

Bernoulli distribution

The simplest non-uniform distribution

p – probability of success (1)

$1-p$ – probability of failure (0)

$$f(x) = P(X = x) = \begin{cases} p & \text{if } x = 1 \\ 1 - p & \text{if } x = 0 \end{cases}$$

Jacob Bernoulli

(1654-1705)

Swiss mathematician (Basel)

- Law of large numbers
- Mathematical constant $e=2.718\dots$



Bernoulli distribution

$$f(x) = P(X = x) = \begin{cases} p & \text{if } x = 1 \\ 1 - p & \text{if } x = 0 \end{cases}$$

$$E(X) = 0 \times P(X = 0) + 1 \times P(X = 1) = 0(1 - p) + 1(p) = p$$

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

Refresher: Binomial Coefficients

$$\binom{n}{k} = C_k^n = \frac{n!}{k!(n-k)!}, \text{ called } n \text{ choose } k$$

$$\binom{10}{3} = C_3^{10} = \frac{10!}{3!7!} = \frac{10 \cdot 9 \cdot 8 \cdot 7!}{3 \cdot 2 \cdot 1 \cdot 7!} = 120$$

Number of ways to choose k objects out of n

without replacement and where the **order does not matter**.

Called binomial coefficients because of the binomial formula

$$(p+q)^n = (p+q) \times (p+q) \dots \times (p+q) = \sum_{x=0}^n C_x^n p^x q^{n-x}$$

Binomial Distribution

- **Binomially-distributed** random variable X equals **sum (number of successes) of n independent Bernoulli trials**
- The probability mass function is:

$$f(x) = C_x^n p^x (1-p)^{n-x} \quad \text{for } x = 0, 1, \dots, n \quad (3-7)$$

$q = 1-p$

- Based on the binomial expansion:

$$1 = (p+q)^n = \sum_{x=0}^n C_x^n p^x q^{n-x}$$

Binomial Mean

X is a binomial random variable
with parameters p and n

Mean:

$$\mu = E(X) = np$$

$$\begin{aligned}\mu &= \sum x C_x^n p^x q^{n-x} = p \frac{\partial}{\partial p} \sum C_x^n p^x q^{n-x} = \\ &= p \frac{\partial}{\partial p} (p + q)^n = np\end{aligned}$$

$$\begin{aligned}
E(X(X-1)) &= \\
&= \sum x(x-1) C_x^n p^x q^{n-x} \\
&= p^2 \frac{\partial^2}{\partial p^2} \sum C_x^n p^x q^{n-x} = \\
&= p^2 \frac{\partial^2}{\partial p^2} (p+q)^n \Big|_{q=1-p} = n(n-1)p^2
\end{aligned}$$

$$\begin{aligned}
E(X^2) &= E(X(X-1)) + E(X) = \\
&= n^2 p^2 - n p^2 + n p = n^2 p^2 + n p (1-p)
\end{aligned}$$

$$\begin{aligned}
V(X) &= E(X^2) - E(X)^2 = n^2 p^2 + n p (1-p) - (np)^2 \\
&= \boxed{np(1-p)}
\end{aligned}$$

Binomial mean, variance and standard deviation

Let X be a binomial random variable with parameters p and n

- Mean:

$$\mu = np$$

- Variance:

$$\sigma^2 = V(X) = np(1-p)$$

- Standard deviation:

$$\sigma = \sqrt{np(1-p)}$$

- Standard deviation to mean ratio

$$\sigma/\mu = \sqrt{np(1-p)}/np = \frac{\sqrt{(1-p)/p}}{\sqrt{n}}$$

Matlab exercise: Binomial distribution

- Generate a **sample of size 100,000** for binomially-distributed random variable X with $n=100$, $p=0.2$
- Tip: generate n Bernoulli random variables and use `sum` to add them up
- Plot the approximation to the **Probability Mass Function** based on this sample
- Calculate the mean and variance of this sample and compare it to **theoretical calculations**:
 $E[X]=n*p$ and $V[X]=n*p*(1-p)$

