# Lecture 4: Sequencing Coverage Analysis
## + Introduction to Genome-Wide Association Studies (GWAS)



ECE 365 - Data Science and Genomics

# Announcements:

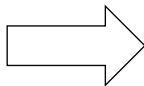- Lab 1 due today (11:59pm)
- Lab 2 released (due March 18)

# Today:

- Finish discussion on sequencing coverage
- Introduction to Genome-Wide Association Studies (GWAS)

# Indexing

16    ↗ read

cgtcagcggacagggc

```
ggtttaatgtggttctgcttggcggtagtcattaagagccccgtggtggccaat
caagaaaatgtcacgccgcttcccagcactttcagctgtttgtcgtagcccat
caccaccgtaagccaagacccagcttcaggccaagtagccttccgccagcggtt
ctgcgtcggcatggattctgcacggcaaagttcacgcgtcggtttgccataatt
aaggacgcgcctggattcaccttgcgatcggcaatcgcaggaatgagagagcag
ataatgaaagcgttgacgtaagaaagccatcgttttcccggtaccggtttttgc
gcctgcccggctacgtcagcgacctcgccagcgtcagcggacagggcgcaagtg
ccgtgaatgggccgtacagttatgaaacccttttttttctaaggggcttctacaa
cccttggatgcagggcgaagtcgggaaaacttctgttctgtttaaaatgtgttt
tgctcatagtgtggtagatctcagcttactattggctttaacgaaagccgtatt
ccggtgaaaataacagtcacgcttttagttgttaatgttacaccaacaacgaaa
ccaacacgccaggcttaattcctgtggagttatatatgagcgtaaatcggatcc
```

all substrings
of length 16

⟹

```
ggtttaatgtggttct   -->   0
gtttaatgtggttctg   -->   1
tttaatgtggttctgc   -->   2
ttaatgtggttctgct   -->   3
taatgtggttctgctt   -->   4
aatgtggttctgcttg   -->   5
atgtggttctgcttgg   -->   6
tgtggttctgcttggc   -->   7
gtggttctgcttggcg   -->   8
tggttctgcttggcgg   -->   9
ggttctgcttggcggt   -->   10
gttctgcttggcggta   -->   11
ttctgcttggcggtag   -->   12
tctgcttggcggtagt   -->   13
ctgcttggcggtagtc   -->   14
tgcttggcggtagtca   -->   15
gcttggcggtagtcat   -->   16
cttggcggtagtcatt   -->   17
ttggcggtagtcatta   -->   18
```

☐ Another way to do indexing: Python dictionary (hash function)

☐ Let's look at this in a notebook

# What if there are errors/mutations on read?

$X =$ **cgtaagcggacatggc**

$Y =$
```
ggtttaatgtggttctgcttggcggtagtcattaagagccccgtggtggccaat
caagaaaatgtcacgccgcttcccagcactttcagctgtttttgtcgtagcccat
caccaccgtaagccaagacccagcttcaggccaagtagccttccgccagcggtt
ctgcgtcggcatggattctgcacggcaaagttcacgcgtcggtttgccataatt
aaggacgcgcctggattcaccttgcgatcggcaatcgcaggaatgagagagcag
ataatgaaagcgttgacgtaagaaagccatcgtttttcccggtaccggtttttgc
gcctgcccggctacgtcagcgacctcgccagcgtcagcggacagggcgcaagtg
ccgtgaatgggccgtacagttatgaaaccctttttttctaaggggcttctacaa
cccttggatgcagggcgaagtcgggaaaacttctgttctgtttaaaatgtgttt
tgctcatagtgtggtagatctcagcttactattggctttaacgaaagccgtatt
ccggtgaaaataacagtcacgctttttagttgttaatgttacaccaacaacgaaa
ccaacacgccaggcttaattcctgtggagttatatatgagcgtaaatcggatcc
```

# What if there are errors/mutations on read?

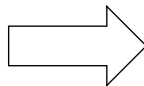$X =$ **cgtaagcggacatggc**

$Y =$
```
ggtttaatgtggttctgcttggcggtagtcattaagagccccgtggtggccaat
caagaaaatgtcacgccgcttcccagcactttcagctgtttgtcgtagcccat
caccaccgtaagccaagacccagcttcaggccaagtagccttccgccagcggtt
ctgcgtcggcatggattctgcacggcaaagttcacgcgtcggtttgccataatt
aaggacgcgcctggattcaccttgcgatcggcaatcgcaggaatgagagagcag
ataatgaaagcgttgacgtaagaaagccatcgtttttcccggtaccggtttttgc
gcctgcccggctacgtcagcgacctcgccagcgtcagcggacagggcgcaagtg
ccgtgaatgggccgtacagttatgaaaccctttttttctaaggggcttctacaa
cccttggatgcagggcgaagtcgggaaaacttctgttctgtttaaaatgtgttt
tgctcatagtgtggtagatctcagcttactattggctttaacgaaagccgtatt
ccggtgaaaataacagtcacgccttttagttgttaatgttacaccaacaacgaaa
ccaacacgccaggcttaattcctgtggagttatatatgagcgtaaatcggatcc
```

