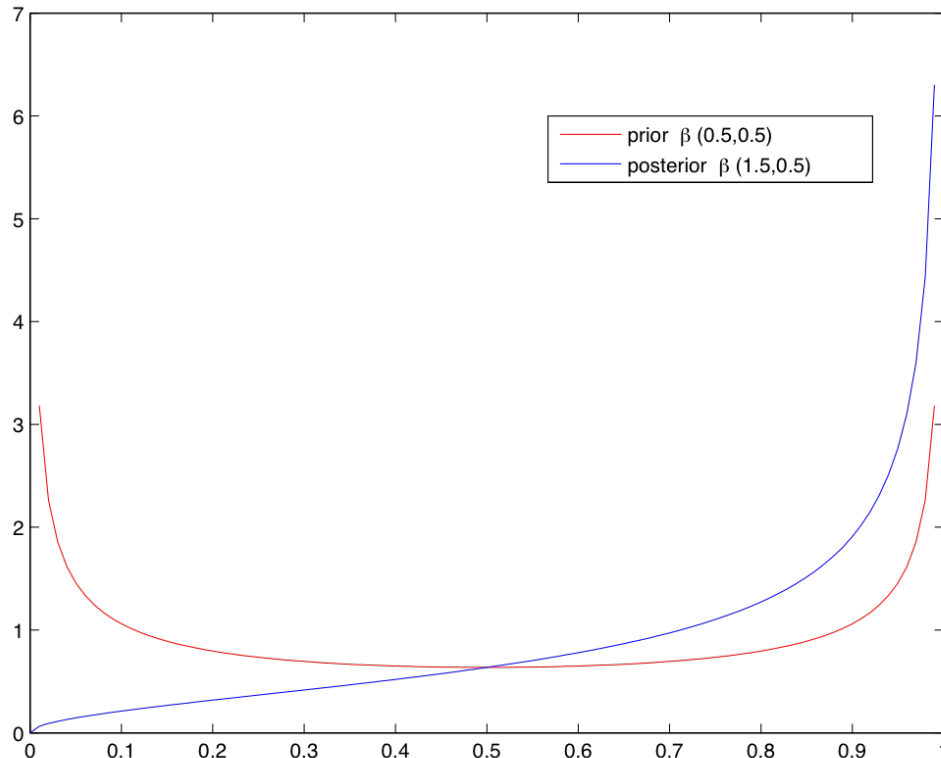
$$\text{Beta}(\theta | \alpha + \text{num.-heads}, \beta + \text{num.-tails}) \propto \text{Beta}(\theta | \alpha, \beta) \prod_i \text{Bernoulli}(x_i | \theta)$$

Same PDF

Example: Beta-Bernoulli

After a single coinflip of heads ($x=1$) the posterior is...

