Discrete Probability

Introduction to Discrete Probability

Unbiased die



Sample Space: $S = \{1,2,3,4,5,6\}$

All possible outcomes

Event: any subset of sample space

$$E_1 = \{3\} \qquad E_2 = \{2,5\}$$

Experiment: procedure that yields events

Throw die

Probability of event $E$:

$$p(E) = \frac{\text{size of event set}}{\text{size of sample space}} = \frac{|E|}{|S|}$$

Note that: $0 \le p(E) \le 1$

since $0 \le |E| \le |S|$

What is the probability that a die brings 3?

Event Space:   $E = \{3\}$

Sample Space:   $S = \{1,2,3,4,5,6\}$

Probability:   $p(E) = \dfrac{|E|}{|S|} = \dfrac{1}{6}$

What is the probability that a die brings 2 or 5?

Event Space:   $E = \{2,5\}$

Sample Space:   $S = \{1,2,3,4,5,6\}$

Probability:   $p(E) = \dfrac{|E|}{|S|} = \dfrac{2}{6}$

Two unbiased dice

Sample Space:   36 possible outcomes

$$S = \{(1,1),(1,2),(1,3),\ldots,(6,6)\}$$

First die      Second die
Ordered pair

What is the probability that two dice bring (1,1)?

Event Space:  $E = \{(1,1)\}$

Sample Space: $S = \{(1,1),(1,2),(1,3),\ldots,(6,6)\}$

Probability:   $p(E) = \dfrac{|E|}{|S|} = \dfrac{1}{36}$

What is the probability that
two dice bring same numbers?

Event Space: $E = \{(1,1), (2,2), (3,3), (4,4), (5,5), (6,6)\}$

Sample Space: $S = \{(1,1), (1,2), (1,3), \ldots, (6,6)\}$

Probability: $p(E) = \dfrac{|E|}{|S|} = \dfrac{6}{36}$

9

Game with unordered numbers

Game authority selects
a set of 6 winning numbers out of 40

Number choices: 1,2,3,...,40
i.e. winning numbers: 4,7,16,25,33,39

Player picks a set of 6 numbers
(order is irrelevant)

i.e. player numbers: 8,13,16,23,33,40

What is the probability that a player wins?

Konstantin Busch - LSU                10

Winning event:
$E = \{\{4,7,16,25,33,39\}\}$     $|E| = 1$
a single set with the 6 winning numbers

Sample space:
$S = \{\text{all subsets with 6 numbers out of 40}\}$
$= \{\{1,2,3,4,5,6\}, \{1,2,3,4,5,7\}, \{1,2,3,4,5,8\}, \ldots\}$

$|S| = \dbinom{40}{6} = 3{,}838{,}380$

Konstantin Busch - LSU                11

Probability that player wins:

$$P(E) = \frac{|E|}{|S|} = \frac{1}{\dbinom{40}{6}} = \frac{1}{3{,}838{,}380}$$

Konstantin Busch - LSU                12

3

## A card game

Deck has 52 cards

13 kinds of cards (2,3,4,5,6,7,8,9,10,a,k,q,j), each kind has 4 suits (h,d,c,s)

Player is given hand with 4 cards

What is the probability that the cards of the player are all of the same kind?

Event: $E = \{\{2_h, 2_d, 2_c, 2_s\}, \{3_h, 3_d, 3_c, 3_s\}, \dots, \{j_h, j_d, j_c, j_s\}\}$

$|E| = 13$    each set of 4 cards is of same kind

Sample Space:

$S = \{$all possible sets of 4 cards out of 52$\}$

$= \{\{2_h, 2_d, 2_c, 2_s\}, \{2_h, 2_d, 2_c, 3_h\}, \{2_h, 2_d, 2_c, 3_d\}, \dots\}$

$$|S| = \binom{52}{4} = \frac{52!}{4!\,48!} = \frac{52 \cdot 51 \cdot 50 \cdot 49}{4 \cdot 3 \cdot 2} = 270{,}725$$

Probability that hand has 4 same kind cards:

$$P(E) = \frac{|E|}{|S|} = \frac{13}{\binom{52}{4}} = \frac{13}{270{,}725}$$

## Game with ordered numbers

Game authority selects from a bin 5 balls in some order labeled with numbers 1…50

Number choices: 1,2,3,…,50
i.e. winning numbers: 37,4,16,33,9

Player picks a set of 5 numbers (order is important)

i.e. player numbers: 40,16,13,25,33

What is the probability that a player wins?