☐ One idea: Consider indexing substrings of length $k \approx 6$

# What if there are errors/mutations on read?

$X =$ **cgtaagcggacatggc**

$Y =$
```
ggtttaatgtggttctgcttggcggtagtcattaagagccccgtggtggccaat
caagaaaatgtcacgccgcttcccagcactttcagctgtttttgtcgtagcccat
caccaccgtaagccaagacccagcttcaggccaagtagccttccgccagcggtt
ctgcgtcggcatggattctgcacggcaaagttcacgcgtcggtttgccataatt
aaggacgcgcctggattcaccttgcgatcggcaatcgcaggaatgagagagcag
ataatgaaagcgttgacgtaagaaagccatcgtttttcccggtaccggtttttgc
gcctgcccggctacgtcagcgacctcgccagcgtcagcggacagggcgcaagtg
ccgtgaatgggccgtacagttatgaaacccttttttttctaaggggcttctacaa
cccttggatgcaggcgaagtcgggaaaacttctgttctgtttaaaatgtgttt
tgctcatagtgtggtagatctcagcttactattggctttaacgaaagccgtatt
ccggtgaaaataacagtcacgcttttagttgttaatgttacaccaacaacgaaa
ccaacacgccaggcttaattcctgtggagttatatatgagcgtaaatcggatcc
```

all substrings
of length $k = 6$

Python Dict ()
```
ggttta --> [0]
gtttaa --> [1, 471]
tttaat --> [2]
ttaatg --> [3, 571]
taatgt --> [4, 572]
aatgtg --> [5, 477]
atgtgg --> [6]
tgtggt --> [7, 495]
gtggtt --> [8]
tggttc --> [9]
ggttct --> [10, 158]
gttctg --> [11, 159, 466]
ttctgc --> [12, 160, 177]
tctgct --> [13]
ctgctt --> [14]
tgcttg --> [15]
gcttgg --> [16]
cttggc --> [17]
ttggcg --> [18]
tggcgg --> [19]
ggcggt --> [20]
```

☐ **One idea: Consider indexing substrings of length $k$**

# What if there are errors/mutations on read?

$X =$ **cgtaagcggacatggc**

take substrings
of length $k = 6$

```
cgtaag --> [114, 286]
gtaagc --> [115]
taagcg -->
aagcgg -->
agcgga --> [359]
gcggac --> [360]
cggaca --> [361]
ggacat -->
gacatg -->
acatgg -->
catggc -->
```

# What if there are errors/mutations on read?

$X =$ **cgt**a**agcgga**ca**t**ggc**

take substrings
of length $k = 6$

```
cgtaag --> [114, 286]
gtaagc --> [115]
taagcg -->
aagcgg -->
agcgga --> [359]
gcggac --> [360]
cggaca --> [361]
ggacat -->
gacatg -->
acatgg -->
catggc -->
```

Smith Waterman $(X, Y)$

genome

cgtaagcggacatggc
                1 x x 1 x 1 x x x x
...tagcccatcaccaccgtaagccaagacccagcttcaggccaagtagccttccgccagcgg...

114            Y
→ bad alignment

cgtaagcggacatggc
1 1 x x 1 1 1 1 1 x x 1 x 1
...ccggctacgtcagcgacctcgccagcgtcagcggacagggcgcaagtgccgtgaatgggc...
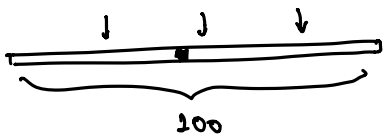
359   Y

smith Waterman $(X, Y)$
good alignment!

# How do we choose $k$?

$$\lceil 3.2 \rceil = 4$$
$$\lceil 5 \rceil = 5$$

□ Based on the sequencing technology (e.g., Illumina: $L = 100$, error rate = 0.1%)

· Suppose at most 1 error per read



100

Claim: there is a segment of length $\dfrac{100-1}{2} = 49.5$ with no errors

· Suppose $\leq t$ errors per read



there is a segment of length $\left\lceil \dfrac{L-t}{t+1} \right\rceil$ with no errors

$\hookrightarrow$ set $= k$

How many reads have $> t$ errors?