Poisson Distribution

- Limit of the binomial distribution when
 - n , the **number of attempts**, is very **large**
 - p , the **probability of success** is very **small**
 - $E(X) = np = \lambda$ is $O(1)$

The annual numbers of deaths from horse kicks in 14 Prussian army corps between 1875 and 1894

Number of deaths	of Observed frequency	Expected frequency
0	144	139
1	91	97
2	32	34
3	11	8
4	2	1
5 and over	0	0
Total	280	280

From von Bortkiewicz 1898



Siméon Denis Poisson
(1781–1840)
French mathematician
and physicist

Let $\lambda = np = E(x)$, so $p = \frac{\lambda}{n}$

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

$$= \frac{n(n-1)\dots(n-x+1)}{x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \sim \frac{n^x}{x!} \left(\frac{\lambda}{n}\right)^x = \frac{\lambda^x}{x!};$$

$$\sum_x \frac{\lambda^x}{x!} = e^\lambda.$$

Normalization requires $\sum_x P(X = x) = 1$.

$$\text{Thus } P(X = x) = \frac{\lambda^x}{x!} e^{-\lambda}$$

Poisson Mean & Variance

If X is a Poisson random variable, then:

- Mean: $\mu = E(X) = \lambda \approx n \cdot p$
- Variance: $\sigma^2 = V(X) = \lambda \approx n \cdot p \cdot (1 - p) \approx n \cdot p$
- Standard deviation: $\sigma = \lambda^{1/2}$

Note: Variance = Mean

Note: Standard deviation/Mean = $\lambda^{-1/2}$
decreases with λ

Are you in class?

- A. Yes
- B. No
- C. I am not sure, I am still asleep

Get your i-clickers

Matlab exercise: Poisson distribution

- Generate a **sample of size 100,000** for Poisson-distributed random variable X with $\lambda = 2$
- Plot the approximation to the **Probability Mass Function** based on this sample
- Calculate the mean and variance of this sample and compare it to **theoretical calculations**:
 $E[X] = \lambda$ and $V[X] = \lambda$

Matlab exercise: Poisson distribution

- **Stats=100000; lambda=2;**
- **r2=random('Poisson',lambda,Stats,1);**
- **mean(r2)**
- **var(r2)**
- **[a,b]=hist(r2, 0:max(r2));**
- **p_p=a./sum(a);**
- **figure; stem(b,p_p);**
- **figure; semilogy(b,p_p,'ko-')**

Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

WHY AREN'T ECONOMISTS RICH

WHY ARE THERE SO MANY CROWS IN ROCHESTER, MN
WHY IS THERE PHLEGM

WHY DO AMERICANS CALL IT SOCCER

WHY IS PSYCHIC WEAK TO BUG

WHY ARE MY EARS RINGING

WHY DO CHILDREN GET CANCER

WHY ARE THERE SO MANY AVENGERS

WHY IS POSEIDON ANGRY WITH ODYSSEUS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS THERE ICE IN SPACE

WHY IS WOLVERINE NOT IN THE AVENGERS

WHY ARE THERE ANTS IN MY LAPTOP

WHY IS EARTH TILTED
WHY IS SPACE BLACK
WHY IS OUTER SPACE SO COLD
WHY ARE THERE PYRAMIDS ON THE MOON
WHY IS NASA SHUTTING DOWN



WHY IS THERE AN OWL IN MY BACKYARD
WHY IS THERE AN OWL OUTSIDE MY WINDOW
WHY IS THERE AN OWL ON THE DOLLAR BILL
WHY DO OWLS ATTACK PEOPLE
WHY ARE AK 47s SO EXPENSIVE
WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE
WHY ARE THERE GODS
WHY ARE THERE TWO SPOCKS

WHY ARE DOGS AFRAID OF FIREWORKS
WHY IS THERE NO KING IN ENGLAND

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK
WHY AREN'T THERE E GRADES
WHY IS ISOLATION BAD
WHY DO BOYS LIKE ME
WHY DON'T BOYS LIKE ME
WHY IS THERE ALWAYS A JAVA UPDATE
WHY ARE THERE RED DOTS ON MY THIGHS



WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS
WHY DO I SAY UH
WHY IS SEA SALT BETTER
WHY ARE THERE TREES IN THE MIDDLE OF FIELDS
WHY IS THERE NOT A POKEMON MMO
WHY IS THERE LAUGHING IN TV SHOWS
WHY ARE THERE DOORS ON THE FREEWAY
WHY ARE THERE SO MANY SVCHOST.EXE RUNNING
WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA
WHY ARE THERE SCARY SOUNDS IN MINECRAFT
WHY IS THERE KICKING IN MY STOMACH
WHY ARE THERE TWO SLASHES AFTER HTTP
WHY ARE THERE CELEBRITIES
WHY DO SNAKES EXIST
WHY DO OYSTERS HAVE PEARLS
WHY ARE DUCKS CALLED DUCKS
WHY DO THEY CALL IT THE CLAP
WHY ARE KYLE AND CARTMAN FRIENDS
WHY IS THERE AN ARROW ON AANG'S HEAD
WHY ARE TEXT MESSAGES BLUE
WHY ARE THERE MUSTACHES ON CLOTHES
WHY ARE THERE MUSTACHES ON CARS
WHY ARE THERE MUSTACHES EVERYWHERE
WHY ARE THERE SO MANY BIRDS IN OHIO
WHY IS THERE SO MUCH RAIN IN OHIO
WHY IS OHIO WEATHER SO WEIRD
WHY ARE THERE MALE AND FEMALE BIKES

WHY AREN'T THERE DINOSAUR GHOSTS

WHY ARE THERE BRIDESMAIDS
WHY DO DYING PEOPLE REACH UP
WHY AREN'T THERE VARICOSE ARTERIES
WHY ARE OLD KLINGONS DIFFERENT



WHY IS PROGRAMMING SO HARD
WHY IS THERE A 0 OHM RESISTOR
WHY DO AMERICANS HATE SOCCER
WHY DO RHYMES SOUND GOOD
WHY DO TREES DIE
WHY IS THERE NO SOUND ON CNN
WHY AREN'T POKEMON REAL
WHY AREN'T BULLETS SHARP
WHY DO DREAMS SEEM SO REAL