$$\text{Beta}(\theta \mid \alpha + x, \beta + 1 - x)$$



The prior (red) is a fair coin,

$$\text{Beta}(\theta \mid \alpha = 0.5, \beta = 0.5)$$

After observing one head, the posterior (blue) concentrates on heads,

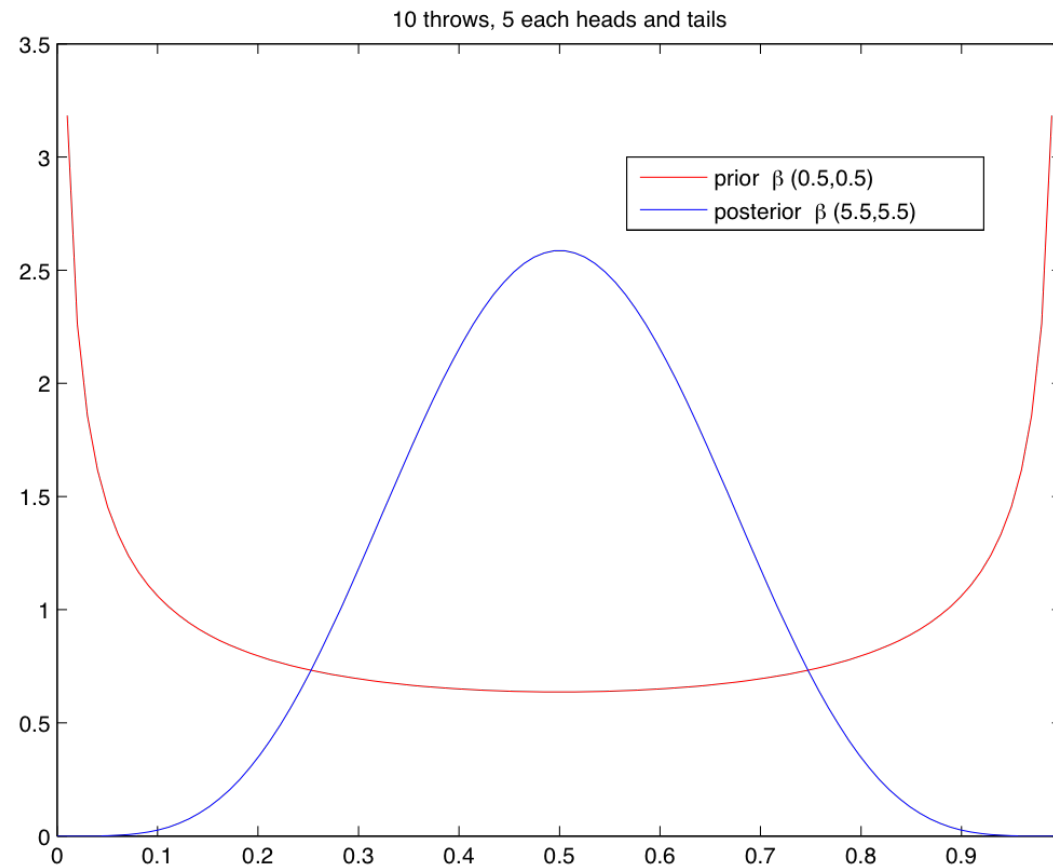
$$\text{Beta}(\theta \mid 1.5, 0.5)$$

What do you expect if we flip $N=10$ times with 5 heads and 5 tails?

Example: Beta-Bernoulli

After a $N=10$ flips (5 heads, 5 tails) we have...

$$\text{Beta}(\theta \mid \alpha + 5, \beta + 5) = \text{Beta}(\theta \mid 5.5, 5.5)$$



Posterior
concentrates on fair
coin $\theta = 0.5$

Example: Beta-Bernoulli

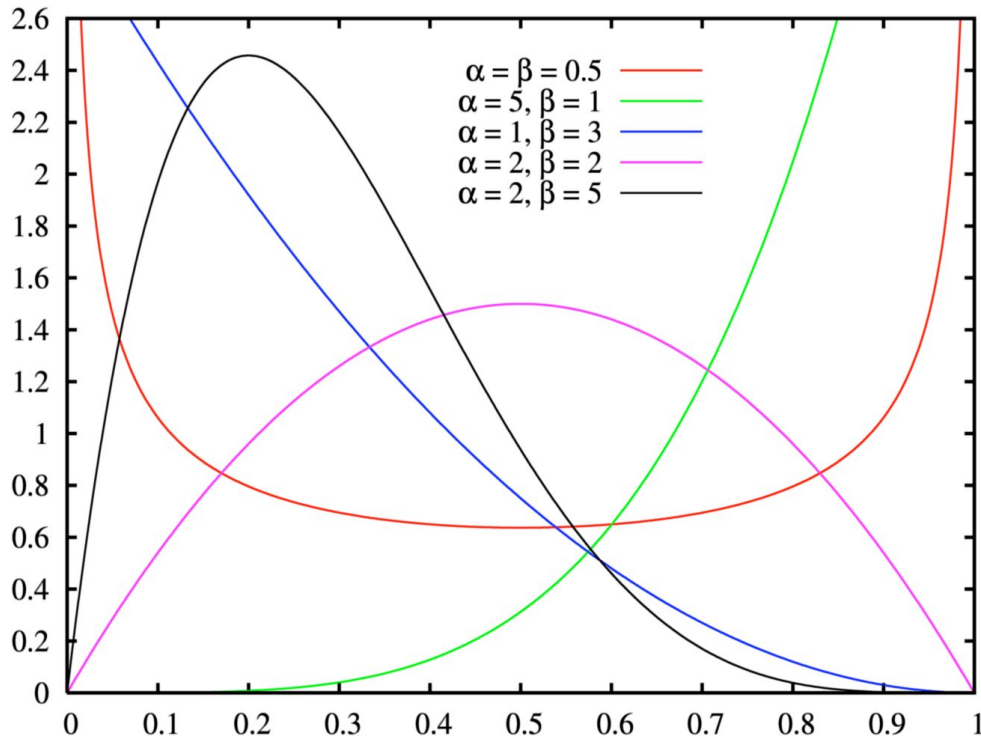
Bernoulli A.k.a. the **coinflip** distribution on binary RVs $X \in \{0, 1\}$

