

Finish

# BLOOM FILTERS

ANNA KARLIN

MOST SLIDES BY SHREYA JAYARAMAN, LUXI WANG, ALEX TSUN

# BLOOM FILTERS: MOTIVATION



- Large universe of possible data items.
- Hash table is stored on disk or in network, so any lookup is expensive.
- Many (if not most) of the lookups return “Not found”.

Altogether, this is bad. You're wasting **a lot of time and space** doing lookups for items that aren't even present.

Examples:

- Google Chrome: wants to warn you if you're trying to access a malicious URL. Keep hash table of malicious URLs.
- Network routers: want to track source IP addresses of certain packets, .e.g., blocked IP addresses.

# BLOOM FILTERS: MOTIVATION



- Probabilistic data structure.
- Close cousins of hash tables.
- Ridiculously space efficient
- To get that, make occasional errors, specifically false positives.

Typical implementation: only 8 bits per element!

# BLOOM FILTERS



- Stores information about a set of elements.
- Supports two operations:
  1. **add(x)** - adds x to bloom filter
  2. **contains(x)** - returns true if x in bloom filter, otherwise returns false
    - a. If return false, **definitely** not in bloom filter.
    - b. If return true, **possibly** in the structure (some false positives).

# BLOOM FILTERS: EXAMPLE

**bloom filter  $t$  with  $m = 5$  that uses  $k = 3$  hash functions**

```
function INITIALIZE( $k, m$ )  
  for  $i = 1, \dots, k$ : do  
     $t_i =$  new bit vector of  $m$  0's
```

Index →	0	1	2	3	4
$t_1$	0	0	0	0	0
$t_2$	0	0	0	0	0
$t_3$	0	0	0	0	0

# BLOOM FILTERS: EXAMPLE

bloom filter  $t$  of length  $m = 5$  that uses  $k = 3$  hash functions

```
function ADD(x)
  for  $i = 1, \dots, k$ : do
     $t_i[h_i(x)] = 1$ 
```

add("thisisavirus.com")

$h_1$ ("thisisavirus.com")  $\rightarrow$  2

$h_2$ ("thisisavirus.com")  $\rightarrow$  1

$h_3$ ("thisisavirus.com")  $\rightarrow$  4

Index $\rightarrow$	0	1	2	3	4
$t_1$	0	0	1	0	0
$t_2$	0	1	0	0	0
$t_3$	0	0	0	0	1

# BLOOM FILTERS: EXAMPLE

bloom filter  $t$  of length  $m = 5$  that uses  $k = 3$  hash functions

```
function CONTAINS(x)
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

True

contains("thisisavirus.com")

$h_1$ ("thisisavirus.com")  $\rightarrow$  2

$h_2$ ("thisisavirus.com")  $\rightarrow$  1

$h_3$ ("thisisavirus.com")  $\rightarrow$  4

Index $\rightarrow$	0	1	2	3	4
$t_1$	0	0	1	0	0
$t_2$	0	1	0	0	0
$t_3$	0	0	0	0	1

# BLOOM FILTERS: EXAMPLE

bloom filter  $t$  of length  $m = 5$  that uses  $k = 3$  hash functions

```
function CONTAINS(x)  
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

True

contains("thisisavirus.com")

$h_1(\text{"thisisavirus.com"}) \rightarrow 2$

$h_2(\text{"thisisavirus.com"}) \rightarrow 1$

$h_3(\text{"thisisavirus.com"}) \rightarrow 4$

Since all conditions satisfied, returns True (correctly)

Index →	0	1	2	3	4
$t_1$	0	0	1	0	0
$t_2$	0	1	0	0	0
$t_3$	0	0	0	0	1



# BLOOM FILTERS: EXAMPLE

bloom filter  $t$  of length  $m = 5$  that uses  $k = 3$  hash functions

**function** CONTAINS( $x$ )  
**return**  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$

**True**

**True**

**True**

**contains("verynormalsite.com")**

$h_1(\text{"verynormalsite.com"}) \rightarrow 2$

$h_2(\text{"verynormalsite.com"}) \rightarrow 0$

$h_3(\text{"verynormalsite.com"}) \rightarrow 4$

Since all conditions satisfied, returns **True** (incorrectly)

