

# STAT 544 - Bayesian Statistics

## Probability and Inference

Jarad Niemi (Dr. J)

Iowa State University

last updated: January 14, 2016

# Outline

- Quick review of probability
  - Kolmogorov's axioms
  - Bayes' rule
  - Application to Down's syndrome screening
- Bayesian statistics
  - Condition on what is known
  - Describe uncertainty using probability
  - Exponential example
- What is probability?
  - Frequency interpretation
  - Personal belief
- Why or why not Bayesian?

# Events

## Definition

The set,  $\Omega$ , of all possible outcomes of a particular experiment is called the **sample space** for the experiment.

## Definition

An **event** is any collection of possible outcomes of an experiment, that is, any subset of  $\Omega$  (including  $\Omega$  itself).

# Craps

Craps:

- $\Omega = \{(1, 1), (1, 2), \dots, (1, 6), (2, 1), (2, 2), \dots, (6, 6)\}$
- Come-out roll win: the sum of the dice is 7 or 11
- Come-out roll loss: the sum of the dice is 2, 3, or 12
- Come-out roll establishes a point: the sum of the dice is 4, 5, 6, 8, 9, or 10
- Events:
  - the come-out roll wins
  - the come-out roll loses
  - the come-out roll establishes a point

# Pairwise disjoint

## Definition

Two events  $A_1$  and  $A_2$  are **disjoint** (or **mutually exclusive**) if both  $A_1$  and  $A_2$  cannot occur simultaneously, i.e.  $A_i \cap A_j = \emptyset$ . The events  $A_1, A_2, \dots$  are **pairwise disjoint** (or **mutually exclusive**) if  $A_i$  and  $A_j$  cannot occur simultaneously for all  $i \neq j$ , i.e.  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ .

Craps pairwise disjoint examples:

- Win ( $A_1$ ), Loss ( $A_2$ )
- Win ( $A_1$ ), Loss ( $A_2$ ), Point ( $A_3$ )
- $A_1 = (1, 1), A_2 = (1, 2), \dots, A_6 = (1, 6), A_7 = (2, 1), \dots, A_{12} = (2, 6), \dots, A_{36} = (6, 6)$

# Partition

## Definition

A set of events,  $\{A_1, A_2, \dots\}$ , is a **partition** of the sample space  $\Omega$  if and only if

- the events in  $\{A_1, A_2, \dots\}$  are pairwise disjoint and
- $\bigcup_{i=1}^{\infty} A_i = \Omega$ .

Craps partition examples:

- Win ( $A_1$ ), Loss ( $A_2$ ), Point ( $A_3$ )
- $A_1 = (1, 1), A_2 = (1, 2), \dots, A_6 = (1, 6), A_7 = (2, 1), \dots, A_{12} = (2, 6), \dots, A_{36} = (6, 6)$

# Kolmogorov's axioms of probability

## Definition

Given a sample space  $\Omega$  and event space  $F$ , a **probability** is a function  $P : F \rightarrow \mathbb{R}$  that satisfies

1.  $P(A) \geq 0$  for any  $A \in F$
2.  $P(\Omega) = 1$
3. If  $A_1, A_2, \dots \in F$  are pairwise disjoint, then
$$P(A_1 \text{ or } A_2 \text{ or } \dots) = \sum_{i=1}^{\infty} P(A_i).$$

# Craps come-out roll probabilities

The following table provides the probability mass function for the sum of the two dice (assuming the probability of each elementary outcome is equal):

Outcome	2	3	4	5	6	7	8	9	10	11	12	Sum
Combinations	1	2	3	4	5	6	5	4	3	2	1	36
Probability	$\frac{1}{36}$	$\frac{2}{36}$	$\frac{3}{36}$	$\frac{4}{36}$	$\frac{5}{36}$	$\frac{6}{36}$	$\frac{5}{36}$	$\frac{4}{36}$	$\frac{3}{36}$	$\frac{2}{36}$	$\frac{1}{36}$	1

Craps probability examples:

- $P(\text{Win}) = P(7 \text{ or } 11) = 8/36 = 2/9$
- $P(\text{Loss}) = P(2, 3, \text{ or } 12) = 4/36 = 1/9$
- $P(\text{Point}) = P(4, 5, 6, 8, 9 \text{ or } 10) = 6/9$



