

Probability, 1

Frank Edwards

10/23/2019

How often, on average, does an event occur?

Probability is a set of tools for describing randomness.

Probability helps us sort signal (patterns) from noise.

Two core theories

Frequentist: Probability is the proportion of times an event occurs if we repeat an experiment under the same conditions many times

If n_x is the number of heads, and n_t is the number of coin flips, then the probability of heads is

$$P(x) \approx \frac{n_x}{n_t}$$

$$P(x) = \lim_{n_t \rightarrow \infty} \frac{n_x}{n_t}$$

Bayesian: Probability is a subjective judgment about the likelihood that an event occurs, with endpoints at 0 (never occurs) and 1 (always occurs)

Deterministic processes do not include randomness. For example, if I drop a ball, it will fall.

Stochastic events include randomness. For example, if I flip a fair coin, it will be heads half of the time.

Nearly all social processes are *stochastic*.

- Experiment: an action that produces stochastic events
- Sample space (Ω): a set of all possible outcomes of the experiment
- Event: a subset of the sample space

Example: coin flips

- Experiment
 1. flip a coin
 2. roll a dice
 3. vote in democratic primary
- Sample space Ω
 1. {Heads, Tails}
 2. {1,2,3,4,5,6}
 3. {abstain, vote for Harris, vote for Warren, vote for Sanders, vote for Castro}
- Event
 1. Heads, tails, not heads, heads or tails, heads and tails
 2. 3, even number, anything but 6
 3. Did not vote, voted for a woman, voted for a senator

If all outcomes are equally likely, and n represents the number of elements in a given set, then probability P of event A is:

$$P(A) = \frac{n_A}{n_\Omega}$$

1. The probability of any event A is non-negative: $P(A) \geq 0$
2. The probability that one of the outcomes in the sample space occurs is
1: $P(\Omega) = 1$
3. Addition rule: If events A and B are mutually exclusive:

$$P(A \text{ or } B) = P(A) + P(B)$$

Probability that an event doesn't occur

$$1 - P(\text{not } A) = P(A)$$

If $P(\text{rolling a } 6) = \frac{1}{6}$ then $P(\text{not rolling a } 6) = \frac{5}{6}$

$$P(A) = P(A \text{ and } B) + P(A \text{ and not } B)$$

If $P(\text{eats pizza}) = 0.5$ and $P(\text{eats pizza, happy}) = 0.4$, then
 $P(\text{eats pizza, unhappy}) = 0.1$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

If $P(\text{happy}) = 0.5$, then $P(\text{eats pizza or happy}) = 0.6$

How many ways can we arrange the set: $\{X,Y,Z\}$?

Permutations, combinations

How many ways can we arrange the set: $\{X,Y,Z\}$?

XYZ, XZY, YXZ, YZX, ZXY, ZYX

Permutations, combinations

How many ways can we arrange the set: {X,Y,Z}?

XYZ, XZY, YXZ, YZX, ZXY, ZYX

Sampling without replacement: how many permutations of n elements are there when we draw k at a time, and can't draw the same element twice

$${}_nP_k = \frac{n!}{(n-k)!}$$

Combinations are the number of selections without regard to their order

$${}_nC_k = \frac{{}_nP_k}{k!}$$

The birthday problem

How many people do we need to have in a class for there to be a 50 percent chance that at least two of them have the same birthday?

Assume each date of birth is equally likely.

$$P(\text{two people have the same birthday}) = 1 - P(\text{nobody has the same birthday})$$

Non-unique birthdays as permutations

We have k students

How many ways can k birthdays be arranged if there are no duplicate birthdays? A sample without replacement tells us:

$${}_nP_k = \frac{n!}{(n-k)!}$$

$${}_{365}P_k = \frac{365!}{(365-k)!}$$

Sampling with replacement, k non-unique birthdays can be arranged:

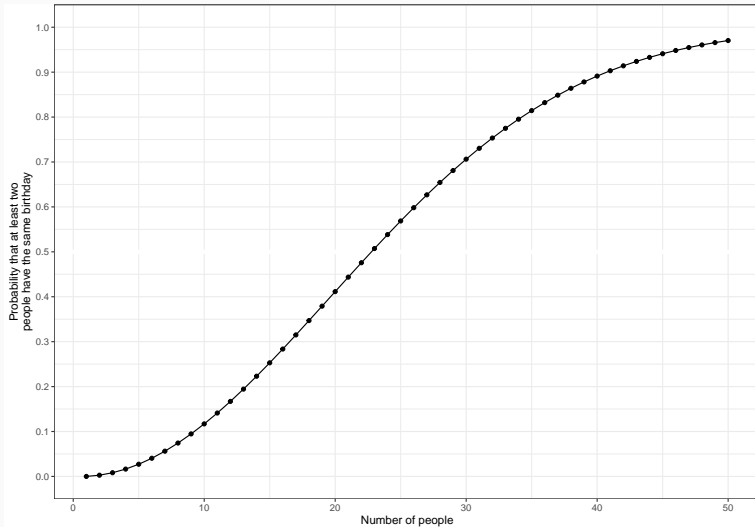
$$365^k$$

As a probability problem

Because $P(A) = \frac{n_A}{n_\Omega}$, we can obtain $1 - P(\text{nobody has the same birthday})$ as the size of the set for non-unique combinations divided by the size of the sample space

$${}_{365}P_k = 1 - \frac{365!}{(365 - k)!} \times \frac{1}{365^k}$$

The results



Simulation

Let's randomly draw k birthdays *with replacement* to estimate how likely a shared birthday is for various class sizes.

```
sim_bdays <- function(k) {  
  ## draw k random birthdays from the vector 1:365  
  days <- sample(1:365, k, replace = TRUE)  
  ## if there are no duplicates, there are k unique birthdays, return TRUE if  
  ## duplicates  
  length(unique(days)) < k  
}
```

```
sim_bdays(1)
```

```
## [1] FALSE
```

```
sim_bdays(366)
```

```
## [1] TRUE
```

Monte Carlo simulation

Repeat our random draw of k birthdays a large number of times to approximate the solution. How likely is a shared birthday for 30 students?

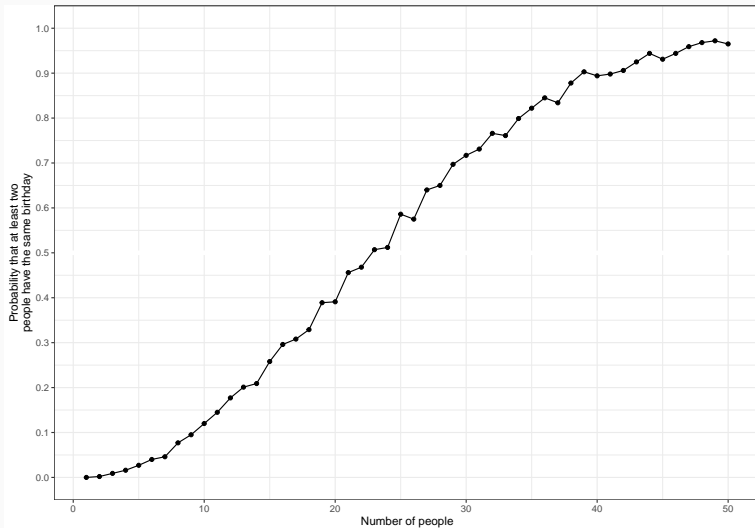
```
n <- 1000
k <- 30
sames <- rep(0, n)
for (i in 1:n) {
  sames[i] <- sim_bdays(k)
}
### monte carlo solution
mean(sames)
```

```
## [1] 0.709
```