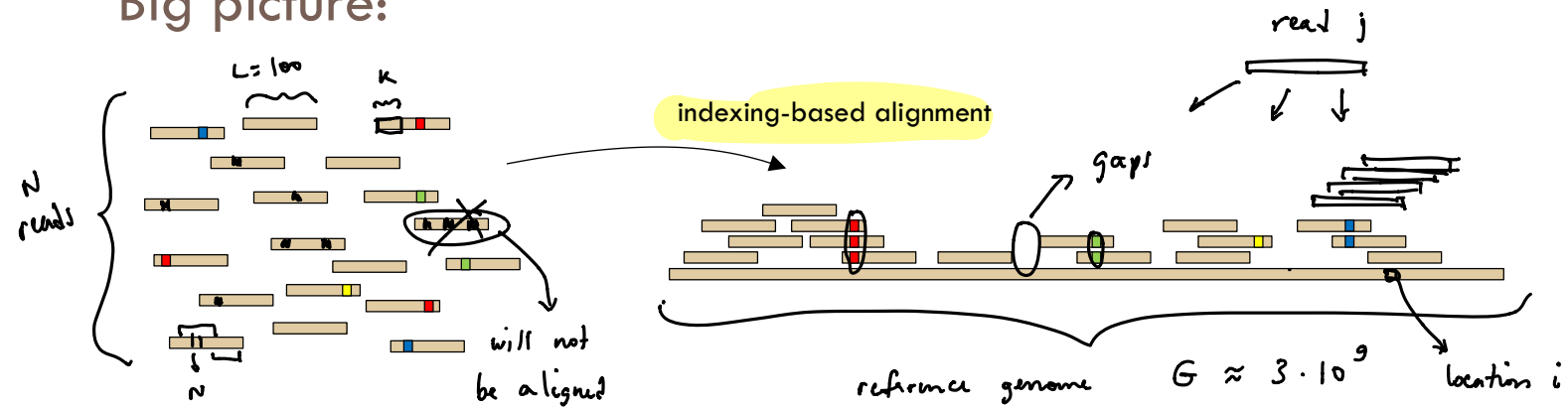# of errors in a read: $X \sim \text{Binomial}(100, 0.001)$

fraction of reads with $> t$ errors: $P(X > t)$

e.g. for $t = 2$.

$P(X > 2) = 0.00015$

$k = \left\lceil \dfrac{100-2}{3} \right\rceil = 33$

# Big picture:



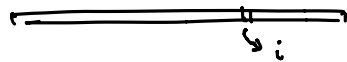How to pick $N$?

Assume that start location of each read is uniformly distributed on genome

$$P(\text{read } j \text{ covers position } i) = \frac{L}{G}$$

$\Rightarrow$ Expected # of reads covering position $i$ = $\frac{NL}{G} \triangleq c$ $\begin{pmatrix} \text{coverage} \\ \text{depth} \end{pmatrix}$

Prob. that there are gaps? Fix $i$.

$$P(i \text{ is uncovered}) = \prod_{j=1}^{N}\left(1 - \frac{L}{G}\right) = \left(1 - \frac{L}{G}\right)^{N} = \underbrace{\left(1 - \frac{1}{\frac{G}{L}}\right)^{\overbrace{\frac{G}{L} \cdot \frac{NL}{G}}^{c}}}_{\approx e^{-1}} \approx e^{-c}$$

Recall: $\displaystyle\lim_{n \to \infty}\left(1 - \frac{1}{n}\right)^{n} = e^{-1}$

$$P(\text{some position is uncovered}) \overset{\underset{\text{union bound (good approximation)}}{\downarrow}}{\leq} \; G e^{-c} = \varepsilon \quad (\text{desired failure prob.})$$

$$\Rightarrow \quad e^{c} = \frac{G}{\varepsilon} \quad \Rightarrow \quad c = \ln\left(\frac{G}{\varepsilon}\right)$$

Since $\quad c = \dfrac{NL}{G}, \qquad N = \dfrac{G}{L} \ln\left(\dfrac{G}{\varepsilon}\right). \quad \left(\text{Lander - Waterman}\right)$