Sampling without replacement:
After a ball is selected
it is not returned to bin

Sample space size: 5-permutations of 50 balls

$$|S| = P(50,5) = \frac{50!}{(50-5)!} = \frac{50}{45!} = 50 \cdot 49 \cdot 48 \cdot 47 \cdot 46 = 245,251,200$$

Probability of success: $P(E) = \frac{|E|}{|S|} = \frac{1}{245,251,200}$

Sampling with replacement:
After a ball is selected
it is returned to bin

Sample space size: 5-permutations of 50 balls
with repetition

$$|S| = 50^5 = 312,500,000$$

Probability of success: $P(E) = \frac{|E|}{|S|} = \frac{1}{312,500,000}$

Probability of Inverse:  $p(\overline{E}) = 1 - p(E)$

Proof:  $\overline{E} = S - E$

$$p(\overline{E}) = \frac{|S - E|}{|S|} = \frac{|S| - |E|}{|S|} = 1 - \frac{|E|}{|S|} = 1 - p(E)$$

End of Proof

Example: What is the probability that
a binary string of 8 bits contains
at least one 0?

$E = \{01111111, 10111111, \ldots, 00111111, \ldots, 00000000\}$

$\overline{E} = \{11111111\}$

$$p(E) = 1 - p(\overline{E}) = 1 - \frac{|\overline{E}|}{|S|} = 1 - \frac{1}{2^8}$$

Probability of Union:    $E_1, E_2 \subseteq S$

$$p(E_1 \cup E_2) = p(E_1) + p(E_2) - p(E_1 \cap E_2)$$

Proof: $|E_1 \cup E_2| = |E_1| + |E_2| - |E_1 \cap E_2|$

$$p(E_1 \cup E_2) = \frac{|E_1 \cup E_2|}{|S|} = \frac{|E_1| + |E_1| - |E_1 \cap E_2|}{|S|}$$
$$= \frac{|E_1|}{|S|} + \frac{|E_2|}{|S|} - \frac{|E_1 \cap E_2|}{|S|}$$
$$= p(E_1) + p(E_2) - p(E_1 \cap E_2)$$

**End of Proof**

Example: What is the probability that
a binary string of 8 bits
starts with 0 **or** ends with 11?

Strings that start with 0:

$E_1 = \{00000000, 00000001, \ldots, 01111111\}$

$|E_1| = 2^7$  (all binary strings with 7 bits 0xxxxxxx)

Strings that end with 11:

$E_2 = \{00000011, 00000111, \ldots, 11111111\}$

$|E_2| = 2^6$  (all binary strings with 6 bits xxxxxx11)

Strings that start with 0 **and** end with 11:

$E_1 \cap E_2 = \{000000011, 00000111, \ldots, 01111111\}$

$|E_1 \cap E_2| = 2^5$ (all binary strings with 5 bits 0xxxxx11)

Strings that start with 0 **or** end with 11:

$$p(E_1 \cup E_2) = p(E_1) + p(E_2) - p(E_1 \cap E_2)$$
$$= \frac{|E_1|}{|S|} + \frac{|E_2|}{|S|} - \frac{|E_1 \cap E_2|}{|S|}$$
$$= \frac{2^7}{2^8} + \frac{2^6}{2^8} - \frac{2^5}{2^8} = \frac{1}{2} + \frac{1}{4} - \frac{1}{8} = \frac{5}{8}$$

Probability Theory

Sample space:   $S = \{x_1, x_2, \ldots, x_n\}$

Probability distribution function $p$:

$$0 \le p(x_i) \le 1$$

$$\sum_{x=1}^{n} p(x_i) = 1$$

Notice that it can be: $p(x_i) \neq p(x_j)$

Example: Biased Coin
    Heads (H) with probability 2/3
    Tails (T) with probability 1/3