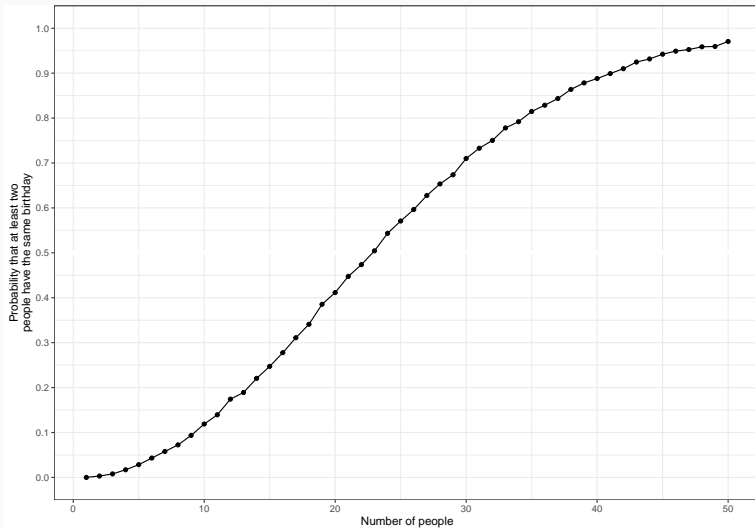
```
### exact solution
birthday(30)
```

```
## [1] 0.7063162
```

Monte Carlo to approximate a curve, 1000 simulations per k



Monte Carlo to approximate a curve, 10000 simulations per k



The joint probability of two events (A and B) occurring is expressed as

$$P(A \text{ and } B)$$

The marginal probability of an event B is

$$P(B)$$

The conditional probability of event A occurring given that event B occurred is the ratio of the joint probability of A and B divided by the marginal probability of B

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

Voter files

```
data("FLVoters")  
voters <- na.omit(FLVoters)  
head(voters)
```

##	surname	county	VTD	age	gender	race
## 1	PIEDRA	115	66	58	f	white
## 2	LYNCH	115	13	51	m	white
## 4	LATHROP	115	80	54	m	white
## 5	HUMMEL	115	8	77	f	white
## 6	CHRISTISON	115	55	49	m	white
## 7	HOMAN	115	84	77	f	white

Marginal probability

What is the probability that a random voter in the population is Black:

$$P(\text{Black}) = ?$$

```
voters %>% count(race) %>% mutate(n = n/sum(n))
```

```
## # A tibble: 6 x 2
##   race      n
##   <chr>    <dbl>
## 1 asian    0.0192
## 2 black    0.131
## 3 hispanic 0.131
## 4 native   0.00318
## 5 other    0.0340
## 6 white    0.682
```

Is a woman: $P(\text{Woman}) = ?$

```
voters %>% count(gender) %>% mutate(n = n/sum(n))
```

```
## # A tibble: 2 x 2
##   gender    n
##   <chr>    <dbl>
## 1 f        0.536
## 2 m        0.464
```

What is the probability that a voter is a Black woman:

$P(\text{Black and woman}) = ?$

```
voters %>% count(gender, race) %>% mutate(n = n/sum(n)) %>% pivot_wider(names_from = gender,  
  values_from = n)
```

```
## # A tibble: 6 x 3  
##   race      f      m  
##   <chr>    <dbl> <dbl>  
## 1 asian    0.00911 0.0101  
## 2 black    0.0744  0.0566  
## 3 hispanic 0.0731  0.0577  
## 4 native   0.00187 0.00132  
## 5 other    0.0173  0.0167  
## 6 white    0.360   0.322
```

What is the probability that a voter is a woman?

```
## # A tibble: 6 x 3
##   race      f      m
##   <chr>    <dbl> <dbl>
## 1 asian    0.00911 0.0101
## 2 black    0.0744  0.0566
## 3 hispanic 0.0731  0.0577
## 4 native   0.00187 0.00132
## 5 other    0.0173  0.0167
## 6 white    0.360   0.322
```