Index →	0	1	2	3	4
$t_1$	0	1	1	0	0
$t_2$	1	1	0	0	0
$t_3$	0	0	0	0	1

# BLOOM FILTERS: SUMMARY



- An empty bloom filter is an empty  $k \times m$  bit array with all values initialized to zeros
  - $k$  = number of hash functions
  - $m$  = size of each array in the bloom filter
- `add(x)` runs in  $O(k)$  time
- `contains(x)` runs in  $O(k)$  time
- requires  $O(km)$  space (in bits!)
- Probability of false positives from collisions can be reduced by increasing the size of the bloom filter

# BLOOM FILTERS: APPLICATION



- Google Chrome has a database of malicious URLs, but it takes a long time to query.
- Want an in-browser structure, so needs to be efficient and be space-efficient
- Want it so that can check if a URL is in structure:
  - If return False, then definitely not in the structure (don't need to do expensive database lookup, website is safe)
  - If return True, the URL may or may not be in the structure. Have to perform expensive lookup in this rare case.

Added functions  
to Bloom filter

## FALSE POSITIVE PROBABILITY

[new item  $x$  shows up.]

### Assumptions

$\left\{ \begin{array}{l} h_i(y_j) \\ h_i(x) \end{array} \right. \quad \begin{array}{l} j=1, \dots, n \\ i=1, \dots, k \end{array}$   
all mutually independent  
& uniformly distributed over  
corresponding array  
 $\underline{h_i(y_j)} \sim \underline{\text{Unif}}\{0, 1, \dots, m-1\}$

$\Pr(\text{false pos on } x)$

$$= \Pr(t_1[\underbrace{h_1(x)}_{a_1}] = 1, t_2[\underbrace{h_2(x)}_{a_2}] = 1, \dots, t_k[\underbrace{h_k(x)}_{a_k}] = 1)$$

$$\stackrel{\text{indep}}{=} \prod_{i=1}^k \Pr(t_i[a_i] = 1)$$

$\Pr(\exists y_j \text{ s.t. } h_i(y_j) = a_i)$

$$= 1 - \Pr(\forall j \quad h_i(y_j) \neq a_i)$$

$$\stackrel{\text{indep}}{=} \prod_{j=1}^n \underbrace{\Pr(h_i(y_j) \neq a_i)}_{1 - \frac{1}{m}}$$

$$(1 - \frac{1}{m})^n$$

$$= 1 - (1 - \frac{1}{m})^n$$

$$= \left[ 1 - (1 - \frac{1}{m})^n \right]^k$$

$$\leq 0.05$$

# COMPARISON WITH HASH TABLES - SPACE



- Google storing 5 million URLs, each URL 40 bytes.
- Bloom filter with k=8 and m = 10,000,000.

Hash Table

$$5,000,000 \times 40 \\ = 200 \text{ MB}$$

Bloom Filter

$$\frac{10,000,000 \cdot 8}{8} = 10 \text{ MB}$$

# COMPARISON WITH HASH TABLES - TIME



- Say avg user visits 100,000 URLs in a year, of which 2,000 are malicious.
- 0.5 seconds to do lookup in the database, 1ms for lookup in Bloom filter.
- Suppose the false positive rate is 2%

## Hash Table

$$100,000 \times 0.5 \text{ sec} \\ = 50,000 \text{ sec}$$

## Bloom Filter

$$100,000 \times 1 \text{ ms} \\ + 2000 \times 0.5 \text{ sec} \\ + 100,000 \cdot 0.02 \cdot 0.5 \text{ sec} \\ = 2100 \text{ sec.}$$

# BLOOM FILTERS: MANY APPLICATIONS



- Any scenario where space and efficiency are important.
- Used a lot in networking
- In distributed systems when want to check consistency of data across different locations, might send a Bloom filter rather than the full set of data being stored.
- Google BigTable uses Bloom filters to reduce the disk lookups for non-existent rows and columns
- Internet routers often use Bloom filters to track blocked IP addresses.
- And on and on...