WHY IS THERE HELL IF GOD FORGIVES

WHY IS GPS FREE

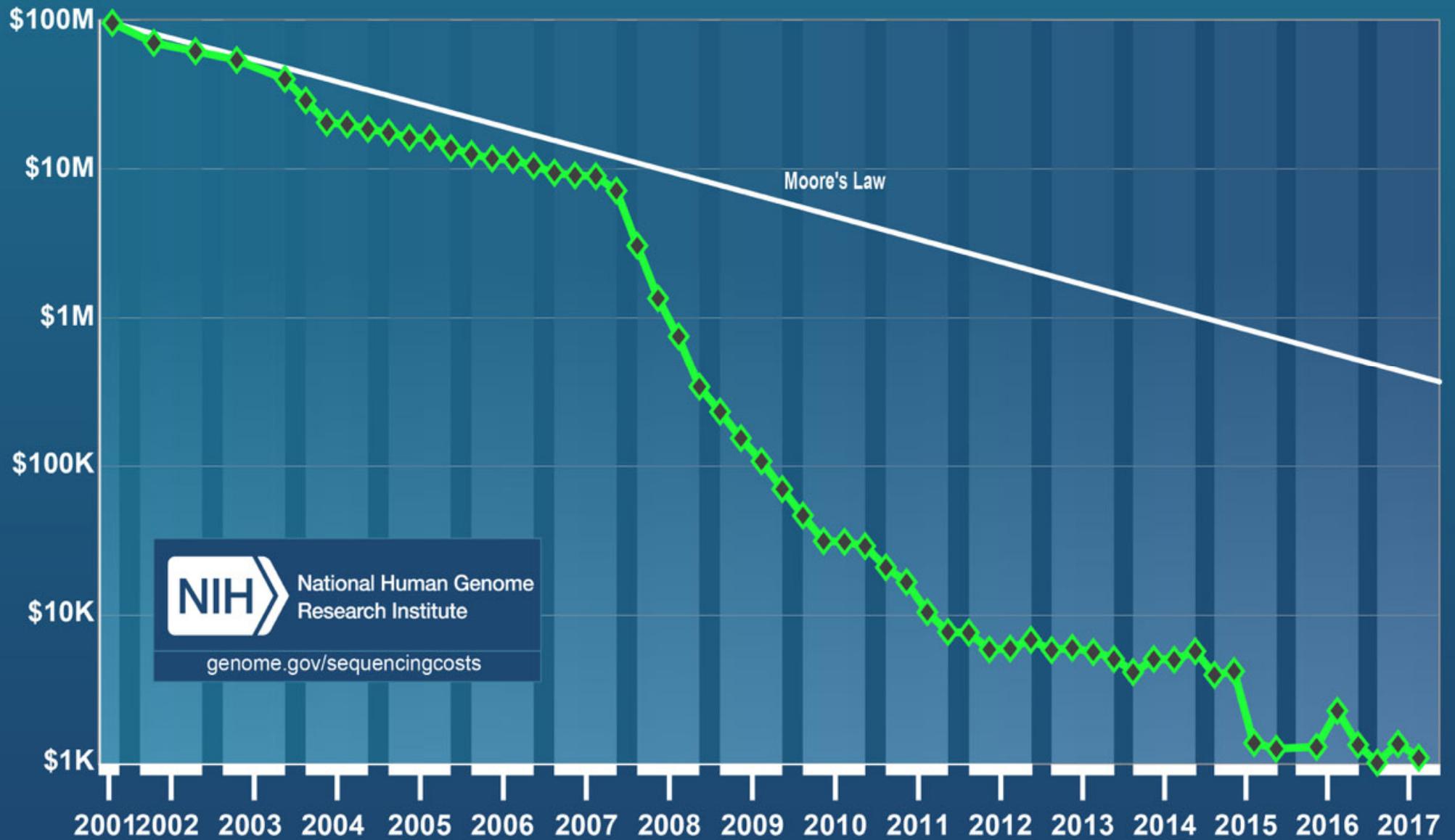
WHY IS LIFE SO BORING



WHY ARE ULTRASOUNDS IMPORTANT
WHY ARE ULTRASOUND MACHINES EXPENSIVE
WHY IS STEALING WRONG
WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

Poisson Distribution in Genome Assembly

Cost per Genome



Poisson Example: Genome Assembly

- **Goal:** DNA sequence (ACTG) of the entire genome
- **Problem:** Sequencers generate random short reads

Sequencer	Sanger 3730xl	454 GS	Ion Torrent	SOLiDv4	Illumina HiSeq 2000	Pac Bio
Mechanism	Dideoxy chain termination	Pyrosequencing	Detection of hydrogen ion	Ligation and two-base coding	Reversible Nucleotides	Single molecule real time
Read length	400-900 bp	700 bp	~400 bp	50 + 50 bp	100 bp PE	>10000 bp
Error Rate	0.001%	0.1%	2%	0.1%	2%	10-15%
Output data (per run)	100 KB	1 GB	100 GB	100 GB	1 TB	10 GB
Approx cost per GB		10,000	1000	100	10	1000

- **Solution:** assemble genome from short reads using computers. Whole Genome Shotgun Assembly.

Table from the course EE 372 taught by David Tse at Stanford

Current sequencing technologies

	Second gen. (Illumina)	Oxford Nanopore (MinIon)	PacBio
read length (bases)	100-500	10K-100K	10K-20K
error rates	< 1%	10-15%	10-15%
speed (time/base)	6 mins/base/strand	250 bases/s	3 bases/s
# of reads in parallel	10^9	2000	150K
throughput (total # of bases/s)	3M	500K	450K

Table from the course EE 372: Data Science for High-Throughput Sequencing.
taught by David Tse at Stanford



MinION, a palm-sized gene sequencer made by UK-based Oxford Nanopore Technologies