Use the law of total probability:

$$P(A) = P(A \text{ and } B) + P(A \text{ and not } B)$$

put differently, for all categories of B i:

$$P(A) = \sum_{i=1}^n P(A \text{ and } B_i)$$

Conditional probability

If a voter is a man, what is the probability that he is Asian: $P(\text{Asian}|\text{man}) = ?$

```
voters %>% filter(gender == "m") %>% count(race) %>% mutate(n = n/sum(n))
```

```
## # A tibble: 6 x 2
##   race      n
##   <chr>    <dbl>
## 1 asian    0.0217
## 2 black    0.122
## 3 hispanic 0.124
## 4 native   0.00284
## 5 other    0.0359
## 6 white    0.693
```

Conditional probability

Alternatively, we can use the definition of conditional probability as the ratio of the joint probability to the marginal probability:

$$P(\text{Asian}|\text{man}) = \frac{P(\text{Asian and man})}{P(\text{man})}$$

```
voters %>% count(gender, race) %>% mutate(n = n/sum(n)) %>% pivot_wider(names_from = gender,  
  values_from = n)
```

```
## # A tibble: 6 x 3  
##   race      f      m  
##   <chr>    <dbl> <dbl>  
## 1 asian    0.00911 0.0101  
## 2 black    0.0744  0.0566  
## 3 hispanic 0.0731  0.0577  
## 4 native   0.00187 0.00132  
## 5 other    0.0173  0.0167  
## 6 white    0.360   0.322
```

Conditioning on more than one variable

What is the probability that a male voter over age 60 is white?

$$P(\text{white} | \text{male and over 60})$$

```
voters %>% mutate(over60 = age > 60) %>% count(over60, gender, race) %>% mutate(n = n/sum(n)) %>%  
  pivot_wider(names_from = gender, values_from = n)
```

```
## # A tibble: 12 x 4  
##   over60 race      f      m  
##   <lgl> <chr>    <dbl>   <dbl>  
## 1 FALSE asian    0.00691 0.00823  
## 2 FALSE black    0.0555  0.0435  
## 3 FALSE hispanic 0.0549  0.0436  
## 4 FALSE native  0.00121 0.000768  
## 5 FALSE other   0.0124  0.0129  
## 6 FALSE white   0.212   0.198  
## 7 TRUE  asian    0.00219 0.00187  
## 8 TRUE  black    0.0189  0.0132  
## 9 TRUE  hispanic 0.0182  0.0142  
## 10 TRUE native  0.000658 0.000549  
## 11 TRUE other   0.00494 0.00373  
## 12 TRUE white   0.148   0.124
```


Conditioning on more than one variable

In general:

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

and

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

Independence

Two events are independent if knowledge of one event gives us no information about the other event.

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

$$A \perp B$$

if and only if

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

Are race and gender independent in voting population in Florida?

If independent, then we should observe

$P(\text{Black and male}) = P(\text{Black})P(\text{male})$ and so on for other groups.

```
ind_test <- voters %>% count(gender, race) %>% mutate(n = n/sum(n)) %>% pivot_wider(names_from = gender,  
  values_from = n)
```

```
ind_test
```

```
## # A tibble: 6 x 3  
##   race      f      m  
##   <chr>    <dbl> <dbl>  
## 1 asian    0.00911 0.0101  
## 2 black    0.0744  0.0566  
## 3 hispanic 0.0731  0.0577  
## 4 native   0.00187 0.00132  
## 5 other    0.0173  0.0167  
## 6 white    0.360   0.322
```

Are race and gender independent in voting population in Florida?

First, calculate the marginal probability of being male by summing over all joint probabilities for male by race

$$P(A) = \sum_{i=1}^n P(A \text{ and } B_i)$$

```
p_male <- sum(ind_test$m)
p_male
```

```
## [1] 0.4641721
```

Are race and gender independent in voting population in Florida?

