

Spring 2025

CS 2840

Pranav Mahableshwarkar

Instructor: Sorin Istrail

LECTURE 1:

1/23/2025

Introduction

- Human genome contains $\sim 3\text{G}$ basepairs $\rightarrow 46$ chromosomes.
- 2 individuals are 99% the same \rightarrow difference = 10M basepairs
- **GWAS Task:** Finding, computationally, specific regions of shifts.

Definition: A SNP – single nucleotide polymorphism – a difference in one basepair occurs once every 600 bp.
Most SNPs are common \rightarrow in 2–3% of the population.

Genomic Foundations

We will consider GWAS, protein folding, Linkage Disequilibrium, and more in this course.

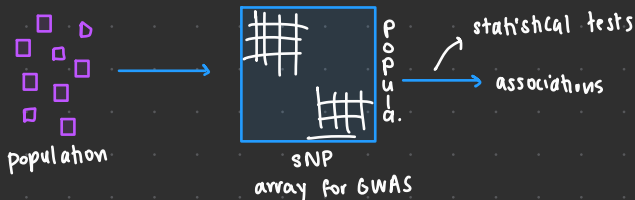
ALGORITHMS

- MLE and E-M Algorithms, covering topics of spectral graph theory, to examine the substructures of populations.
- ex: there are 2600+ subpopulations within in India
- protein folding/misfolding, Alpha Fold, \rightarrow applied to mad cow disease, for example

AN INTRO TO Genome Wide Association Studies (GWAS)

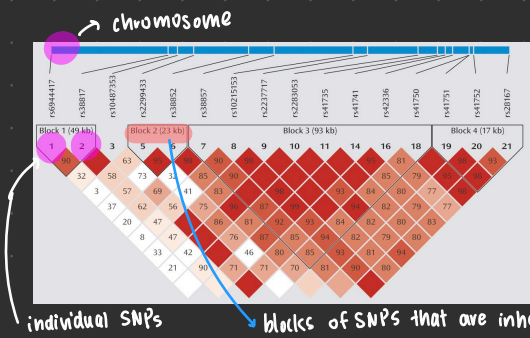
Before GWAS, disease context was random and non-rigorous.

- we now can have 100,000s of data points \rightarrow low cost of sequencing
- Can be applied to unrelated people over many years, allowing for more subtle detections



Defining GWAS

- method for interrogating all 10 million variable points
- **Infinite Allele Statement:** There can only ever be 4 letters at a position, every SNP is binary
this is because the genome is so large, multiple mutations @ 1 pt. uncommon



our SNPs are arranged in blocks (haplotype blocks) that allow us to group them

↓
further studied in linkage disequilibrium.

- After GWAS we can use SNPs to inform the risk for disease.
- We can also identify gene pathways from significant GWAS SNP hits.

What is Disease?

The naive thought is a broken protein causing a cascade of issues (monogenic). This is NOT the case for most complex diseases.

SNPs in the Genome

Definition: A position in the genome at which two or more diff bases in the genome occur, each with a frequency $\sim 1\%$

The most abundant type of polymorphism.

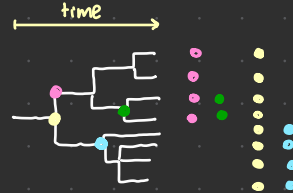
Haplotype: G C T C G A C A A C A G } the haplotype is the sequence of the SNPs: CAG

From this, haplotypes are our alleles that we can binarize due to the infinite sites assumption. → our matrix is entirely binary

Recombination: When new sequences are generated from 2 sequences: $P_1P_2, B_1B_2 \rightarrow P_1B_2, B_1P_2$

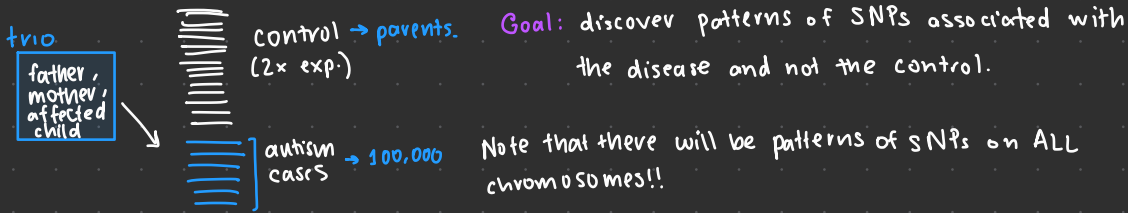
LINKAGE DISEQUILIBRIUM

- variations in chromosome populations over time
- easily visualized through phylogenetic trees



- **definition:** a way of statistically defining disease associations.