Sample space: $S = \{H, T\}$

$p(H) = \dfrac{2}{3}$ $\quad$ $p(T) = \dfrac{1}{3}$ $\qquad$ $p(H) + p(T) = \dfrac{2}{3} + \dfrac{1}{3} = 1$

Uniform probability distribution:

$$p(x_i) = \frac{1}{n}$$

Sample space: $\quad S = \{x_1, x_2, \ldots, x_n\}$

Example: Unbiased Coin
  Heads (H) or Tails (T) with probability 1/2

$S = \{H, T\}$ $\qquad$ $p(H) = \dfrac{1}{2}$ $\qquad$ $p(T) = \dfrac{1}{2}$

Probability of event $E$:

$$E = \{x_1, x_2, \ldots, x_k\} \subseteq S$$

$$p(E) = \sum_{i=1}^{k} p(x_i)$$

For uniform probability distribution: $p(E) = \dfrac{|E|}{|S|}$

Example: Biased die $\quad S = \{1,2,3,4,5,6\}$

$$p(1) = p(2) = p(3) = p(4) = p(5) = \frac{1}{7} \quad p(6) = \frac{2}{7}$$

What is the probability
that the die outcome is 2 or 6? $E = \{2,6\}$

$$p(E) = p(2) + p(6) = \frac{1}{7} + \frac{2}{7} = \frac{3}{7}$$

Combinations of Events:

Complement: $p(\overline{E}) = 1 - p(E)$

Union: $p(E_1 \cup E_2) = p(E_1) + p(E_2) - p(E_1 \cap E_2)$

Union of disjoint events: $p\left(\bigcup_i E_i\right) = \sum_i p(E_i)$

Conditional Probability

Three tosses of an unbiased coin

Tails     Heads     Tails

Condition:   first coin is Tails

Question:   What is the probability that there is an odd number of Tails, given that first coin is Tails?

Sample space:
$$S = \{TTT, TTH, THT, THH, HTT, HTH, HHT, HHH\}$$

Restricted sample space given condition:
$$F = \{TTT, TTH, THT, THH\}$$

   first coin is Tails

Event without condition:
$$E = \{TTT, THH, HTH, HHT\}$$
   Odd number of Tails

Event with condition:
$$E_F = E \cap F = \{TTT, THH\}$$
   first coin is Tails

$$F = \{TTT, TTH, THT, THH\}$$

$$E_F = E \cap F = \{TTT, THH\}$$

**Given condition,**
**the sample space changes to** $F$

$$p(E_F) = \frac{|E \cap F|}{|F|} = \frac{|E \cap F| / |S|}{|F| / |S|} = \frac{p(E \cap F)}{p(F)} = \frac{2/8}{4/8} = 0.5$$

(the coin is unbiased)

**Notation of event with condition:**

$$E_F = E \mid F$$

event $E$ given $F$

$$p(E_F) = p(E \mid F) = \frac{p(E \cap F)}{p(F)}$$

**Conditional probability definition:**
(for arbitrary probability distribution)

**Given sample space** $S$ **with**
**events** $E$ **and** $F$ **(where** $p(F) > 0$ **)**
**the conditional probability of** $E$ **given** $F$ **is:**

$$p(E \mid F) = \frac{p(E \cap F)}{p(F)}$$

**Example:** What is probability that a family
of two children has two boys
given that one child is a boy

Assume equal probability to have boy or girl

**Sample space:** $S = \{BB, BG, GB, GG\}$

**Condition:** $F = \{BB, BG, GB\}$
one child is a boy

Event: $E = \{BB\}$
both children are boys

Conditional probability of event:

$$p(E \mid F) = \frac{p(E \cap F)}{p(F)} = \frac{p(\{BB\})}{p(\{BB, BG, GB\})} = \frac{1/4}{3/4} = \frac{1}{3}$$