Short Reads assemble into Contigs

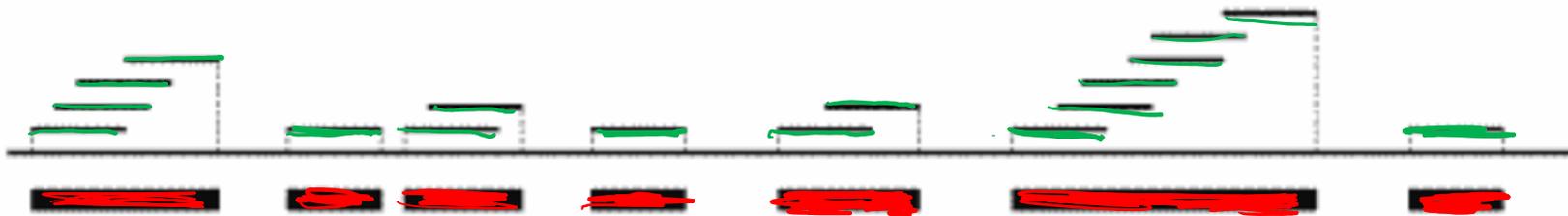
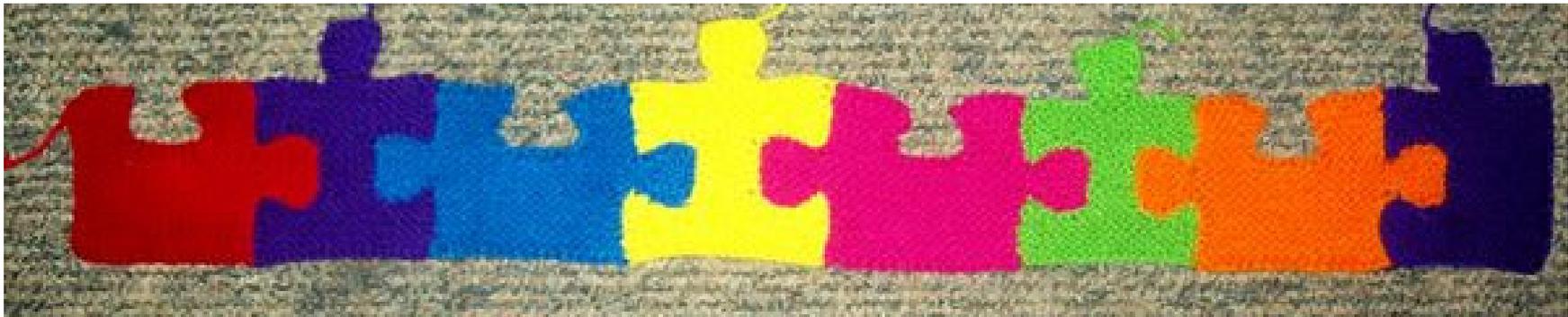


Figure 5.1.