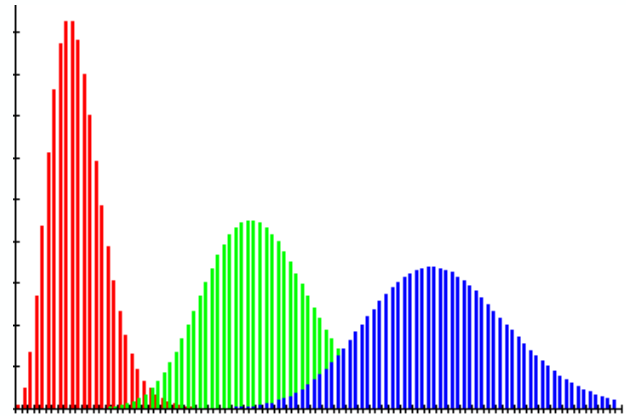


# BLOOM FILTERS TYPICAL EXAMPLE...

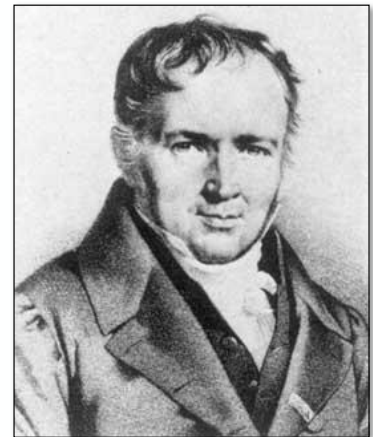
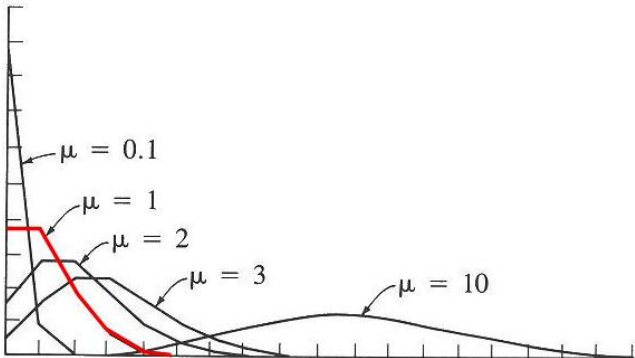


of randomized algorithms and randomized data structures.

- Simple
  - Fast
  - Efficient
  - Elegant
  - Useful!
- 
- You'll be implementing Bloom filters on pset 4. Enjoy!



# a zoo of (discrete) random variables



## discrete uniform random variables

---

A discrete random variable  $X$  equally likely to take any (integer) value between integers  $a$  and  $b$ , inclusive, is *uniform*.

$$\text{Range}(X) = \{a, a+1, \dots, b\}$$

Notation:  $\text{Unif}(a, b)$

Probability mass function:

$$P_X(k) = \begin{cases} \frac{1}{b-a+1} & k \in \text{Range}(X) \\ 0 & \text{o.w.} \end{cases}$$

Mean:  $\frac{a+b}{2}$

Variance:

## discrete uniform random variables

A discrete random variable  $X$  **equally likely** to take any (integer) value between integers  $a$  and  $b$ , inclusive, is **uniform**.

Notation:  $X \sim \text{Unif}(a,b)$

Probability:  $P(X = i) = \frac{1}{b - a + 1}$

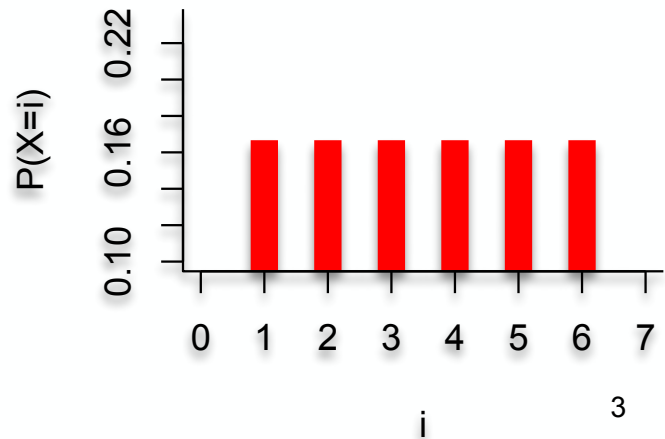
Mean, Variance:  $E[X] = \frac{a + b}{2}$ ,  $\text{Var}[X] = \frac{(b - a)(b - a + 1)}{12}$

Example: value shown on one roll of a fair die is  $\text{Unif}(1,6)$ :

$$P(X=i) = 1/6$$

$$E[X] = 7/2$$

$$\text{Var}[X] = 35/12$$



# indicator

# Bernoulli random variables

An experiment results in "Success" or "Failure"