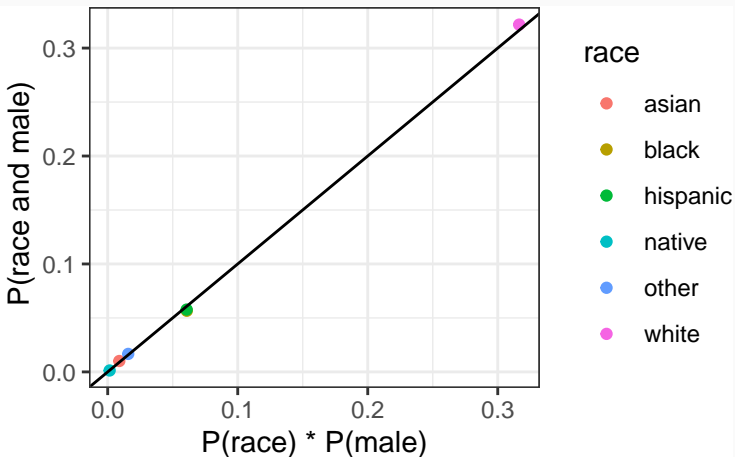
Next, calculate the marginal probability for each racial group by summing over joint probabilities of sex by race

```
p_race <- ind_test %>% mutate(p_race = m + f)
p_race
```

```
## # A tibble: 6 x 4
##   race      f      m p_race
##   <chr>    <dbl> <dbl> <dbl>
## 1 asian  0.00911 0.0101 0.0192
## 2 black  0.0744  0.0566 0.131
## 3 hispanic 0.0731 0.0577 0.131
## 4 native  0.00187 0.00132 0.00318
## 5 other   0.0173  0.0167 0.0340
## 6 white   0.360   0.322 0.682
```

Are race and gender independent in voting population in Florida?

Now, examine whether the joint probability of sex and race is equal to the product of the marginal probability of being a man times the marginal probability of being in each racial group.



Recall that a Bayesian perspective treats probability as a subjective opinion about how likely an event is. How should we change our beliefs after we make observations about the world?

Bayes' rule formalizes how we should update our beliefs based on evidence:

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

If we have a *prior* belief that event A has $P(A)$ chance of occurring, then we observe some data, represented as event B , we update our beliefs and obtain *posterior probability* $P(A|B)$.

Example: Detecting breast cancer

How good is a mammogram at detecting breast cancer?

What we know: One percent of women have breast cancer. 80 percent of people who have cancer and take a mammogram test positive. 9.6 percent of people who take a mammogram get a positive result when they do not have breast cancer.

If you take a mammogram and get a positive result, what is the probability that you have breast cancer?

$$P(\text{Cancer}) = 0.01$$

$$P(\text{Test positive}|\text{Cancer}) = 0.8$$

$$P(\text{Test positive}|\text{No cancer}) = 0.096$$

Let cancer be event A, and testing positive be event B.

We wish to know $P(A|B)$

Using Bayes' rule

The prior probability of having cancer is 0.01. How should we update our belief that someone has cancer based on a positive test?

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

Using the law of total probability, we can rewrite the denominator of Bayes' rule as:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\text{not } A)P(\text{not } A)}$$

Using Bayes' rule

We can apply Bayes' rule for A = Cancer, B = positive test:

$$P(\text{Cancer}|\text{Test positive}) =$$

$$\frac{P(\text{Test positive}|\text{Cancer})P(\text{Cancer})}{P(\text{Test positive}|\text{Cancer})P(\text{Cancer}) + P(\text{Test positive}|\text{No cancer})P(\text{No cancer})}$$

$$P(\text{Cancer}|\text{Test positive}) = \frac{0.8 \times 0.01}{0.8 \times 0.01 + 0.096 \times 0.99}$$

```
(0.8 * 0.01)/(0.8 * 0.01 + 0.096 * 0.99)
```

```
## [1] 0.07763975
```

Given these probabilities, the posterior likelihood that someone has cancer given a prior probability of one percent and a positive test is about 0.078

- No class next week, ASC
- Homework: 6.6.1 questions 1-4. Try 5 and 6 if you like, but not required
- Lab today: brief mid-semester assessment