$$\text{Bernoulli}(X | \theta) = \theta^X (1 - \theta)^{(1-X)}$$



Beta distribution on $\theta \in (0, 1)$ with $\alpha, \beta > 0$ has PDF,

$$\text{Beta}(\theta | \alpha, \beta) \propto \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$



For N coinflips x_1, \dots, x_N the posterior is,

$$\text{Beta}(\theta | \alpha + \sum_i x_i, \beta + N - \sum_i x_i)$$

Example: Beta-Bernoulli

$$\begin{aligned} \text{Beta}(\theta \mid \alpha, \beta) \prod_{i=1}^N \text{Bernoulli}(x_i \mid \theta) &\propto \\ &\propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \prod_i \theta^{x_i} (1-\theta)^{1-x_i} \\ &= \theta^{\alpha-1} (1-\theta)^{\beta-1} \theta^{\sum_i x_i} (1-\theta)^{\sum_i (1-x_i)} \\ &= \theta^{\alpha-1} (1-\theta)^{\beta-1} \theta^{\sum_i x_i} (1-\theta)^{(N-\sum_i x_i)} \\ &= \theta^{\alpha-1+\sum_i x_i} (1-\theta)^{\beta-1+N-\sum_i x_i} \\ &\propto \text{Beta}(\theta \mid \alpha + \sum_i x_i, \beta + N - \sum_i x_i) \end{aligned}$$

Other Conjugate Pairs

Likelihood	Model Parameters	Conjugate Prior
Normal	Mean	Normal
Normal	Mean / Variance	Normal-Inv-Gamma
Multivariate Normal	Mean / Variance	Normal-Inv-Wishart
Multinomial	Probability vector	Dirichlet
Gamma	Rate	Gamma
Poisson	Rate	Gamma
Exponential	Rate	Gamma

Wikipedia has a nice list of standard conjugate forms...

https://en.wikipedia.org/wiki/Conjugate_prior

Priors in AI / ML / Data Science

- Priors are often used as *regularizers* (promote smoothing)
 - Reduces overfitting as random noise is not smooth
 - Often regularizers can be of simple form, even conjugate
- Priors often house sophisticated domain knowledge
 - Possibly from earlier encounters with data
 - Possibly problem constraints (e.g. θ must be nonnegative)
 - World knowledge is complex, so good priors are often complex and **not conjugate**

Choosing a Prior

- Conjugate priors can keep posteriors in closed form
 - This can speed up our codes (a lot!)
- The conjugate priors for standard distributions are fairly expressive
 - Often they can serve the purpose
- They are cool (better than doing nothing or the wrong thing)
- But they require that the likelihood is of a standard form
 - This is often a lot to hope for!
- Simply expressed functions may not be able to encode what you know
 - Constraints, non-local relationships

Prediction

Can make predictions of unobserved \tilde{y} before seeing any data,

$$p(\tilde{y}) = \sum_k p(\theta = k)p(\tilde{y} | \theta = k)$$

Similar calculation to marginal likelihood

*This is the **prior predictive** distribution*

For continuous parameters sum turns into integral,

$$p(\tilde{y}) = \int p(\theta)p(\tilde{y} | \theta) d\theta$$

*This is a prediction based on **no observed data***

Prediction

When we observe y we can predict future observations \tilde{y} ,

$$p(\tilde{y}) = \sum_k \underbrace{p(\theta = k | y)}_{\text{This is now the posterior}} p(\tilde{y} | \theta = k)$$

This is now the posterior

*This is the **posterior predictive distribution***

Again, for continuous parameters sum turns into integral,

$$p(\tilde{y} | y) = \int p(\theta | y) p(\tilde{y} | \theta) d\theta$$

Prediction Example

About **29%** of American adults have high blood pressure (BP). Home test has **30% false positive** rate and no false negative error.