Independent Events

Events $E_1$ and $E_2$ are independent iff:

$$p(E_1 \cap E_2) = p(E_1) p(E_2)$$

Equivalent definition (if $p(E_2) > 0$ ):

$$p(E_1 \mid E_2) = p(E_1)$$

Example: 4 bit uniformly random strings
$E_1$: a string begins with 1
$E_2$: a string contains even 1

$E_1 = \{1000, 1001, 1010, 1011, 1100, 1101, 1110, 1111\}$
$E_2 = \{0000, 0011, 0101, 0110, 1001, 1010, 1100, 1111\}$
$E_1 \cap E_2 = \{1111, 1100, 1010, 1001\}$

$|E_1| = |E_2| = 8$       $p(E_1) = p(E_2) = \frac{8}{16} = \frac{1}{2}$

$|E_1 \cap E_2| = 4$   $\boxed{p(E_1 \cap E_2) = \frac{4}{16} = \frac{1}{4} = \frac{1}{2} \cdot \frac{1}{2} = p(E_1) p(E_2)}$

Events $E_1$ and $E_2$ are independent

Bernoulli trial: Experiment with two outcomes:
success or failure

Success probability: $p$

Failure probability: $q = 1 - p$

Example:  Biased Coin

Success = Heads          Failure = Tails
$p = p(H) = \frac{2}{3}$          $q = p(T) = \frac{1}{3}$

Independent Bernoulli trials:

the outcomes of successive Bernoulli trials do not depend on each other

Example:   Successive coin tosses

Throw the biased coin 5 times



What is the probability to have 3 heads?

Heads probability: $p = \dfrac{2}{3}$   (success)

Tails probability:  $q = \dfrac{1}{3}$   (failure)

 HHHTT

HTHHT

HTHTH

THHTH

$\vdots$

Total numbers of ways to arrange in sequence 5 coins with 3 heads: $\dbinom{5}{3}$

Probability that any particular sequence has 3 heads and 2 tails is specified positions:

$$p^3 q^2$$

For example:



HHHTT
$pppqq = p^3 q^2$

HTHHT
$pqppq = p^3 q^2$

HTHTH
$pqpqp = p^3 q^2$

Probability of having 3 heads:

$$p^3q^2 + p^3q^2 + \cdots + p^3q^2 = \binom{5}{3}p^3q^2$$

1st
sequence
success
(3 heads)

2nd
sequence
success
(3 heads)

$\binom{5}{3}$st
sequence
success
(3 heads)

Throw the biased coin 5 times

Probability to have exactly 3 heads:

$$\binom{5}{3}p^3q^2 \quad = \frac{5!}{3!2!}\left(\frac{2}{3}\right)^3\left(\frac{1}{3}\right)^2 \approx 0.0086$$

Probability to have 3 heads and 2 tails
in specified sequence positions

All possible ways to arrange in sequence 5 coins with 3 heads

Theorem: Probability to have $k$ successes
in $n$ independent Bernoulli trials:

$$\binom{n}{k}p^k q^{n-k}$$

Also known as
binomial probability distribution:

$$b(k;n,p) = \binom{n}{k}p^k q^{n-k}$$

Proof:

$$\binom{n}{k}p^k q^{n-k}$$

Total number of
sequences with
$k$ successes and
$n-k$ failures

Probability that
a sequence has
$k$ successes and
$n-k$ failures
in specified positions

Example:
SFSFFS...SSF

End of Proof

**Example:** Random uniform binary strings
probability for 0 bit is 0.9
probability for 1 bit is 0.1

What is probability of 8 bit 0s out of 10 bits?

i.e. 0100001000

$$p = 0.9 \quad q = 0.1 \quad k = 8 \quad n = 10$$

$$b(k;n,p) = \binom{n}{k} p^k q^{n-k} = \binom{10}{8}(0.9)^8 (0.1)^2 = 0.1937102445$$

---

**Birthday Problem**

Birthday collision: two people have birthday
in same day

Problem:

How many people should be in a room
so that the probability of birthday collision
is at least ½?

