# Binomial Distribution

- **Binomially-distributed** random variable $X$ equals sum (number of successes) of n independent Bernoulli trials

- The probability mass function is:

$q = 1 - p$

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

- Based on the binomial expansion:

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

# 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

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

# 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=λ* is O(1)

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

| Number 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 $p = \dfrac{\lambda}{n}$

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

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

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

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

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

# Poisson Mean & Variance

If X is a Poisson random variable, then:

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

Note: Variance = Mean

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

# Matlab exercise: Poisson distribution

- Generate a sample of size 100,000 for Poisson-distributed random variable X with λ =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]= λ and V[X]=λ

# Matlab exercise: Poisson distribution

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

# QUESTIONS
## FOUND IN GOOGLE AUTOCOMPLETE

Credit: XKCD comics

# Poisson Distribution in Genome Assembly

Cost per Raw Megabase of DNA Sequence

Moore's Law

genome.gov/sequencingcosts

# Poisson Example: Genome Assembly

- Goal: DNA sequence of the entire genome of an organism
- Problem: Sequencers generate short reads of random portions of a genome
- Solution: assemble genome from short reads using computers
- Whole Genome Shogun Assembly pioneered by Craig Venter in 1990s
- The human genome was jointly announced in 2001 by the Human Genome Project (public) and Celera Genomics (Craig Venter's company)

# Short Reads assemble into Contigs



Figure 5.1.

# Current sequencing technologies

| Technology | Read Length | Error Rate | Cost per Gbase |
|---|---|---|---|
| Illumina NovaSeq | 75-500 bp | ~0.1% | $5-$150 |
| BGI DNBSEQ | 35-300 bp | ~0.1% | $5-$120 |
| Ion Torrent | 200-600 bp | ~0.5% | $70-$1000 |
| PacBio | 10,000-25,000 bp | 13% | $7-$40 |
| Oxford Nanopore | 10,000-100,000+ bp | 3-10% | $30-$60 |

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

# Promise of Genomics



I think I found the corner piece!

# How many short reads do we need?

**Input**

**Output**

**Low coverage:**

A few pieces to assemble 🙂

many contigs, many gaps ☹

**High coverage:**

many pieces to assemble ☹

a few contigs, a few gaps 🙂

# 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: GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
35bp

Copy GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
by   GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
PCR: GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT
     GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT

Fragment:  GGCGTCTA    TATCTCGG    CTCTAGGCCCTC    ATTTTTT
           GGC      GTCTATAT    CTCGGCTCTAGGCCCTCA    TTTTTT
           GGCGTC  TATATCT    CGGCTCTAGGCCCT    CATTTTTT
           GGCGTCTAT    ATCTCGGCTCTAG    GCCCTCA    TTTTTT

Courtesy of Ben Langmead. Used with permission.

# 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

Reconstruct this

CTAGGCCCTCAATTTTT
GGCGTCTATATCT
CTCTAGGCCCTCAATTTTT
TCTATATCTCGGCTCTAGG
GGCTCTAGGCCCTCATTTTTT
CTCGGCTCTAGCCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
GGCGTCGATATCT
TATCTCGACTCTAGGCC
GGCGTCTATATCTCG

From these

GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT

Courtesy of Ben Langmead. Used with permission.

# Assembly

Overlaps between short reads help to put them together

<pre>
              CTAGGCCCTCAATTTTT
             CTCTAGGCCCTCAATTTTT
            GGCTCTAGGCCCTCATTTTTT
           CTCGGCTCTAGCCCCTCATTTT
          TATCTCGACTCTAGGCCCTCA           177 nucleotides
          TATCTCGACTCTAGGCC
         TCTATATCTCGGCTCTAGG
       GGCGTCTATATCTCG
       GGCGTCGATATCT
       GGCGTCTATATCT
       GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTTT   35 nucleotides
</pre>

# 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): λ = NL/G is O(1).



Ewens, Grant, Chapter 5.1

# 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: $G*[1- exp(- \lambda)]$*

| $\lambda$ | 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 $\lambda$

# How long should be the length $L_{ov}$ of the overlap to connect two short reads into a contig?

**L**

**L$_o$**

**G**

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

$L_{ov} \sim 16\text{-}20$ would be enough

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

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

# How many contigs?



$$\text{P(short read can be extended by another short read)} = \frac{L - L_o}{G} = \text{p}$$

$$\text{P(short read cannot be extended by any short reads)} = e^{-pN} \approx Ne^{-\lambda}$$

$$\text{number of contigs} = Ne^{-pN} \approx Ne^{-\lambda}$$

# How many contigs?

- A given short read is the
  right end of a contig if and only if
  no left ends of other short reads fall within it.
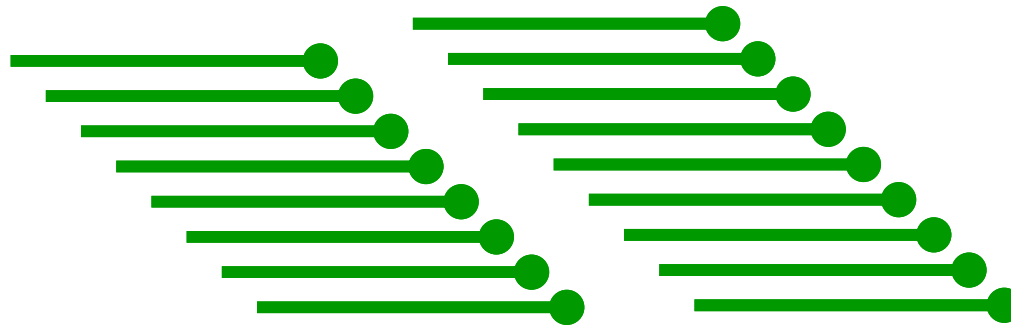- The left end of another short read has the probability
  $p=(L-1)/G$ to fall within a given read. There are
  $N-1$ other reads. Hence the expected number of left
  ends inside a given shot read is
  $p \cdot (N-1)=(N-1) \cdot (L-1)/G \approx \lambda$
- If significant overlap required to merge two short reads
  is $L_{ov}$, modified $\lambda$ is given by $(N-1) \cdot (L- L_{ov})/G$
- Probability that no left ends fall inside a short read is
  $exp(-\lambda)$. Thus the Number of contigs is $N_{contigs}=Ne^{-\lambda}$:

| $\lambda$ | 0.5 | 0.75 | 1 | 1.5 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|
| Mean number of contigs | 60.7 | 70.8 | 73.6 | 66.9 | 54.1 | 29.9 | 14.7 | 6.7 | 3.0 | 1.3 |

Table 5.2. The mean number of contigs for different levels of coverage, with
$G = 100,000$ and $L = 500$.

# Average length of a contig?

- Length of a genome covered:
$G_{covered}=G \cdot P(X>0)=G \cdot (1- exp(- \lambda))$

- *Number of contigs $N_{contigs}=N \cdot e^{-\lambda}$*

- Average length of a contig =

$<L>=\sum_i L_i/N_{contigs}=G_{covered}/N_{contigs}=$

$G \cdot (1- exp(- \lambda))/ N \cdot e^{-\lambda}=L \cdot (1- exp(- \lambda))/ \lambda \cdot e^{-\lambda}$

| $\lambda$ | 2 | 4 | 6 | 8 | 10 |
|---|---|---|---|---|---|
| Mean contig size | 1,600 | 6,700 | 33,500 | 186,000 | 1,100,000 |

Table 5.3. The mean contig size for different values of $a$ for the case $L = 500$.