What is the likelihood of *another* positive measurement?

$$p(\tilde{y} = true \mid y = true) = \sum_{\theta \in \{true, false\}} p(\theta \mid y = true) p(\tilde{y} = true \mid \theta)$$

$$= 0.42 * 0.30 + 0.58 * 1.00 \approx 0.71$$

What conclusions can be drawn from this calculation?

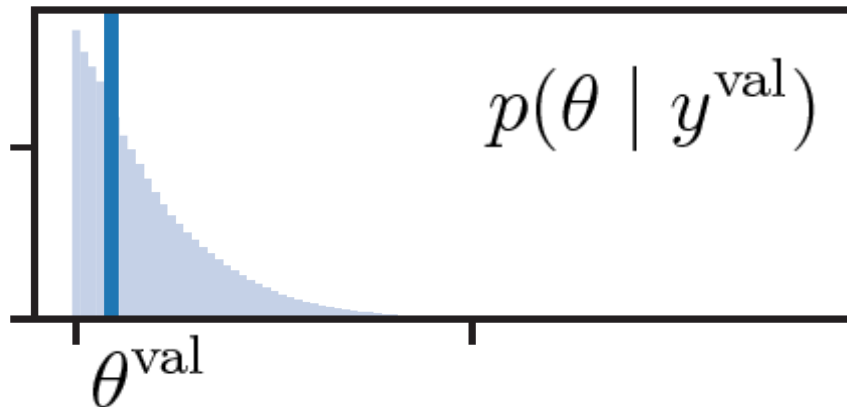
Model Validation

How do we know if the model $p(\theta, y)$ is good?

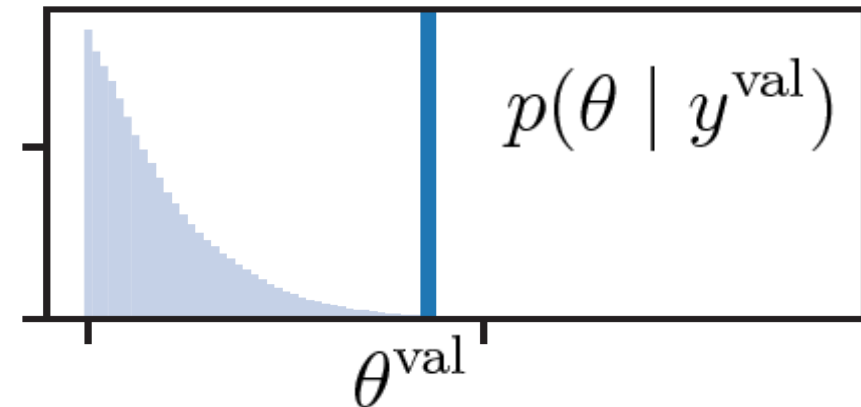
Supervised Learning

Validation set $\{(\theta^{\text{val}}, y^{\text{val}})\}$ consists of known θ^{val} . Are true values typically preferred under the posterior?

Good (maybe lucky)



Not Good (maybe unlucky)



Repeat trials over validation set for more certainty

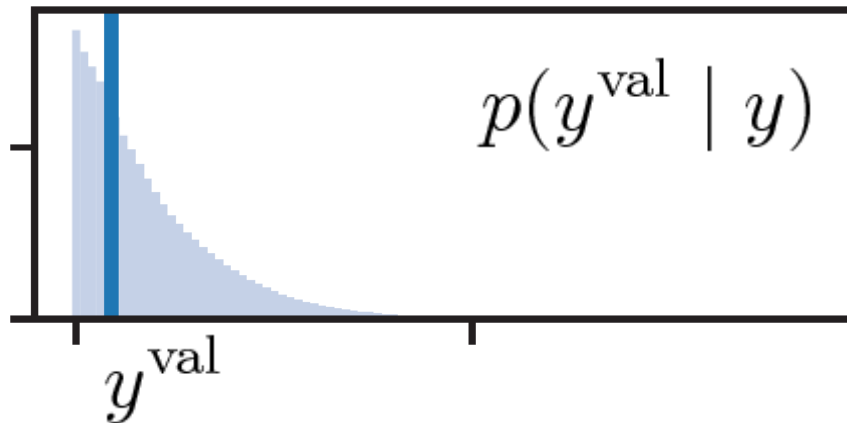
Model Validation

How do we know if the model $p(\theta, y)$ is good?