Assumption: equal probability to be born in any day

---

366 days in a year

If the number of people is 367 or more
then birthday collision is guaranteed by
pigeonhole principle

Assume that we have $n \leq 366$ people

---

We will compute

$p_n$ :probability that $n$ people have
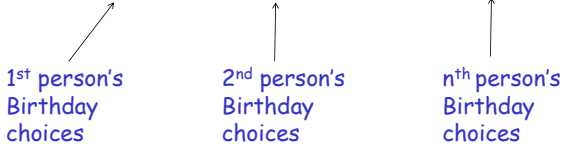all different birthdays

It will helps us to get

$1 - p_n$ :probability that there is
a birthday collision among $n$ people

**Sample space:**   Cartesian product

$$S = \{1, 2, \ldots, 366\} \times \{1, 2, \ldots, 366\} \times \cdots \times \{1, 2, \ldots, 366\}$$

1st person's Birthday choices        2nd person's Birthday choices        nth person's Birthday choices
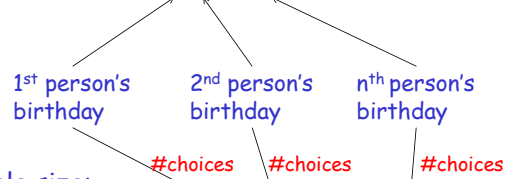
$$S = \{(1, 1, \ldots, 1), (2, 1, \ldots, 1), \ldots (366, 366, \ldots, 366)\}$$

**Sample space size:**   $|S| = 366 \cdot 366 \cdots 366 = 366^n$

**Event set:** each person's birthday is different

$$E = \{(1, 2, \ldots, 366), (366, 1, \ldots, 365), \ldots, (366, 365, \ldots, 1)\}$$

1st person's birthday        2nd person's birthday        nth person's birthday

#choices   #choices        #choices

**Sample size:**

$$|E| = P(366, n) = \frac{366!}{(366 - n)!} = 366 \cdot 365 \cdot 364 \cdots (366 - n + 1)$$

**Probability of no birthday collision**

$$p_n = \frac{|E|}{|S|} = \frac{366 \cdot 365 \cdot 364 \cdots (366 - n + 1)}{366^n}$$

**Probability of birthday collision:** $1 - p_n$

$$n = 22 \qquad 1 - p_n \approx 0.475$$

$$n = 23 \qquad 1 - p_n \approx 0.506$$

**Probability of birthday collision:** $1 - p_n$

$$n = 23 \qquad 1 - p_n \approx 0.506$$

**Therefore:** $n \geq 23$ people have probability at least ½ of birthday collision

The birthday problem analysis can be used to determine appropriate hash table sizes that minimize collisions

Hash function collision:   $h(k_1) = h(k_2)$

Randomized algorithms:
   algorithms with randomized choices (Example: quicksort)

Las Vegas algorithms:
   randomized algorithms whose output is always correct (i.e. quicksort)

Monte Carlo algorithms:
   randomized algorithms whose output is correct with some probability (may produce wrong output)

A Monte Carlo algorithm

```
Primality_Test( n,k ) {
    for(i = 1 to k) {
        b ← random_num ber(1,…,n)
        if (Miller_Test(n,b ) == failure)
            return(false)   // n is not prime
    }
    return(true) // most likely n is prime
}
```

```
Miller_Test( n,b ) {
```
$$n\text{-}1 = 2^s t$$
$$s,t \geq 0, \quad s \leq \log n, \quad t \text{ is odd}$$
```
    for ( j = 0 to s-1 ) {
```
if ( $b^t \equiv 1 (\bmod n)$ or $b^{2^j t} \equiv -1 (\bmod n)$ )
```
            return(success)
    }
    return(failure)
}
```

A prime number $n$ passes the Miller test for every $1 \le b \le n$

If the primality test algorithm returns false then the number is not prime for sure