$X$  is an *indicator random variable* (1 = success, 0 = failure)

$$P(X=1) = p \quad \text{and} \quad P(X=0) = 1-p$$

$X$  is called a *Bernoulli* random variable:  $X \sim \text{Ber}(p)$

Mean:

Variance:

$$\text{Var}(X) = E(X^2) - (E(X))^2 = p - p^2$$
$$E(X^2) = p$$
$$E(X) = 1 \cdot p + 0 \cdot (1-p)$$

	mean	variance
a)	$p$	$p$
b)	$p$	$1-p$
c)	$p$	$p(1-p)$
d)	$1-p$	$p(1-p)$

# Bernoulli random variables

An experiment results in “Success” or “Failure”

$X$  is an *indicator random variable* (1 = success, 0 = failure)

$$P(X=1) = p \quad \text{and} \quad P(X=0) = 1-p$$

$X$  is called a *Bernoulli* random variable:  $X \sim \text{Ber}(p)$

$$E[X] = E[X^2] = p$$

$$\text{Var}(X) = E[X^2] - (E[X])^2 = p - p^2 = p(1-p)$$

Examples:

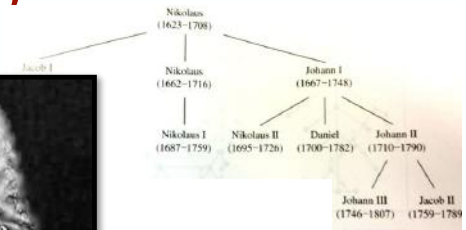
coin flip

random binary digit  $\leftarrow$

whether a disk drive crashed  $\uparrow$



Jacob (aka James, Jacques)  
Bernoulli, 1654 – 1705



# binomial random variables

Consider  $n$  independent random variables  $Y_i \sim \text{Ber}(p)$

$X = \sum_i Y_i$  is the number of successes in  $n$  trials *indep.*

$X$  is a Binomial random variable:  $X \sim \text{Bin}(n, p)$

p.m.f.?

$$P_X(k) = P(X=k) =$$

	$P_X(k)$
a)	$p^k (1-p)^{n-k}$
b)	$np$
c)	$\binom{n}{k} p^k (1-p)^{n-k}$
d)	$\binom{n}{n-k} p^k (1-p)^{n-k}$

## Examples

- # of heads in  $n$  coin flips
- # of 1's in a randomly generated length  $n$  bit string
- # of disk drive crashes in a 1000 computer cluster
- # bit errors in file written to disk
- # of typos in a book
- # of elements in particular bucket of large hash table
- # of server crashes per day in giant data center

# binomial random variables

Consider  $n$  independent random variables  $Y_i \sim \text{Ber}(p)$

$X = \sum_i Y_i$  is the number of successes in  $n$  trials

$X$  is a *Binomial* random variable:  $X \sim \text{Bin}(n,p)$

$$\underline{X = Y_1 + Y_2 + \dots + Y_n}$$

$$E(Y_i) = p$$
$$\underline{\text{Var}(Y_i) = p(1-p)}$$

Probability mass function:

$$P_X(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$k = 0, 1, \dots, n$$

Mean:

Variance:

	<u>mean</u>	<u>variance</u>
a)	$p$	$p(1-p)$
b)	$np$	$np(1-p)$
c)	$np$	$np^2$
d)	$np$	$n^2p$



## mean, variance of the binomial (II)

---

If  $Y_1, Y_2, \dots, Y_n \sim \text{Ber}(p)$  and independent,

then  $X = \sum_{i=1}^n Y_i \sim \text{Bin}(n, p)$ .