# Big picture:

- Map sequencing data to a reference genome


indexing-based alignment

# Big picture:

- Map sequencing data to a reference genome



indexing-based alignment

variants

- Using aligned data, we can identify *variants*

# Big picture:

□ Map sequencing data to a reference genome



indexing-based alignment

□ Using aligned data, we can identify *variants*

□ Variants allow us to create genotype matrices

different variants

$$
\begin{array}{c}
\text{person 1} \\
\\
\\
\vdots \\
\\
\\
\text{person } n
\end{array}
\begin{bmatrix}
0 & 1 & 1 & 0 & 0 & 1 \\
1 & 0 & 0 & 1 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 1 \\
1 & 1 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 \\
1 & 1 & 0 & 1 & 0 & 1
\end{bmatrix}
$$

genotype

# Genotype data

- Focus on a set of common variants in the genome

# Genotype data

□ Focus on a set of common variants in the genome



SNP

not a SNP

□ Typically, we focus on Single-Nucleotide Polymorphisms (SNPs, read snips)

# Genotype data

- Focus on a set of common variants in the genome



- Typically, we focus on Single-Nucleotide Polymorphisms (SNPs, read snips)
- Most SNPs only allow two possible values: REF and ALT $\left( e.g. \;\; REF = A, \; ALT = C \right)$

# Genotype data

☐ Focus on a set of common variants in the genome



SNP

☐ Typically, we focus on Single-Nucleotide Polymorphisms (SNPs, read snips)

☐ Most SNPs only allow two possible values: REF and ALT

☐ More info on SNPs: SNPedia!

　　☐ Search for rs429358



**SNPedia**

💬 Talk　❓ Contributions　👤 Create account　🔑 Log in
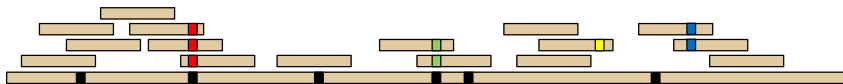
Search SNPedia 🔍

Navigation ▾

📄 Page　💬 Discussion　📝 Edit with form　📝 Edit　🕓 History

Have questions? Visit https://www.reddit.com/r/SNPedia

**Empower yourself**　MyHeritage DNA　Order now

### rs429358

This SNP, located in the fourth exon of the ApoE gene, affects the amino acid at position 130 of the resulting protein. The more common **rs429358** allele is (T). If the allele is (C) and the same chromosome also harbors the rs7412(C) allele, the combination is known as an APOE-ε4 allele. The APOE-ε4 allele has a strong influence on the risk of Alzheimer's disease.

Many studies have estimated the level of risk, and it varies depending on age, sex, ethnicity,

| Orientation | plus |
| Stabilized | plus |

| Geno ⬍ | Mag ⬍ | Summary ⬍ |
|---|---|---|
| (C;C) | 1.2 | one of 2 snps relevant to classifying APOE genotype |

# Don't forget: humans are diploid

*two copies of each chromosome*

- ☐ Two "copies" of each chromosome (different variants)
- ☐ Aligned data is a combination from both:

# Variant data: VCF files

- Standard file format to list variants of one or many individuals
- Let's take a look!

# Variant data: VCF files

- Standard file format to list variants of one or many individuals
- Let's take a look!
- Genotype matrix:

person 1                                    person n

$$
\begin{array}{l}
\text{SNP 1} \\
\text{SNP 2} \\
\\
\\
\\
\\
\\
\text{SNP m}
\end{array}
\begin{bmatrix}
0|0 & 1|0 & 1|0 & 0|0 & 0|0 & 0|1 \\
0|1 & 0|0 & 1|1 & 1|0 & 0|0 & 1|1 \\
0|1 & 0|1 & 0|1 & 0|0 & 0|0 & 0|0 \\
0|0 & 0|0 & 0|0 & 1|1 & 0|1 & 1|0 \\
1|0 & 0|1 & 1|1 & 0|0 & 0|0 & 0|0 \\
0|0 & 1|1 & 0|0 & 0|0 & 0|1 & 0|0 \\
0|1 & 0|1 & 0|0 & 1|0 & 0|0 & 1|1
\end{bmatrix}
$$