A composite number $n$ passes the Miller test in range $1 \le b \le n$ for fewer than $\dfrac{n}{4}$ numbers

false positive with probability: $\dfrac{1}{4}$

If the algorithm returns true then the answer is correct (number is prime) with high probability:

$$1 - \left(\frac{1}{4}\right)^k = 1 - \frac{1}{n} \approx 1$$

for $k = \log_4 n$

## Bayes' Theorem

$$p(E) \ne 0 \qquad p(F) \ne 0$$

$$p(F \mid E) = \frac{p(E \mid F)\, p(F)}{p(E \mid F)\, p(F) + p(E \mid \overline{F})\, p(\overline{F})}$$

Applications: Machine Learning
Spam Filters

Bayes' Theorem Proof:

$$p(F \mid E) = \frac{p(E \cap F)}{p(E)}$$

$$p(E \mid F) = \frac{p(E \cap F)}{p(F)}$$

$$p(E \cap F) = p(F \mid E)\, p(E)$$

$$p(E \cap F) = p(E \mid F)\, p(F)$$

$$p(F \mid E)\, p(E) = p(E \mid F)\, p(F)$$

$$p(F \mid E) = \frac{p(E \mid F)\, p(F)}{p(E)}$$

16

$$E = (E \cap F) \cup (E \cap \overline{F})$$

$$(E \cap F) \cap (E \cap \overline{F}) = \varnothing$$

$$\Rightarrow p(E) = p(E \cap F) + p(E \cap \overline{F})$$

$$p(E \cap F) = p(E \mid F) p(F)$$

$$p(E \cap \overline{F}) = p(E \mid \overline{F}) p(\overline{F})$$

$$p(E) = p(E \mid F) p(F) + p(E \mid \overline{F}) p(\overline{F})$$

$$p(F \mid E) = \frac{p(E \mid F) p(F)}{p(E)}$$

$$p(E) = p(E \mid F) p(F) + p(E \mid \overline{F}) p(\overline{F})$$

$$p(F \mid E) = \frac{p(E \mid F) p(F)}{p(E \mid F) p(F) + p(E \mid \overline{F}) p(\overline{F})}$$

End of Proof

Example: Select random box
        then select random ball in box

Box 1          Box 2



Question:
If a red ball is selected,
what is the probability it was taken from box 1?

$E$: select red ball          $F$: select box 1

$\overline{E}$: select green ball     $\overline{F}$: select box 2

Question probability: $P(F \mid E) = ?$

Question:
If a red ball is selected,
what is the probability it was taken from box 1?

17

Bayes' Theorem:

$$p(F \mid E) = \frac{p(E \mid F)p(F)}{p(E \mid F)p(F) + p(E \mid \overline{F})p(\overline{F})}$$

We only need to compute:

$$p(F) \qquad p(\overline{F}) \qquad p(E \mid F) \qquad p(E \mid \overline{F})$$

$E$: select red ball $\qquad$ $F$: select box 1

$\overline{E}$: select green ball $\qquad$ $\overline{F}$: select box 2

Box 1 $\qquad\qquad$ Box 2



$p(F) = 1/2 = 0.5$ $\qquad$ $p(\overline{F}) = 1/2 = 0.5$

Probability to select box 1 $\qquad$ Probability to select box 2

$E$: select red ball $\qquad$ $F$: select box 1

$\overline{E}$: select green ball $\qquad$ $\overline{F}$: select box 2

Box 1 $\qquad\qquad$ Box 2



$p(E \mid F) = 7/9 = 0.777...$ $\qquad$ $p(E \mid \overline{F}) = 3/7 = 0.428....$

Probability to select red ball from box 1 $\qquad$ Probability to select red ball from box 2

$$p(F \mid E) = \frac{p(E \mid F)p(F)}{p(E \mid F)p(F) + p(E \mid \overline{F})p(\overline{F})}$$

$$p(F) = 1/2 = 0.5 \qquad p(\overline{F}) = 1/2 = 0.5$$

$$p(E \mid F) = 7/9 = 0.777... \qquad p(E \mid \overline{F}) = 3/7 = 0.428....$$

$$p(F \mid E) = \frac{0.777 \times 0.5}{0.777 \times 0.5 + 0.428 \times 0.5} = \frac{0.777}{0.777 + 0.428} = 0.644$$