LECTURE 2: GWAS computational Pipeline



Chapter 1: Haplotype Phasing Problem

Every study participant → genotyped → SNP Array

Humans are diploid: mother chrom 1, 2, ..., 22
father chrom 1, 2, ..., 22

sex chromosome: female XX 23
male XY 24

SNP = single nucleotide polymorphism

↳ SNP has 2 alleles: the allele from the mother and the allele from the father

Definition: A haplotype is a single DNA sequence representing the SNPs that occur every ~500-600 bp or so

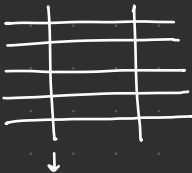
Say from the mother we have: $\overset{\text{SNP 1}}{\overbrace{A \ A \ T \ T}^{\text{in practice it is very challenging to construct the haplotypes because we only have the pairs, we don't know which parent the individual nucleotides originated from.}}}$
father we have: $\overbrace{G \ A \ G \ A}$

Haplotype Phasing comes down to inferring the mother father haplotypes.

★ WITH ONE INDIVIDUAL THIS IS IMPOSSIBLE ★

COMBINATORICS of HAPLOTYPES:

Definitions: a genotype is a string over $\{0, 1, 2\}$, a haplotype: is a string over $\{0, 1\}$



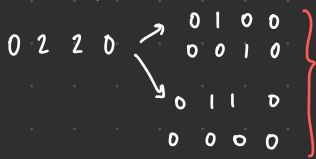
the genotype is made up of two haplotypes!

$$\text{two haplotypes} = \begin{Bmatrix} 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{Bmatrix} = 0 \ 1 \ 2 \ 0$$

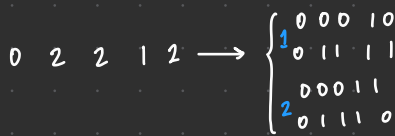
ambiguous, different.

this is a SNP if: there is high % of people that have it. Recall that due to the infinite sites model, only two nucleotides are possible at a given position. i.e. AA, GA, AG, or GG

How does phasing get complicated?



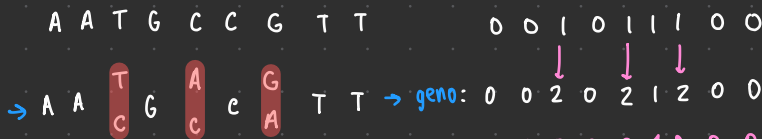
there are 2 explanations!, if there are 3 2s in the genotype we will have 2^2 possibilities.



3 ... \rightarrow this generalizes to 2^{n-1} explanations for n ambiguities

The Haplotype Phasing Problem

- Input: g_1, g_2, \dots, g_n
- Output: $h_1, h_2, \dots, h_n, h_{n+1}$ \rightarrow all the haplotypes for all of the given genomes.



0 0 0 0 0 1 0 0 0, the Clark Rule haplotype
 \downarrow
 now we have two haplotypes

Where do we start haplotyping?

When you have unambiguous genome snippets or single ambiguity in the snippet, you have good start pts

- $g_1 = \text{ATTCT} \begin{cases} \text{ATTCT} \\ \text{ATTCT} \end{cases}$
- $g_2 = \text{AT} \begin{matrix} G \\ A \end{matrix} \text{CT} \begin{cases} \text{ATGCT} \\ \text{ATACT} \end{cases}$

3 ISSUES IN THE CLARK REGIME

- No full homozygote or single SNP
- We arrive at a point where there are no new updates we can make, and there are no remaining candidates. The remaining genomes are called **orphans**.
- We find incorrect haplotypes. This rule is heuristic and does not guarantee exceptions. These errors are called **anomalous haplotypes**.

$$g_1 = 0 \ 1 \ 2 \ 2 \ 2$$

$$g_2 = 0 \ 0 \ 2 \ 1 \ 2$$

$$\begin{array}{l}
 g_3 = 0 \ 0 \ 0 \ 0 \ 0 \\
 g_4 = 0 \ 2 \ 1 \ 1 \ 1
 \end{array}
 \begin{array}{l}
 \leftarrow h_1: 0 \ 0 \ 0 \ 0 \ 0 \\
 h_1: 0 \ 0 \ 0 \ 0 \ 0 \\
 h_2: 0 \ 0 \ 1 \ 1 \ 1 \\
 h_3: 0 \ 1 \ 1 \ 1 \ 1
 \end{array}$$

these are automatically phased.