# Conditional probability

## Definition

If  $A$  and  $B$  are events in  $F$ , and  $P(B) > 0$ , then the **conditional probability of  $A$  given  $B$** , written  $P(A|B)$ , is

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

Craps conditional probability example:

$$P(7|\text{Win}) = \frac{P(7 \text{ and Win})}{P(\text{Win})} = \frac{P(7)}{P(\text{Win})} = \frac{6/36}{8/36} = \frac{6}{8}$$

# Bayes' rule

## Theorem

If  $A$  and  $B$  are events in  $F$ , then *Bayes' rule* states

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

Craps Bayes' rule example:

$$P(7|\text{Win}) = \frac{P(\text{Win}|7)P(7)}{P(\text{Win})} = \frac{1 \cdot P(7)}{P(\text{Win})} = \frac{6/36}{8/36} = \frac{6}{8}$$

# Down Syndrome screening

If a pregnant woman has a test for Down syndrome and it is positive, what is the probability that the child will have Down syndrome? Let  $D$  indicate a child with Down syndrome and  $D^c$  the opposite. Let '+' indicate a positive test result and  $-$  a negative result.

$$\text{sensitivity} = P(+|D) = 0.94$$

$$\text{specificity} = P(-|D^c) = 0.77$$

$$\text{prevalence} = P(D) = 1/1000$$

$$\begin{aligned} P(D|+) &= \frac{P(+|D)P(D)}{P(+)} = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|D^c)P(D^c)} = \frac{0.94 \cdot 0.001}{0.94 \cdot 0.001 + 0.23 \cdot 0.999} \\ &\approx 1/250 \end{aligned}$$

$$P(D|-) \approx 1/10,000$$

# A Bayesian statistics

Let

- $y$  be the data we will collect from an experiment,
- $K$  be everything we know for certain about the world (aside from  $y$ ), and
- $\theta$  be anything we don't know for certain.

My definition of a Bayesian statistician is an individual who makes decisions based on the probability distribution of those things we don't know conditional on what we know, i.e.

$$p(\theta|y, K).$$

# Bayesian statistics (with explicit conditioning)

- Parameter estimation:

$$p(\theta|y, M)$$

where  $M$  is a model with parameter (vector)  $\theta$  and  $y$  is data assumed to come from model  $M$  with true parameter  $\theta_0$ .

- Hypothesis testing/model selection:

$$p(M_j|y, \mathcal{M})$$

where  $\mathcal{M}$  is a set of models with  $M_j \in \mathcal{M}$  for  $j = 1, 2, \dots$  and  $y$  is data assumed to come from some model  $M_0 \in \mathcal{M}$ .

- Prediction:

$$p(\tilde{y}|y, M)$$

where  $\tilde{y}$  is unobserved data and  $y$  and  $\tilde{y}$  are both assumed to come from  $M$ . Alternatively,

$$p(\tilde{y}|y, \mathcal{M})$$

where  $y$  and  $\tilde{y}$  are both assumed to come from some  $M_0 \in \mathcal{M}$ .

# Bayesian statistics (with implicit conditioning)

- Parameter estimation:

$$p(\theta|y)$$

where  $\theta$  is the unknown parameter (vector) and  $y$  is the data.

- Hypothesis testing/model selection:

$$p(M_j|y)$$

where  $M_j$  is one of a set of models under consideration and  $y$  is data assumed to come from one of those models.

- Prediction:

$$p(\tilde{y}|y)$$

where  $\tilde{y}$  is unobserved data and  $y$  and  $\tilde{y}$  are both assumed to come from the same (set of) model(s).

# Bayes' Rule

Bayes' Rule applied to a partition  $P = \{A_1, A_2, \dots\}$ ,

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{P(B)} \frac{P(B|A_i)P(A_i)}{\sum_{i=1}^{\infty} P(B|A_i)P(A_i)}$$

Bayes' Rule also applies to probability density (or mass) functions, e.g.

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

where the integral plays the role of the sum in the previous statement.

# Parameter estimation

Let  $y$  be data from some model with unknown parameter  $\theta$ . Then

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

and we use the following terminology

Terminology	Notation
Posterior	$p(\theta y)$
Prior	$p(\theta)$
Model	$p(y \theta)$
Prior predictive distribution (marginal likelihood)	$p(y)$

If  $\theta$  is discrete (continuous),

then  $p(\theta)$  and  $p(\theta|y)$  are probability mass (density) functions.