# Variant data: VCF files

- Standard file format to list variants of one or many individuals
- Let's take a look!
- Genotype matrix:

person 1         person n

|       | person 1 | | | | | person n |
|-------|---|---|---|---|---|---|
| SNP 1 | 0\|0 | 1\|0 | 1\|0 | 0\|0 | 0\|0 | 0\|1 |
| SNP 2 | 0\|1 | 0\|0 | 1\|1 | 1\|0 | 0\|0 | 1\|1 |
|       | 0\|1 | 0\|1 | 0\|1 | 0\|0 | 0\|0 | 0\|0 |
|       | 0\|0 | 0\|0 | 0\|0 | 1\|1 | 0\|1 | 1\|0 |
|       | 1\|0 | 0\|1 | 1\|1 | 0\|0 | 0\|0 | 0\|0 |
|       | 0\|0 | 1\|1 | 0\|0 | 0\|0 | 0\|1 | 0\|0 |
| SNP m | 0\|1 | 0\|1 | 0\|0 | 1\|0 | 0\|0 | 1\|1 |

*additive effect*

⟹ count ALT occurrences

| 0 | 1 | 1 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 1 | 0 | 2 | 1 | 0 | 2 |
| 1 | 1 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 2 | 1 | 1 |
| 1 | 1 | 1 | 0 | 0 | 0 |
| 0 | 2 | 0 | 0 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 | 2 |

# Genome-Wide Association Studies (GWAS)

$$
\begin{array}{c}
\phantom{\text{person 1}} \quad \text{SNP 1} \qquad\qquad\qquad \text{SNP m} \\
\begin{array}{r}
\text{person 1} \\ \\ \\ \\ \\ \\ \text{person n}
\end{array}
\left[
\begin{array}{cccccc}
0 & 1 & 1 & 0 & 0 & 1 \\
1 & 0 & 2 & 1 & 0 & 2 \\
1 & 1 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 2 & 1 & 1 \\
1 & 1 & 1 & 0 & 0 & 0 \\
0 & 2 & 0 & 0 & 1 & 0 \\
1 & 1 & 0 & 1 & 0 & 2
\end{array}
\right]
\end{array}
$$

# Genome-Wide Association Studies (GWAS)

|          | SNP 1 |   |   |   |   | SNP m |
|----------|-------|---|---|---|---|-------|
| person 1 | 0     | 1 | 1 | 0 | 0 | 1     |
|          | 1     | 0 | 2 | 1 | 0 | 2     |
|          | 1     | 1 | 1 | 0 | 0 | 0     |
|          | 0     | 0 | 0 | 2 | 1 | 1     |
|          | 1     | 1 | 1 | 0 | 0 | 0     |
|          | 0     | 2 | 0 | 0 | 1 | 0     |
| person n | 1     | 1 | 0 | 1 | 0 | 2     |

phenotype
$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

e.g. has Diabetes

(doesn't need to be binary. e.g., height)

genome-wide

# Genome-Wide Association Studies (GWAS)



associated

|           | SNP 1 | | SNP 200 | | SNP m | | | phenotype |
|-----------|---|---|---|---|---|---|---|---|
| person 1  | 0 | 1 | 1 | 0 | 0 | 1 | | 0 |
|           | 1 | 0 | 2 | 1 | 0 | 2 | | 0 |
|           | 1 | 1 | 1 | 0 | 0 | 0 | | 1 |
|           | 0 | 0 | 0 | 2 | 1 | 1 | | 1 |
|           | 1 | 1 | 1 | 0 | 0 | 0 | | 1 |
|           | 0 | 2 | 0 | 0 | 1 | 0 | | 0 |
| person n  | 1 | 1 | 0 | 1 | 0 | 2 | | 1 |
| new person | 1 | 0 | 1 | 1 | 2 | 0 | | ? (predict) |

☐ Which SNPs are associated with a given phenotype?

☐ Given a new individual's genotype, can you predict their phenotype?

# Revisiting Logistic Regression

☐ Predict binary variable from real-valued features

$$p(1 \mid \underline{x}) = \frac{e^{\beta_0 + \underline{\beta}^T \underline{x}}}{1 + e^{\beta_0 + \underline{\beta}^T \underline{x}}} = \frac{1}{1 + e^{-(\beta_0 + \underline{\beta}^T \underline{x})}}$$

$$\overset{''}{p}$$

$$1 + e^{-(\beta_0 + \underline{\beta}^T \underline{x})} = \frac{1}{p} \implies e^{-(\beta_0 + \underline{\beta}^T \underline{x})} = \frac{1}{p} - 1 = \frac{1-p}{p}$$

$$\longrightarrow \ln\left(\frac{p}{1-p}\right) = \underbrace{\beta_0 + \underline{\beta}^T \underline{x}}_{\text{linear model}}$$

$$\frac{p}{1-p} : \text{odds ratio}$$

$$\ln\frac{p}{1-p} : \text{log odds ratio}$$