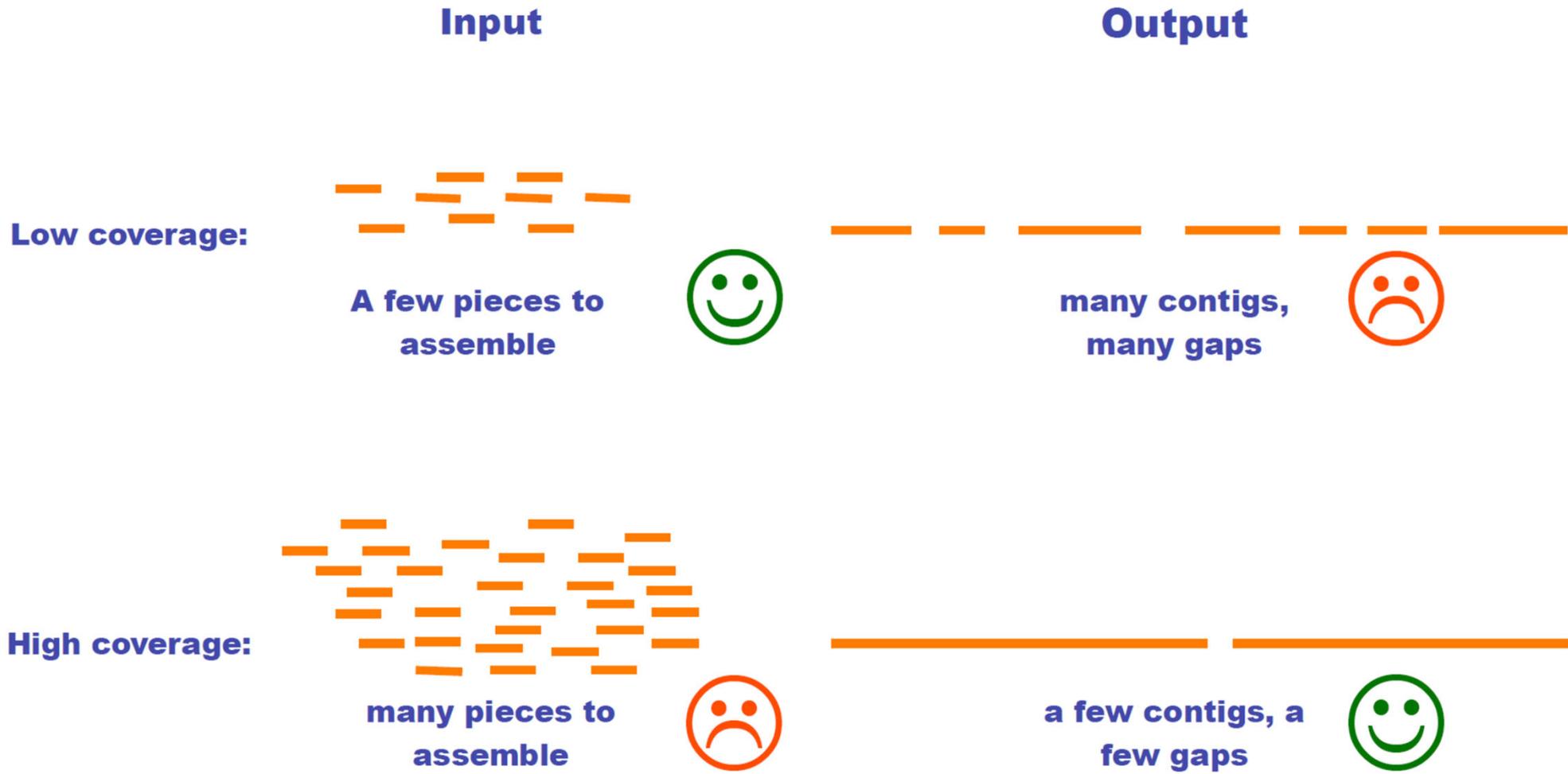
Promise of Genomics



Drew Sheneman, New Jersey -- The Newark Star Ledger, [E-mail Drew](#).

I think I found the corner piece!

How many short reads do we need?



Genome Assembly

Whole-genome “shotgun” sequencing starts by copying and fragmenting the DNA

(“Shotgun” refers to the random fragmentation of the whole genome; like it was fired from a shotgun)

Input: GCGTCTATATCTCGGCTCTAGGCCCTCATTTTTTT
35bp

Copy GCGTCTATATCTCGGCTCTAGGCCCTCATTTTTTT
by GCGTCTATATCTCGGCTCTAGGCCCTCATTTTTTT
PCR: GCGTCTATATCTCGGCTCTAGGCCCTCATTTTTTT
GCGTCTATATCTCGGCTCTAGGCCCTCATTTTTTT

Fragment: GCGTCTA TATCTCGG CTCTAGGCCCTC ATTTTTT
GGC GTCTATAT CTCGGCTCTAGGCCCTCA TTTTTT
GGCGTC TATATCT CGGCTCTAGGCCCT CATTTTTTT
GCGTCTAT ATCTCGGCTCTAG GCCCTCA TTTTTT