$$\text{Applying } h_3 \rightarrow g_1 \rightarrow \begin{array}{c} 0 \ 1 \ 1 \ 1 \ 1 \\ 0 \ 1 \ 2 \ 2 \ 2 \end{array} \rightarrow h_4: 0 \ 1 \ 0 \ 0 \ 0$$

$$\text{Applying } h_2 \rightarrow g_2 \rightarrow \begin{array}{c} 0 \ 0 \ 1 \ 1 \ 1 \\ 0 \ 0 \ 2 \ 1 \ 2 \end{array} \rightarrow h_5: 0 \ 0 \ 0 \ 1 \ 0$$

now we have explained all of the genotypes using Clark's Rule.

LECTURE 3: The Haplotype Phasing Problem

- Chapter 1 Outline
- 1.1 Clark Algorithm (greedy)
 - 1.2 The E-M Algorithm (ML, maximum likelihood)
 - 1.3 The Parsimony Algorithm
 - 1.4 The ML-Phasing and Parsimony Phasing

Formalizing the CLARK PROCEDURE

- 1 Identify all homozygotes and single hetero-zygotes (Aa, aA) haplotypes. Phase them and call these **resolved genotypes**. [See Lecture 2 Notes for examples]
- 2 Determine whether any of the unsolved haplotypes can be used to solve any of the remaining ambiguous genotypes using the **Clark Rule**.
- 3 Iteratively continue until we have resolved all genotypes.

3 ISSUES: See Lecture 2 Notes for these enumerated issues

E-M Algorithm (1.2)

Goal: Finding the haplotype frequencies in a population and the maximum likelihood haplotype phasing.

Define: The **E-M Algorithm** is an **iterative algorithm** to compute successive sets of haplotype frequencies.

- 1 p_1, p_2, \dots, p_T starting on initial $p_1^0, p_2^0, \dots, p_T^0$ \rightarrow 2ⁿ haplotypes, there are a lot!
- 2 These initial values are used as if they are unknown true frequencies to estimate the **explanation frequencies** $\Rightarrow p(h_k | g_i)^0$, expectation step
 \hookrightarrow this is a genome explanation for a genotype.
- 3 The expected explanation frequencies are used in turn to estimate haplotype frequencies at the next step - maximization step $\rightarrow p_1^1, p_2^1, \dots, p_T^1$
continue until convergence \rightarrow when consecutive steps yield minimal changes.

CONSTRUCTION + ALGORITHM PROCESS

definition: a **genotype** is a multi-locus genotype whose multi-locus haplotype phase is unknown
an **explanation** is a particular combination of 2 haplotypes, explaining a genotype

In E-M, we ultimately are trying to define a model that **maximizes the probability** or **likelihood** of observing the given data (**genotypes**)

Likelihood Function

Ultimately, we can/will only estimate frequencies over known/seen phenotypes due to computational limitations.

ASSUMPTION: The distribution of the given sample is **multinomial** w.r.t latent and unknown frequencies P_1, P_2, \dots

$$\text{Thus } IP(\text{sample} \mid P_1, \dots, P_m) = \frac{n!}{n_1! n_2! \dots n_m!} \times \prod_{i=1}^m P_i^{n_i}$$

num individuals m

Here, n_i : represent counts for the m unique genotypes.

The number of **explanations** (C_j) for S_j heterozygous **loci** can be defined as follows:

note that this **UNIQUE** to 1 genotype.

$$C_j = 2^{S_j - 1} \rightarrow \text{this is clearly seen above in our Clark's Rule examples (see above)}$$

Putting this together, P_j : the probability of the j^{th} **genotype** is the following. Let us also add the following notation: $H_j = \{h_k, h_\ell\}$ where $|H_j| = C_j$ is the set of all explan. for a given phenotype.

$$P_j = \sum_{i=1}^{C_j} IP(\text{explanation } i) = \sum_{h_k, h_\ell \in H_j} IP(h_k h_\ell) \rightarrow \text{recall that 2 haplotypes make up every explanation.}$$

Here: $IP(h_k h_\ell) = p_k^2$ if $k = \ell$, $2p_k p_\ell$ if $k \neq \ell$.