If  $y$  is discrete (continuous),

then  $p(y|\theta)$  and  $p(y)$  are probability mass (density) functions.



## Example: exponential model

Let  $Y|\theta \sim \text{Exp}(\theta)$ , then this defines the likelihood, i.e.

$$p(y|\theta) = \theta e^{-\theta y}.$$

Let's assume a convenient prior  $\theta \sim \text{Ga}(a, b)$ , then

$$p(\theta) = \frac{b^a}{\Gamma(a)} \theta^{a-1} e^{-b\theta}.$$

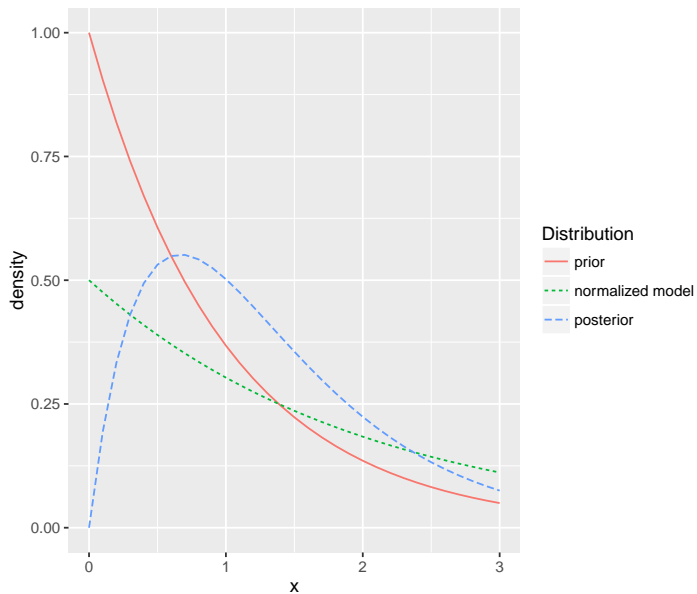
The prior predictive distribution is

$$p(y) = \int p(y|\theta)p(\theta)d\theta = \frac{b^a}{\Gamma(a)} \frac{\Gamma(a+1)}{(b+y)^{a+1}}.$$

The posterior is

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)} = \frac{(b+y)^{a+1}}{\Gamma(a+1)} \theta^{a+1-1} e^{-(b+y)\theta},$$

thus  $\theta|y \sim \text{Ga}(a+1, b+y)$ .



## A shortcut

If

$$p(y) = \int p(y|\theta)p(\theta)d\theta < \infty,$$

then we can actually use the following to find the posterior

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

where the  $\propto$  signifies that terms not involving  $\theta$  (or anything on the left of the conditioning bar) are irrelevant and can be dropped.

In the exponential example

$$p(\theta|y) \propto p(y|\theta)p(\theta) = \theta e^{-\theta y} \theta^{a-1} e^{-b\theta} = \theta^{a+1-1} e^{-(b+y)\theta}$$

where we can recognize  $p(\theta|y)$  as the **kernel** of a  $Ga(a+1, b+y)$  distribution and thus  $\theta|y \sim Ga(a+1, b+y)$  and  $p(y) < \infty$ .

# Independent data

Suppose  $Y_i|\theta \stackrel{ind}{\sim} \text{Exp}(\theta)$  for  $i = 1, \dots, n$  and  $y = (y_1, \dots, y_n)$ , then

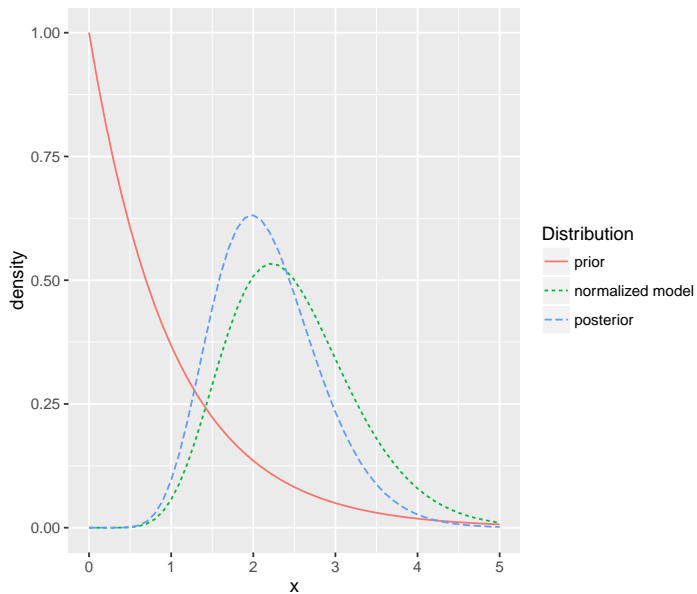
$$p(y|\theta) = \prod_{i=1}^n p(y_i|\theta) = \theta^n e^{-\theta n\bar{y}}$$