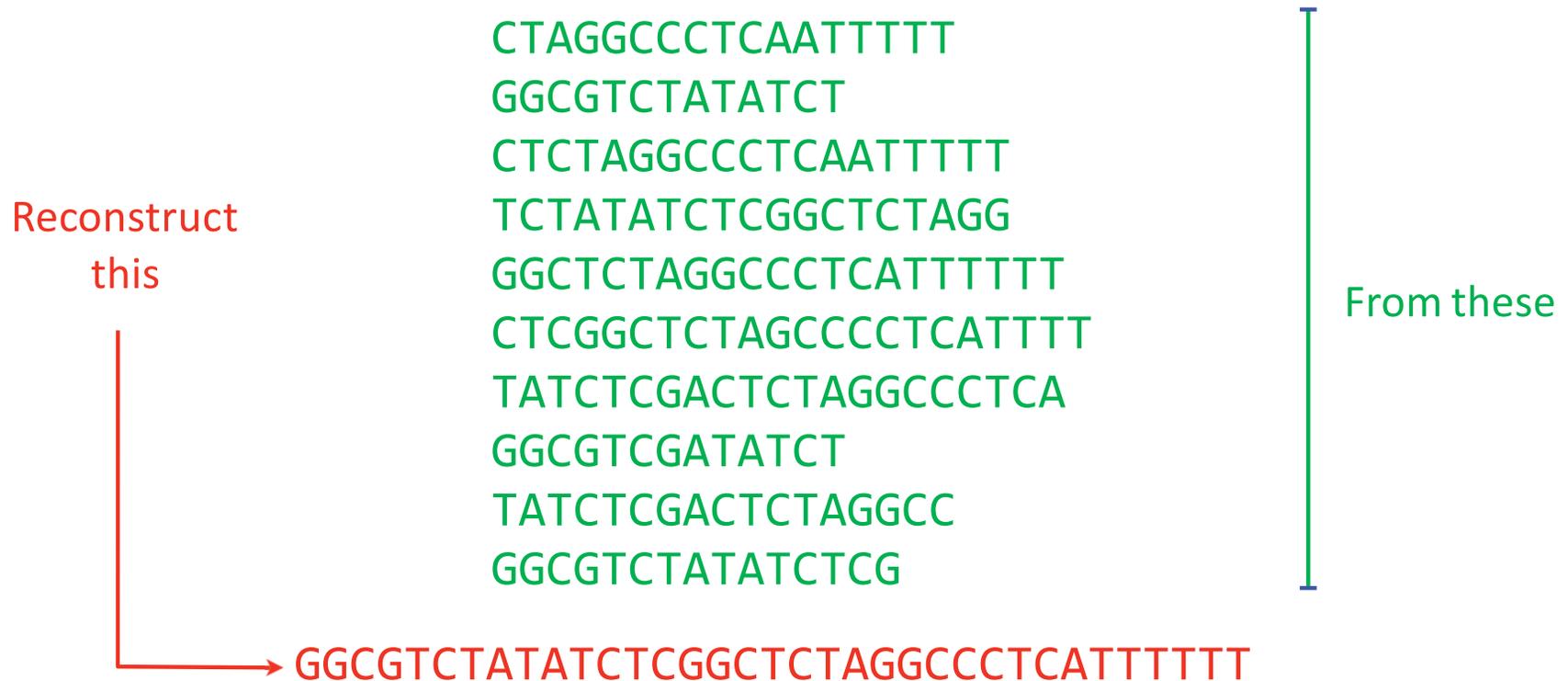
Courtesy of [Ben Langmead](http://www.langmead-lab.org). Used with permission.

<http://www.langmead-lab.org/teaching-materials/>

Assembly

Assume sequencing produces such a large # fragments that almost all genome positions are *covered* by many fragments...

...but we don't know what came from where



Courtesy of [Ben Langmead](http://www.langmead-lab.org/). Used with permission.