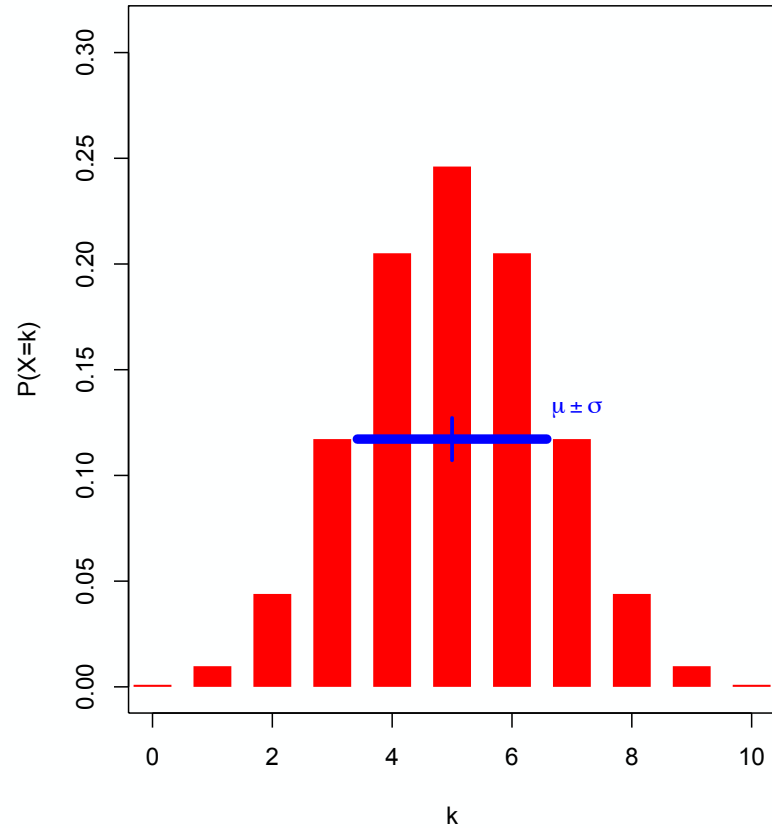
$$E[X] = np$$

$$E[X] = E \left[ \sum_{i=1}^n Y_i \right] = \sum_{i=1}^n E[Y_i] = nE[Y_1] = np$$

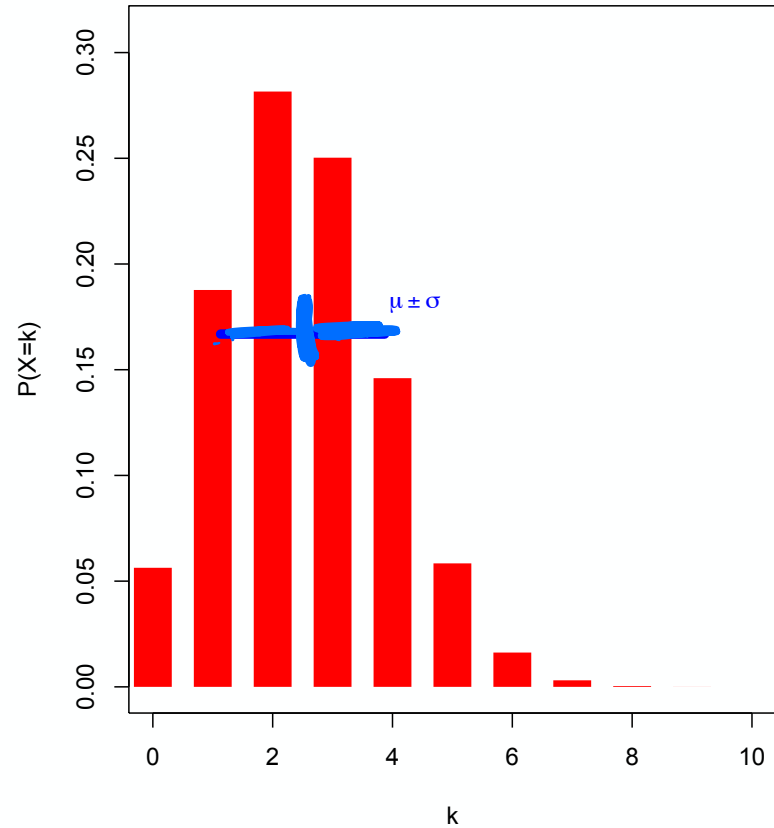
$$\text{Var}[X] = np(1 - p)$$

$$\text{Var}[X] = \text{Var} \left[ \sum_{i=1}^n Y_i \right] = \sum_{i=1}^n \text{Var}[Y_i] = n\text{Var}[Y_1] = np(1 - p)$$