Putting everything together:

$$\text{Likelihood} = L(p_1, p_2, \dots, p_h) = \underbrace{a, \prod_{j=1}^m \left(\sum_{h_k, h_\ell \in H_j} IP(h_k h_\ell) \right)}_{\substack{\text{haplotype freqs} \\ \text{sum to 1.}}} \quad \begin{matrix} \nearrow n_j \text{ number of occ. of genotype } j \\ \searrow \text{multinomial constant} \end{matrix}$$

How do we optimize this?

Normally, with Maximum Likelihood Estimation, we take partials w.r.t our model, set it to 0 and find our optima. Here, this would result in $h-1$ equations:

$$U_t = \text{score of the } t^{\text{th}} \text{ haplotype} = \frac{\partial \log L}{\partial p_t} \xrightarrow{\text{log for numerical stability}} = \sum_{j=1}^m \frac{n_j}{P_j} \frac{\partial P_j}{\partial p_t} = \sum_{j=1}^m \frac{n_j}{P_j} \sum_{h_\ell \in H_j} (2 - \mathbb{1}_{\ell=t}) p_\ell$$

Setting these all to 0 and solving would be very tedious. Thus, enter the

E-M Algorithm!

E-M Algorithm

The goal, as aforementioned, is to compute successive sets of haplotype frequencies p_1, \dots, p_m .

- These frequencies are used to estimate **explanation** frequencies [E-STEP]
- Then, these are used to update haplotype frequencies for the next iteration, **M-STEP**

First, we have to initialize the haplotype frequencies: $p_1^{(0)}, p_2^{(0)}, \dots, p_m^{(0)}$

- This could be uniformly distributed: $IP(\text{expl. } h_k h_e \text{ for } g_j) = P_j(h_k h_e)^{(0)} = \frac{1}{c_j}$
- Other initial conditions include $p_i = p_j$ for all i, j (equal IP over haplotypes) which represents **complete linkage equilibrium**
- haplotype frequencies chosen at random.

For these notes, we take the first approach.

Expectation Step

At the t^{th} iteration, we use the previous iterations' h frequencies to determine the IP of resolving each genotype into possible explanations.

$$IP_j(h_k h_e)^{(t)} = \frac{n_j}{n} \frac{IP(h_k h_e)^{(t)}}{IP_j(\cdot)} \quad \begin{array}{l} \text{interpretation: weighted cond.} \\ \text{probability} \\ \text{function of } h_k, h_e \dots \end{array}$$

We can understand this

as calculating the expected frequency of every explanation for each genotype.

Maximization Step

We then use the gene counting method for optimization.

$$p_i^{(t+1)} = \frac{1}{2} \sum_{j=1}^m \sum_{\substack{h_k h_e \in H_j \\ \text{z iterable}}} P_j(h_k h_e)^{(t)} \delta_{zi}$$

Here, $\delta_{zi} \in \{0, 1, 2\}$ represents the no. of times haplotype i is present in explan z

★ derivation of this will appear in HW1 ★

This derivation will later appear in the notes.

DERIVATION OF M-STEP USING LAGRANGE MULT.

Recall: Classical Expectation Maximization

We have a statistical model that generates data \mathbf{X} , a set of unobserved data \mathbf{Z} , and a set of unknown or missing parameters θ .

These yield the likelihood function: $L(\theta; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z} | \theta)$. The maximum likelihood estimate is determined by maximizing the probability of observing the given data.

$\max_{\theta} L(\theta; \mathbf{X}) = p(\mathbf{X} | \theta) = \int p(\mathbf{X}, \mathbf{Z} | \theta) d\mathbf{z} \rightarrow$ or by taking partial derivatives and using analytical optimization methods

Expectation Step:

Define $Q(\theta | \theta^{(t)}) = \mathbb{E} [\log \text{likelihood function}]$ w.r.t \mathbf{Z} given $\mathbf{X}, \theta^{(t)}$

- $Q(\theta | \theta^{(t)}) = \mathbb{E}_{\mathbf{Z} \sim p(\cdot | \mathbf{X}, \theta^{(t)})} [\log p(\mathbf{X}, \mathbf{Z} | \theta)]$

think of this as your expected likelihood over all possible latent variables \mathbf{Z}

Maximization Step:

$\theta^{(t+1)} = \arg\max_{\theta} Q(\theta | \theta^{(t)}) \rightarrow$ NOTE: this step is often performed by analytically computing the maximum value via Lagrange multipliers.