Then

$$p(\theta|y) \propto p(y|\theta)p(\theta) \propto \theta^{a+n-1} e^{-(b+n\bar{y})\theta}$$

where  $n\bar{y} = \sum_{i=1}^n y_i$ . We recognize this as the kernel of a gamma, i.e.

$$\theta|y \sim \text{Ga}(a+n, b+n\bar{y}).$$



# Bayesian learning (in parameter estimation)

So, Bayes' rule provides a formula for updating from prior beliefs to our posterior beliefs based on the data we observe, i.e.

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

Suppose we gather  $y_1, \dots, y_n$  sequentially (and we assume  $y_i$  independent conditional on  $\theta$ ), then we have

$$p(\theta|y_1) \propto p(y_1|\theta)p(\theta)$$

and

$$p(\theta|y_1, \dots, y_i) \propto p(y_i|\theta)p(\theta|y_1, \dots, y_{i-1})$$

So Bayesian learning is

$$p(\theta) \rightarrow p(\theta|y_1) \rightarrow p(\theta|y_1, y_2) \rightarrow \dots \rightarrow p(\theta|y_1, \dots, y_n).$$

# Model selection

Formally, to select a model (or average over models), we use

$$p(M_j|y) \propto p(y|M_j)p(M_j)$$

where

- $p(y|M_j)$  is the likelihood of the data when model  $M_j$  is true
- $p(M_j)$  is the prior probability for model  $M_j$
- $p(M_j|y)$  is the posterior probability for model  $M_j$

Thus, a Bayesian approach provides a natural way to learn about models, i.e.  $p(M_j) \rightarrow p(M_j|y)$ .

# Prediction

Let  $y$  be observed data and  $\tilde{y}$  be unobserved data from a model with parameter  $\theta$ , then

$$\begin{aligned} p(\tilde{y}|y) &= \int p(\tilde{y}, \theta|y) d\theta \\ &= \int p(\tilde{y}|\theta, y) p(\theta|y) d\theta \\ &= \int p(\tilde{y}|\theta) p(\theta|y) d\theta \end{aligned}$$

where  $p(\theta|y)$  is the posterior we obtained using Bayesian parameter estimation techniques.



## Example: exponential distribution

From previous, let  $y_i \stackrel{\text{ind}}{\sim} \text{Exp}(\theta)$  and  $\theta \sim \text{Ga}(a, b)$ , then  $\theta|y \sim \text{Ga}(a + n, b + n\bar{y})$ . Suppose we are interested in predicting a new value  $\tilde{y} \sim \text{Exp}(\theta)$  (conditionally independent of  $y = (y_1, \dots, y_n)$  given  $\theta$ ). Then we have

$$\begin{aligned}
 p(\tilde{y}|y) &= \int p(\tilde{y}|\theta)p(\theta|y)d\theta \\
 &= \int \theta e^{-\theta\tilde{y}} \frac{(b+n\bar{y})^{a+n}}{\Gamma(a+1)} \theta^{a+n} e^{-\theta(b+n\bar{y})} d\theta \\
 &= \frac{(b+n\bar{y})^{a+n}}{\Gamma(a+n)} \int \theta^{a+n+1} e^{-\theta(b+n\bar{y}+\tilde{y})} d\theta \\
 &= \frac{(b+n\bar{y})^{a+n}}{\Gamma(a+n)} \frac{\Gamma(a+n+1)}{(b+n\bar{y}+\tilde{y})^{a+n+1}} \\
 &= \frac{(a+n)(b+n\bar{y})^{a+n}}{(\tilde{y}+b+n\bar{y})^{a+n+1}}
 \end{aligned}$$

This is the Lomax distribution for  $\tilde{y}$  with parameters  $a + n$  and  $b + n\bar{y}$ .