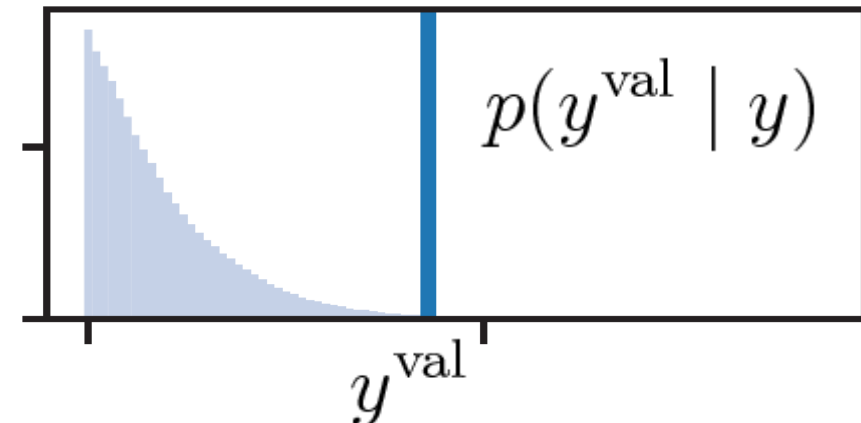
Unsupervised Learning

Validation set $\{y^{\text{val}}\}$ only contains observable data. Check validation data against posterior-predictive distribution.

Good (maybe lucky)



Not Good (maybe unlucky)



Repeat trials over validation set for more certainty

Likelihood and Odds Ratios

Which parameter value θ_1 or θ_2 is more likely to have generated the observed data y ?

The **posterior odds ratio** is:

$$\frac{p(\theta_1 | y)}{p(\theta_2 | y)} = \frac{p(\theta_1) p(y | \theta_1) \cancel{p(y)}}{p(\theta_2) p(y | \theta_2) \cancel{p(y)}}$$

Prior Odds
Ratio

Likelihood
Ratio

Observe: the marginal likelihood $p(y)$ cancels!

Posterior Summarization

Ideally we would report the full posterior distribution as the result of inference...but this is not always possible

Summary of Posterior Location:

Point estimates: mean (MMSE), mode, median (min. absolute error)

Summary of Posterior Uncertainty:

Credible intervals / regions, posterior entropy, variance

Bayesian analysis should report uncertainty when possible

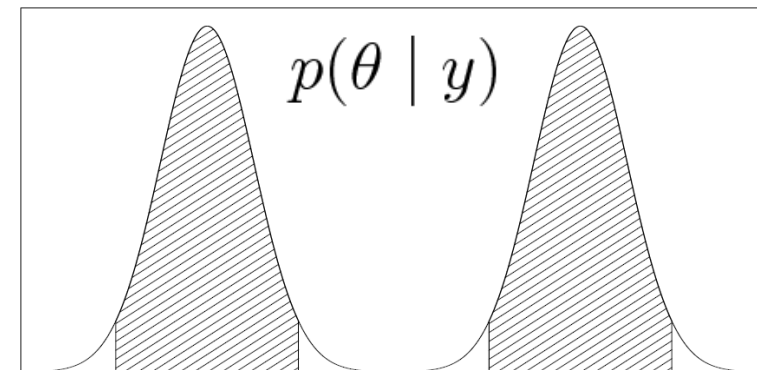
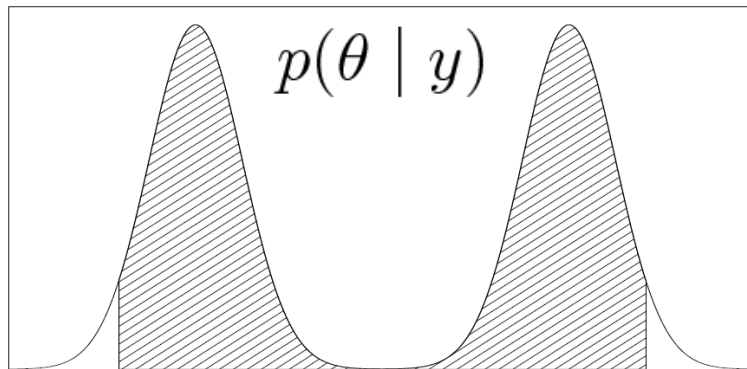
Credible Interval

Def. For parameter $0 < \alpha < 1$ the $100(1 - \alpha)\%$ credible interval $(L(y), U(y))$ satisfies,

$$p(L(y) < \theta < U(y) \mid y) = \int_{L(y)}^{U(y)} p(\theta \mid y) = 1 - \alpha$$

Interval containing fixed percentage of posterior probability density.