APPLICATION TO HAPLOTYPE PHASING

Here, $\mathbf{X} \rightarrow$ our observed genotypes $\mathbf{Z} \rightarrow$ the explanations for our genotypes (unobserved)

$\theta \rightarrow p_0, \dots, p_T$ for T haplotypes

Recall our likelihood function:

$$L(p_1, \dots, p_T) = a_i \prod_{j=1}^m \left(\sum_{h_k, h_e \in H_j} P(h_k h_e) \right)^{n_j}$$

But what does this really represent?

LECTURE 4:

02/04/2025

For the complete worked through EM example, please see the ^{slideshows} course website.

Example:

Say we only have genotypes: 22, 02, 20, 00, 11

↳ notice that there are missing genotypes → this is okay!!! You rarely observe all possible g.

This means the only possible haplotypes are: 00, 01, 10, 11

↳ let us denote these as $\theta_{00}, \theta_{01}, \theta_{10}, \theta_{11}$ → which will be iteratively determined in E-M.

↳ there are n_A genotypes of 22, n_B for 02, etc.

Now, for group A, for example:

$$\begin{aligned}
 IP(Y_A)^{(t)} = P_A^{(t)} &= \sum_{h_k, h_e \in H_A} IP(h_k h_e) = \sum_{h_k, h_e \in H_A} (2 - \delta_{ke}) p_k p_e \\
 &= 2\theta_{00}^{(t)}\theta_{11}^{(t)} + 2\theta_{01}^{(t)}\theta_{10}^{(t)}
 \end{aligned}$$

used 11 in Lecture 3.

Diagram: 22 branches to 00/11 and 01/10

Now in the E-Step, we want to calculate the Expected Number of Each Haplotypes

$$\bullet n_{00}^{(t+1)} = n_A P\left(\frac{00}{11} \mid Y_A\right) + n_B + n_C + 2n_D + 0n_E$$

$$\frac{2\theta_{00}^{(t)}\theta_{11}^{(t)}}{P_A}$$

$$\bullet \theta_{00}^{(t+1)} = \frac{n_{00}^{(t+1)}}{2n} \quad \text{where } 2n = \text{total number of haplotypes}$$

1.3 MAXIMUM LIKELIHOOD PHASING

Define Haplotypes: Sequences over $\{0, 1\}$ Genotypes: $\{0, 1, 2\}$

Let us consider the Likelihood Function:

$$L(p_1, p_2, \dots, p_T) = (p_1 p_2 + p_3 p_4)(p_1 p_5 + p_4 p_6)(p_1 p_7)(p_7 p_8)$$

T = total number of haplotypes, example over 4 genotypes.

We want to find the $p_{1:8}^{opt}$ that maximizes the Likelihood function. This is a computationally infeasible problem.

Parsimony: the smallest number of haplotypes involved. Note that in practice finding these values is done by simplifying the polynomial optimization.

LECTURE 5

02/06/2025

TWO OPTIMIZATIONS (that we can exactly solve):

1. We have $x_1, x_2, \dots, x_r \in \mathbb{R}, > 0$, $\sum_{i=1}^r x_i = \text{Some constant}$, if x_i are probabilities \hookrightarrow sum to 1.

Problem: the solution maximizing the product $P = \prod x_i$

Solution: optimal $\rightarrow x_1 = x_2 = \dots = x_r$

This used in the idea of equal likelihood \rightarrow if we wanted to maximize.

2. Find the solution to: $\max(p' = x_1^{n_1} x_2^{n_2} \dots x_r^{n_r}) \rightarrow$ optimal when $\frac{x_1}{n_1} = \frac{x_2}{n_2} = \dots = \frac{x_r}{n_r}$

This is used in renormalization within some E-M applications... see more in CS1820!

NP-COMPLETENESS PROOF for HP/PARSIMONY

Theorem: parsimony haplotype phasing is NP-hard and maximum-likelihood hap phasing is NP-hard

SNPs are biallelic $\rightarrow \{0, 1\}$: Hubbel denotes genomes differently \rightarrow

h_1 :	0	1	0
h_2 :	0	1	1

$\xrightarrow{\oplus} g: 021$

where 1 is now the ambiguous symbol.

PARSIMONY PHASING PROBLEM:

Given $G = \{g_1, g_2, \dots, g_n\}$ of observed genotypes, find the minimal size of inferred haplotypes

$H = \{h_1, h_2, \dots, h_k\}$ such that every genotype can be represented as a sum of two haplotypes.