# What is probability?

Consider the following three typical uses of the word “probability”:

- What is the probability I will win on the come-out roll in craps?
- What is the probability my unborn child has Down’s syndrome given that they tested positive in an initial screening?
- What is the probability the Green Bay Packers will win this year’s superbowl?

## Win on the come-out roll in craps

To win on the come-out roll in craps requires that the sum of two fair six-sided die is either a 7 or an 11. We calculated this probability earlier (based on equal probabilities of all simple outcomes) to be  $2/9$ . We likely meant that if we were to repeatedly roll the die, the long term proportion of wins (7s and 11s) would be  $2/9$ , i.e.

$$\text{if } X_i = \begin{cases} 1 & \text{if win on roll } i \\ 0 & \text{otherwise} \end{cases} \quad \text{then} \quad \lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n X_i}{n} \rightarrow \frac{2}{9}.$$

### Definition

The **frequency** interpretation of probability is based on the relative frequency of an event (assumed to be performed in an identical manner).

# Win on the come-out roll in craps

## Definition

The **frequency** interpretation of probability is based on the relative frequency of an event (assumed to be performed in an identical manner).

Two problems with this frequency interpretation:

- You cannot possibly throw the dice in an identical manner.
- If I knew enough physics, I could model each throw and tell you exactly what the result would be, i.e. the only randomness is because they throws are not identical.

# Down's syndrome

What is the probability my unborn child has Down's syndrome given that they tested positive in an initial screening?

Here the frequency interpretation makes no sense for two reasons:

- There is only one child and thus no repeat of the experiment.
- There is no randomness: either the child has Down's syndrome or does not.

Instead, we only have our own uncertainty about whether the child has Down's syndrome or not.

## Down's syndrome

Also, why are we only conditioning on the positive test result, shouldn't we condition on everything else that is important, e.g. age. Then the probability we care about is

$$P(D|+, \text{mother is 33}) = \frac{P(+|D, \text{mother is 33})P(D|\text{mother is 33})}{P(+|\text{mother is 33})}$$

Now the specificity, sensitivity, and prevalence are all the relative frequency of the event for this subpopulation.

But what about other measured variables, e.g. Caucasian, leaves in MN, of Scandanavian descent, etc. Taken to its logical extreme, each probability becomes a statement about one single event, e.g. for this individual.

# Superbowl Champions

What is the probability the Green Bay Packers win the Superbowl?

By similar arguments:

- There is only one Superbowl this year and only one Green Bay Packers.
- Is the world random? i.e. do we have free will? If not, then (with enough time, computing power, money, etc) we could model the world and know what the result will be. If yes, is there an objective probability that we could be estimating?

# Personal belief

## Definition

A **subjective probability** describes an individual's personal judgement about how likely a particular event is to occur.

<http://www.stats.gla.ac.uk/glossary/?q=node/488>

**Remark** Coherence of bets. The probability  $p$  you assign to an event  $E$  is the fraction at which you would exchange  $p$  for a return of 1 if  $E$  occurs.

Rational individuals can differ about the probability of an event by having different knowledge, i.e.  $P(E|K_1) \neq P(E|K_2)$ . But given enough data, we might have  $P(E|K_1, y) \approx P(E|K_2, y)$ .



# Personal belief

Using a personal belief definition of probability, it is easy to reconcile the use of probability in common language:

- What is the probability I will win on the come-out roll in craps?
- What is the probability my unborn child has Down's syndrome given that they tested positive in an initial screening?
- What is the probability the Green Bay Packers will win this year's superbowl?
- What is the probability that global climate change is primarily driven by human activity?
- What is the probability the Higgs Boson exists?

and in the mathematical notation:

- $p(\theta) \rightarrow p(\theta|y)$
- $p(H_1) \rightarrow p(H_1|y)$
- $p(\tilde{y}|y)$

# Why or why not Bayesian?

Why do a Bayesian analysis?

- Incorporate prior knowledge via  $p(\theta)$
- Coherent, i.e. everything follows from specifying  $p(\theta|y)$
- Interpretability of results, e.g. the probability the parameter is in  $(L, U)$  is 95%

Why not do a Bayesian analysis?

- Need to specify  $p(\theta)$
- Computational cost
- Does not guarantee coverage, i.e. how well do the procedures work over all their uses