<http://www.langmead-lab.org/teaching-materials/>

Assembly

Overlaps between short reads help to put them together

```
          CTAGGCCCTCAATTTTT
         CTCTAGGCCCTCAATTTTT
        GGCTCTAGGCCCTCATTTTT
       CTCGGCTCTAGCCCCTCATTTT
      TATCTCGACTCTAGGCCCTCA
     TATCTCGACTCTAGGCC
    TCTATATCTCGGCTCTAGG
   GCGTCTATATCTCG
  GCGTCGATATCT
 GCGTCTATATCT
GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTT
```

177 nucleotides

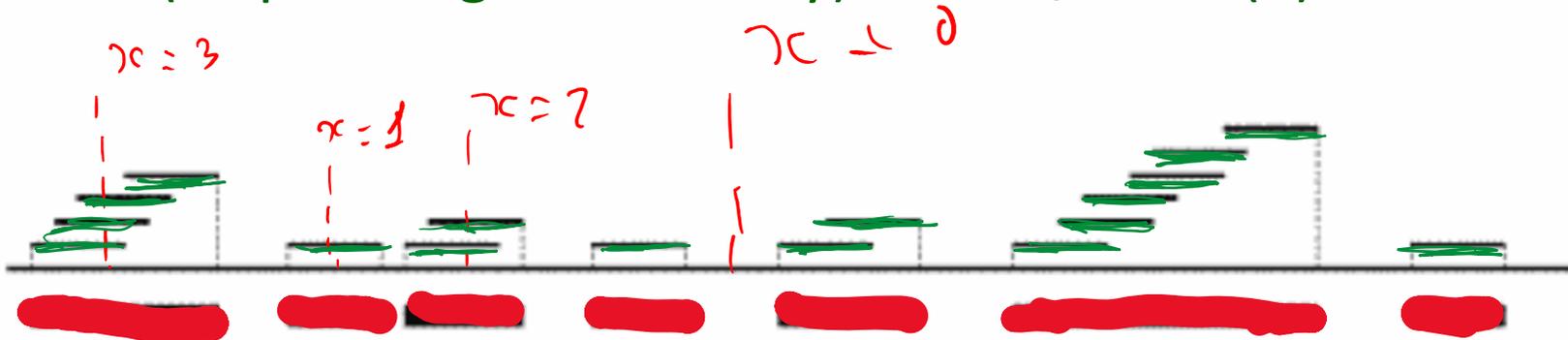
35 nucleotides

Where is the Poisson?

- G - genome length (in bp)
- L - short read average length
- N - number of short read sequenced
- λ - sequencing coverage redundancy = LN/G
- x - number of short reads covering a given site on the genome

$$P(x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

Poisson as a limit of Binomial: For a given site on the genome for each short read Prob(site covered): $p=L/G$ is very small. Number of attempts (short reads): N is very large. Their product (sequencing redundancy): $\lambda = NL/G$ is $O(1)$.



What fraction of the genome is missing?

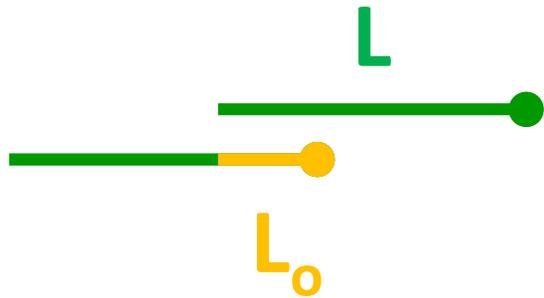
What fraction of genome is covered?

- Coverage: $\lambda = NL/G$,
X – random variable equal to the number of times a given site is covered by short reads.
Poisson: $P(X=x) = \lambda^x \exp(-\lambda) / x!$
 $P(X=0) = \exp(-\lambda)$, $P(X>0) = 1 - \exp(-\lambda)$
- Total length covered: $G * [1 - \exp(-\lambda)]$

λ	2	4	6	8	10	12
Mean proportion of genome covered	.864665	.981684	.997521	.999665	.999955	.999994

Table 5.1. The mean proportion of the genome covered for different values of λ

How long should the overlap be to connect two short reads?



If DNA was a random chain with $p_A = p_C = p_G = p_T = 1/4$

$L_0 \sim 16-20$ would be enough

$$2 \cdot G \cdot 4^{-L_0} = 2 \cdot 3 \times 10^9 \cdot 4^{-16} = 1.4$$

$$2 \cdot 3 \times 10^9 \cdot 4^{-20} = 0.0055 \ll 1$$