We show that every instance of a general graph R may be connected into a set of genotypes so that a minimum set of inferred haplotypes corresponds to a minimum clique partition of R . \leftarrow this is an NP-hard problem!

What is a Clique? $R_1 \cdot \quad R_2 \text{ --- } R_3 \triangle \quad R_4$ 

A graph where every node is connected to every node.

Clique Partition? A partition of a graph into cliques that cover the entire graph. Minimizing the number of cliques is NP-hard. If we can reduce our problem to this, then our parsimony phasing must also be NP-hard

REDUCTION

First we must generate the graph $Z = (V, E) \Leftrightarrow G = \{g_1, g_2, \dots, g_n\}$ the genotype set

Algorithm: Given genotypes of length $2n$.

Notation: $g_{ij} = g_i[j]$ the j th coordinate of the i th genotype.

- We will have $\bigvee_{i=1:n}$ vertices,
- $i = 1:n$: Then, $g_{ij} = 2$ if $i = j$, $g_{ij} = 1$ if V_i is connected by an edge to V_j
- $i = n+1:2n$: Create an identity matrix where
 - $g_{ij} = 1$ if $j = n+i$, 0 otherwise, $1 \leq i \leq n$, $n+1 \leq j \leq 2n$

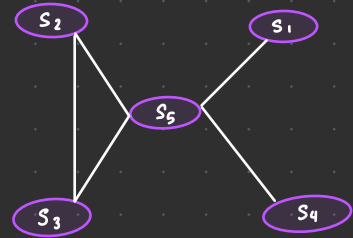
Example:

Vertex	Genotype	identity matrix
s_1	$g(s_1) = 20001$	1 0 0 0 0
s_2	$g(s_2) = 02101$	0 1 0 0 0
s_3	$g(s_3) = 01201$	0 0 1 0 0
s_4	$g(s_4) = 00021$	0 0 0 1 0
s_5	$g(s_5) = 11112$	0 0 0 0 1

these are the original genotypes and define edges

these ensure uniqueness for each $g(s_i)$

★ GRAPH CAME FIRST ★



Now let us define the haplotype set $H = h_1, h_2, \dots, h_T$ where we can write that $g_{i,n+i} = h_{e,n+i} + h_{m,n+i} \rightarrow$ recall that $n+i \rightarrow 2n$ is our identity matrix.

- Then, $h_{e,n+i}$ or $h_{m,n+i} = 1$ ★ BUT NOT BOTH ★

WLOG (without loss of generality) we assume that $h_{m,n+i} = 1 \rightarrow$ thus $h_{m,i}$ for $i=1:n$ must ALL BE UNIQUE!!!

ID	H	identity
h_{m1}	20001	1 0 0 0 0
h_{m2}	02101	0 1 0 0 0
h_{m3}	01201	0 0 1 0 0
h_{m4}	00021	0 0 0 1 0
h_{m5}	11112	0 0 0 0 1
h_{e1}	20001	0
h_{e2}	02101	0
h_{e3}	01201	0
h_{e4}	00021	0
h_{e5}	11112	0

the 2 blocks are the same

Let $h_{e,i}$ be the haplotype consistent with Q genotypes $\{g_{s1}, \dots, g_{s_n}\}$. We want to show that $\{s_1, \dots, s_n\}$ is a clique in \mathcal{Z} .

Recall that we can explain g_i with $g_i = h_{e,i} + h_{m,i}$ ↑ unique

For example, let $D = \{s_2, s_3, s_5\}$ for some $h_{e,i}$.

For $g_{i,i} = 2$ for $i = 1 \leq i \leq n$, which implies that $h_{e,i} = h_{m,i} = 1$. Thus, we must have $h_{e,k} = 1$ for $k = 1, \dots, n$ to satisfy the connections (see example below).

$g_{i,j} = 1$ occurs on all diagonals for $1 \leq j \leq n$ since we need 2 on the diagonal.

otherwise $g_{i,j} = 1$ means the nodes s_i, s_j are connected, 0 if not.

Vertex

Genotype

Haplotypes

s_1

$g(s_1) = 20001 \ 10000$

s_2

$g(s_2) = 02101 \ 61000$

s_3

$g(s_3) = 01201 \ 00100$

s_4

$g(s_4) = 00021 \ 00010$

s_5

$g(s_5) = 11112 \ 00001$

hm_1

10000

10000

hm_2

61000

61000

hm_3

00100

00100

hm_4

00010

00010

hm_5

00001

00001

he_1

10001

00000

he_2

01101

00000

he_3

01101

00000

he_4

00011

00000

he_5

11111

00000

notice that
these add to g_i .