Note: This is not unique -- consider the 95% intervals below:



Frequentist Inference

Example: Suppose we observe the outcome of N coin flips.
 $y = \{y_1, \dots, y_N\}$. What is the probability of heads θ (coin bias)?

- Coin bias θ is not random (e.g. there is some *true* value)
- Uncertainty reported as confidence interval (typically 95%)

Correct Interpretation: On repeated trials of N coin flips θ will fall inside the confidence interval 95% of the time (in the limit)

- Inferences are valid for multiple trials, **never on single trials**

Wrong Interpretation: For *this trial* there is a 95% chance θ falls in the confidence interval

Bayesian Inference

Posterior distribution is complete representation of uncertainty

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

Prior Belief
Likelihood
Marginal Likelihood
(more on this later)

- Must specify a prior belief $p(\theta)$ about coin bias
- Coin bias θ is a random quantity
- Interval $p(l(y) < \theta < u(y) | y) = 0.95$ can be reported in lieu of full posterior, and takes intuitive interpretation for a single trial

Interval Interpretation: For this experiment there is a 95% chance that θ lies in the interval

Summary

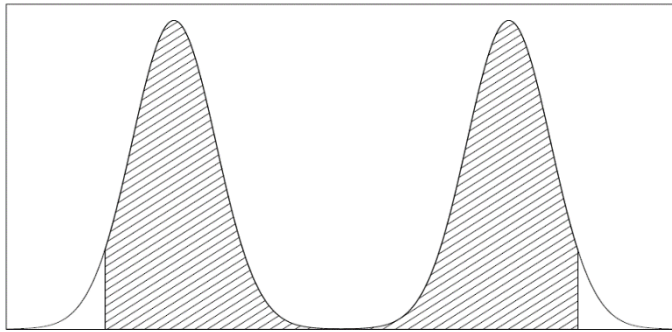
- Bayesian statistics interprets probability differently than classical stats
 - Frequentist: Probability \rightarrow Long run odds in repeated trials
 - Bayesian: Probability \rightarrow Belief of outcome that captures all uncertainty
- Bayesian models treat unknown parameter as random, with a prior
- Bayesian inference via the *posterior distribution* using Bayes' rule

$$p(\theta | y) = \frac{p(\theta)p(y | \theta)}{p(y)}$$

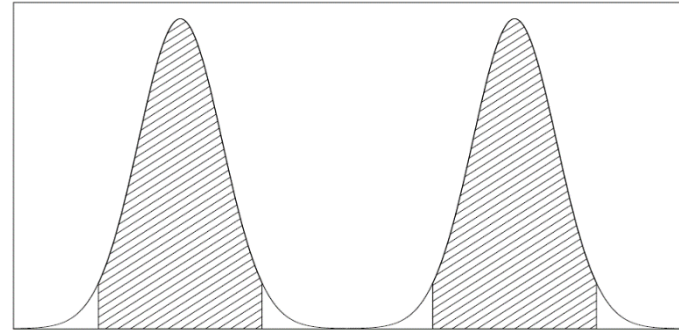
- Bayesian estimators minimize expected risk (e.g. MMSE)
- Maximum a posteriori (MAP) estimate maximizes posterior probability

Summary

- Conjugate prior-posterior pairs ensure closed-form posterior inference
- Posterior uncertainty can be characterized by credible intervals



Not necessarily
unique



- Selecting models can be done via posterior odds ratio

$$\frac{p(\theta_1 | y)}{p(\theta_2 | y)} = \frac{p(\theta_1) p(y | \theta_1) \cancel{p(y)}}{p(\theta_2) p(y | \theta_2) \cancel{p(y)}}$$

- Parameter can be marginalized out via prior/posterior predictive dist'n