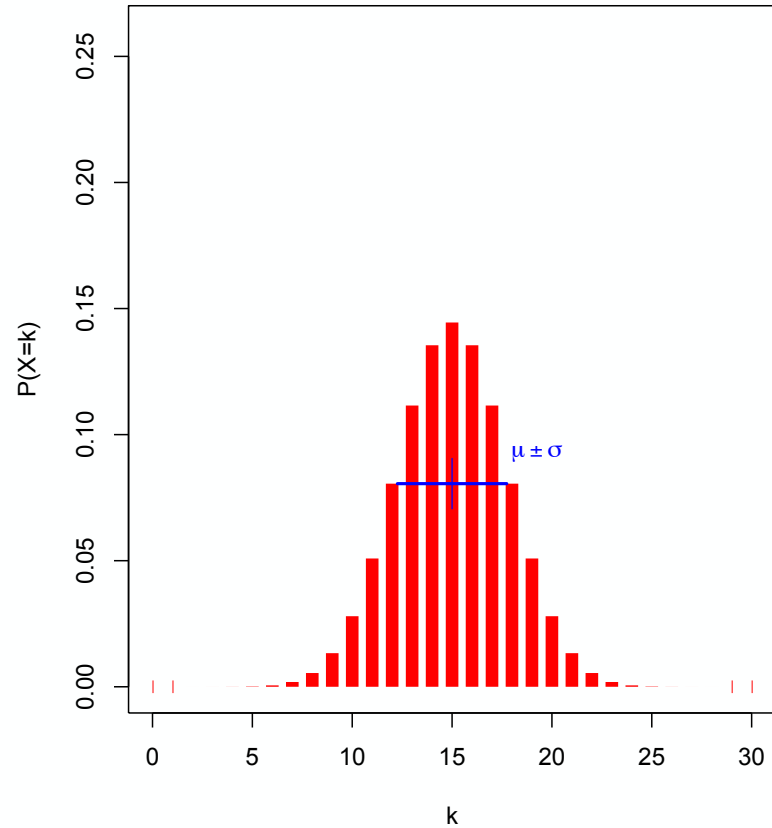
PMF for  $X \sim \text{Bin}(10, 0.5)$



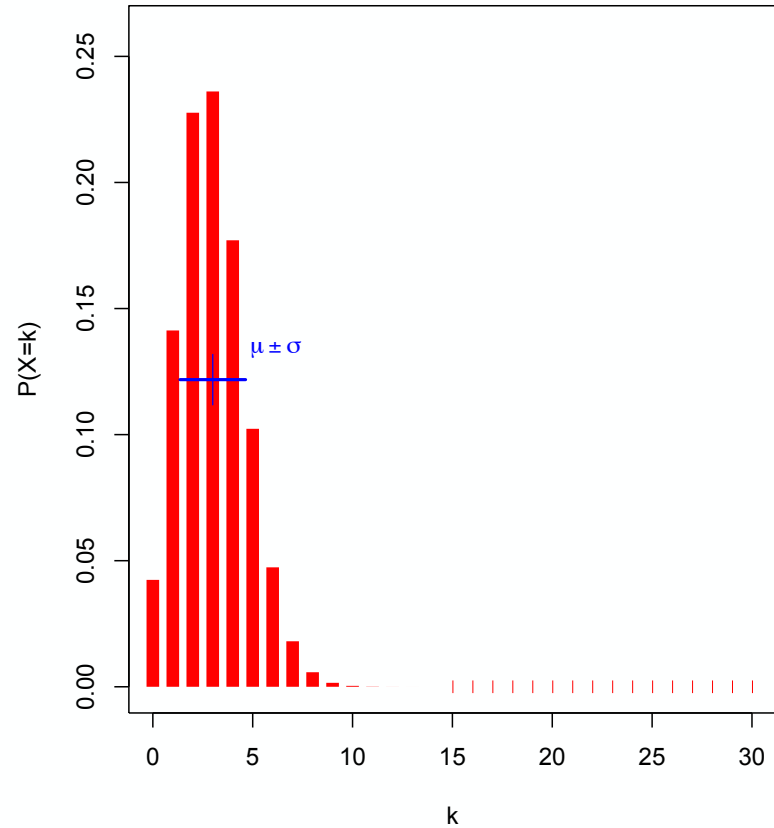
PMF for  $X \sim \text{Bin}(10, 0.25)$



PMF for  $X \sim \text{Bin}(30, 0.5)$



PMF for  $X \sim \text{Bin}(30, 0.1)$



Sending a bit string over the network

$n = 4$  bits sent, each corrupted with probability  $0.1$

$X = \#$  of corrupted bits,  $X \sim \text{Bin}(4, 0.1)$

In real networks, large bit strings (length  $n \approx 10^4$ )