WRAPPING UP THE PROOF

Therefore the set of genotypes such that $g_i, j = 1$ for all $i \neq j$ corresponds to a set of vertices that are a **CLIQUE**!

s_2, s_3, s_5 : 02101 , 01201 , $11112 \rightarrow$ CLIQUE!

Thus, the cliques generated in this graph each have their own common haplotype! Thus the number of haplotypes is a factor of the number of edges.

★ So to minimize the number of haplotypes, we must minimize the number of cliques! ★

This completes the reduction! PPHP is NP-hard.

LECTURE 6

2/11/2025

★ SEE COURSE WEBSITE ★ FOR SLIDES ★

LECTURE 7

2/13/2025

We need to explore and understand population genetics and linkage disequilibrium.

- Mutations / Recombinations } cause genetic variation in a population.
- Random mating
- genotype \rightarrow phenotype

LINKAGE DISEQUILIBRIUM / EQUILIBRIUM

Random Mating allows us to consider \rightarrow recall that humans are diploid \rightarrow ~ 2 alleles per gene

HARDY-WEINBERG EQUILIBRIUM

- With random mating, the alleles of \forall gene are combined at random according to HW-principles
- Given genotypes: A_1A_1 , A_1A_2 , A_2A_2

If $IP(A_1) = p_1$, $IP(A_2) = p_2$, $1 = p_1 + p_2$ and then HWE: A_1A_1 is p_1^2 , A_1A_2 is $2p_1p_2$, A_2A_2 is p_2^2 .

"Random Association" = the frequency of a gamete carrying any particular combination of alleles equals the products of the frequency of the alleles.

Definitions:

- Genes that are in random association are in linkage equilibrium. When they are not they are in linkage disequilibrium.
- With random mating and simplifying assumptions \rightarrow i.e. no mutations, \rightarrow + no selection, no migration, large population size \rightarrow we converge to linkage equilibrium

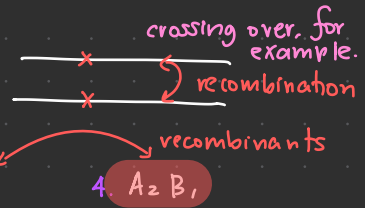
The rate of approach to LE depends on the rate of recombination. in genotypes heterozygous for both genes. There are 2 types of double heterozygotes

- A_1B_1 / A_2B_2
- A_1B_2 / A_2B_1

SOME GENETICS

Given A_1B_1 / A_2B_2 , there are four possible genotypes

1. A_1B_1
2. A_2B_2
3. A_1B_2
4. A_2B_1



By Mendelian Segregation, the frequency of genomic type 1 = the frequency of genomic type 2, " " 3 = " " type 4. because genomic recombination rates must imply type 3 = type 4, vice versa.

Definition: The recombination fraction is proportional to the number of recombinational gametes (type 3 + type 4) produced by a double heterozygote.

$0 \leq r \leq 0.5 \rightarrow 0.5$ means different genes or VERY far apart.

Genes where $r < 0.5$ are linked.

Altogether: type III = IV = $\frac{r}{2}$

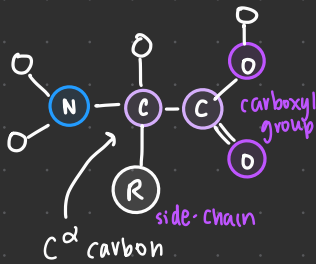
LECTURE 8: PROTEIN FOLDING

2/20/2025

Generally speaking, for a protein of length 50, there are an exponentially large number of ways to fold a protein. So how do we predict native structure?

PRIMER on PROTEINS

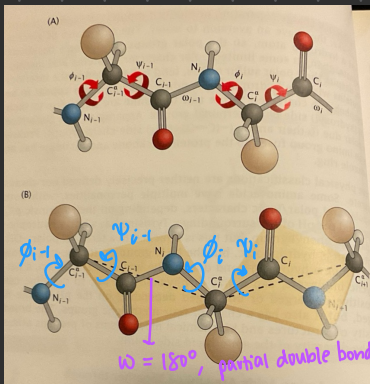
Proteins are polymers, sequences of individual peptides (of which there are 20).



proteins are either **L-isomer** or **D-isomer**, where L is more common

L-isomer is when functional group points away if N-C-C were shown to be on the same plane.