Final result

18

### What if we had more boxes?

Generalized Bayes' Theorem:

$$p(F_j \mid E) = \frac{p(E \mid F_j)\,p(F_j)}{\displaystyle\sum_{i=1}^{n} p(E \mid F_i)\,p(F_i)}$$

Sample space $\quad S = F_1 \cup F_2 \cup \cdots \cup F_n$

mutually exclusive events

### Spam Filters

Training set:    Spam (bad) emails  $B$

Good emails    $G$

A user classifies each email
in training set as good or bad

### Find words that occur in $B$ and $G$

$n_B(w)$                    $n_G(w)$

number of spam emails          number of good emails
that contain word $w$          that contain word $w$

$$p(w) = \frac{n_B(w)}{|B|}$$            $$q(w) = \frac{n_G(w)}{|G|}$$

Probability that              Probability that
a spam email                  a good email
contains $w$                  contains $w$

### A new email X arrives

$S$:  Event that X is spam

$E$:  Event that X contains word $w$

What is the probability that X is spam
given that it contains word $w$?

$$P(S \mid E) = ?$$

Reject if this probability is at least 0.9

$$p(S \mid E) = \frac{p(E \mid S)p(S)}{p(E \mid S)p(S) + p(E \mid \bar{S})p(\bar{S})}$$

Training set for word "Rolex":

"Rolex" occurs in 250 of 2000 spam emails

"Rolex" occurs in 5 of 1000 good emails

We only need to compute:

$$p(S) \qquad p(\bar{S}) \qquad p(E \mid S) \qquad p(E \mid \bar{S})$$

$$\big| \qquad\quad \big| \qquad\quad \big| \qquad\qquad\qquad \backslash$$

0.5      0.5     $p(w) = \dfrac{n_B(w)}{|B|}$     $q(w) = \dfrac{n_G(w)}{|G|}$

simplified
assumption      Computed from training set

If new email contains word "Rolex"
what is the probability that it is a spam?

"Rolex" occurs in 250 of 2000 spam emails

$$n_B(Rolex) = 250$$

$$p(Rolex) = \frac{n_B(Rolex)}{|B|} = \frac{250}{2000} = 0.125$$

"Rolex" occurs in 5 of 1000 good emails

$$n_G(Rolex) = 5$$

$$q(Rolex) = \frac{n_G(Rolex)}{|G|} = \frac{5}{1000} = 0.005$$

If new email X contains word "Rolex"
what is the probability that it is a spam?

$S:$    Event that X is spam

$E:$    Event that X contains word "Rolex"
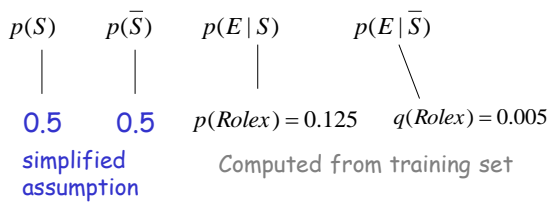
$$P(S \mid E) = ?$$

$$p(S \mid E) = \frac{p(E \mid S)p(S)}{p(E \mid S)p(S) + p(E \mid \bar{S})p(\bar{S})}$$

$$p(S \mid E) = \frac{0.125 \cdot 0.5}{0.125 \cdot 0.5 + 0.005 \cdot 0.5} = \frac{0.125}{0.13} = 0.961...$$

We only need to compute:

$p(S)$    $p(\bar{S})$    $p(E \mid S)$     $p(E \mid \bar{S})$

0.5     0.5    $p(Rolex) = 0.125$    $q(Rolex) = 0.005$

simplified
assumption       Computed from training set

New email is considered to be spam because:

$$p(S \mid E) = 0.961 \; > 0.9 \quad \text{spam threshold}$$

Better spam filters use two words:

$$p(S \mid E_1 \cap E_2) = \frac{p(E_1 \mid S)p(E_2 \mid S)}{p(E_1 \mid S)p(E_2 \mid S) + p(E_1 \mid \bar{S})p(E_2 \mid \bar{S})}$$

Assumption: $E_1$ and $E_2$ are independent

the two words appear
independent of each other