Corruption probability is very small:  $p \approx 10^{-6}$

$X \sim \text{Bin}(10^4, 10^{-6})$  is unwieldy to compute

$n, p$

$n \cdot p$

small.

heads up

Extreme  $n$  and  $p$  values arise in many cases

# bit errors in file written to disk

# of typos in a book

# of elements in particular bucket of large hash table

# of server crashes per day in giant data center

# geometric distribution

In a series  $X_1, X_2, \dots$  of Bernoulli trials with success probability  $p$ , let  $Y$  be the index of the first success, i.e.,

$$X_1 = X_2 = \dots = X_{Y-1} = 0 \text{ \& } X_Y = 1$$

# trials up till & including first success

Then  $Y$  is a *geometric* random variable with parameter  $p$ .

$$Y \sim \text{Geo}(p) \\ \text{Geom}(p)$$

Examples:

Number of coin flips until first head

Number of blind guesses on SAT until I get one right

Number of darts thrown until you hit a bullseye

Number of random probes into hash table until empty slot

Number of wild guesses at a password until you hit it

$$\Pr(Y=k) = \Pr(\text{first success is on } k^{\text{th}} \text{ trial})$$

Probability mass function:  $= p_x(k)$

$$\frac{p_x(k)}{p_x(k)}$$

a)  $(1-p)^k p$

b)  $(1-p)^{k-1} p$

c)  $(1-p)^k$

d)  $p^k$

Mean:

$$\frac{1}{p}$$

Variance:

## geometric distribution

In a series  $X_1, X_2, \dots$  of Bernoulli trials with success probability  $p$ , let  $Y$  be the index of the first success, i.e.,

$$X_1 = X_2 = \dots = X_{Y-1} = 0 \text{ \& } X_Y = 1$$

Then  $Y$  is a *geometric* random variable with parameter  $p$ .

Examples:

Number of coin flips until first head

Number of blind guesses on SAT until I get one right

Number of darts thrown until you hit a bullseye

Number of random probes into hash table until empty slot

Number of wild guesses at a password until you hit it

$$P(Y=k) = (1-p)^{k-1}p;$$

Mean  $1/p;$

Variance  $(1-p)/p^2$

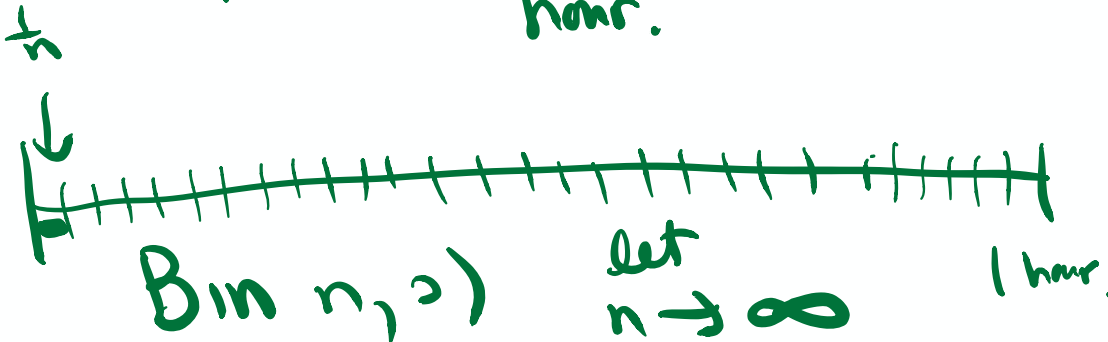
# Poisson motivation

Model: # events that occur in an hour

- rate 3 events per hour.
- exp # events that occur in interval of  $t$  hours  
 $3t$
- occurrence of events in disjoint time intervals is indep.



$Y$  # events that occur in 1 hour.



- $p = ?$
- a)  $\frac{3}{n}$
  - b)  $3n$
  - c) 3
  - d)  $\frac{3}{60}$