Amino acids are marked by C-N peptide bonds \rightarrow **N-terminus** (1) to **C-terminus** (n^{th}) these bonds form planes.



torsion angles ϕ and ψ . rigidity in backbone planes between C^{α} 's \rightarrow this is due to the covalent bond nature within N-C-C

\downarrow
"succession of planes that rotate about the planes between C^{α} indices. See (B) to the left.

Ramachandran Plots

The torsional angles vary \rightarrow depending on angles we can classify as α -helices, β -strands due to the types of angles.

$\nabla \rightarrow$ smaller \angle 's tend to be α , \searrow , more obtuse reflect β -sheets.

Self Avoiding Walks / Contacts: \rightarrow protein backbones are self-avoiding walks.

Models must not overlap while also considering the pairwise energies that result from interactions between peptides and even smaller, the component atoms.

Note: We want to minimize the energy of a system \rightarrow this is when it is the most stable!

HP-model + Contact maps

A simple way to visualize this is the HP-model.

\rightarrow we binarize all of the amino acids to be either H (hydrophobic), P (hydrophilic)

\rightarrow we want to maximize the number of H-H contacts to minimize energy

example:



We can generalize this to using **CONTACT MAPS!**

We create graphical/matrix representations connecting and relating the individual amino acids in a particular protein sequence.

PIPELINE: Structure Similarity \rightarrow Structure Alignment \rightarrow Fold Recognition \rightarrow Fold Alignment

Measuring Protein Similarity

How do we understand similarity for proteins of different sizes?

- RMSD \rightarrow root-mean square distance
 - Difference of distance matrices
 - Contact map overlap
 - Ad hoc scoring schemes
- often done by choosing 100 aa regions for fixing length

Skipping Colin and Pranav's ★ PERFECT POWERPOINTS ★ ...

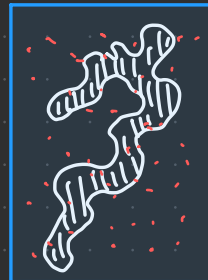
LECTURE 11 Markov Chain Monte Carlo Methods

3/04/25

3.1 Markov Chains + Markov Chain Monte Carlo (MCMC) Introduction



How do we compute the area of such an irregular shape?



Box the shape in a figure of area 1.

If we sample random pts, then by determining the ratio:

$$\frac{\text{pts in}}{\text{pts out}}$$

We \approx the area!!

Define: A Markov Chain has states $\{E_1, E_2, \dots, E_s\}$ with timing $t = 1, 2, \dots$

We also define a transition probability matrix indexed by the states:

$$\begin{array}{c} E_0 \\ \vdots \\ E_s \end{array} \begin{array}{|c|} \hline \\ \hline \end{array} \begin{array}{c} E_0 \dots E_s \end{array} \rightarrow \text{sums to 1 along each row!}$$

AXIOM 1: the Markov Property states that if at some point t the process is in state E_j , the probability that one timestep later we are at state E_i depends only on us currently being at E_j and not our previous trajectory to get to E_j .

Mathematically: $P(X_k = E_k | X_{k-1} = E_{k-1}, \dots, X_0 = E_0) = P(X_k = E_k | X_{k-1} = E_{k-1})$

AXIOM 2: The transition probability is independent of $t \Rightarrow$ temporally homogenous.

Definition: Where do we start our chain? $\lambda: \{ \lambda_1, \dots, \lambda_s \}$ where $\sum \lambda_i = 1$ is our initial state probability distribution.

Aperiodicity: There is no such state such that we revisit it at every t_0 (multiple) steps. If there are guaranteed cycles, we will not visit all the nodes or approximate the stationary distribution.

Define: the stationary distribution π is the probability we are in any state i .

Suppose we have a transition matrix $P \rightarrow$ if $\pi_j = \sum_k \pi_k P_{kj}$ at any timestep we observe that π does not change. It is stationary.

Mathematically: $\pi = \pi P$ must hold!!

irreducibility dictates that every state can be visited from every other state in some number of steps.

★ If we have disconnected components:



we will actually end up with two distinct stationary distributions! (think about why that might be :))

Theorem: Every finite and irreducible Markov chain has a unique stationary distribution. If π is the SD, of $M = (P, \lambda)$ then

$$\pi = \pi P$$

$$\sum_{k=1}^s \pi_k = 1$$

We can use linear algebra to solve $\pi = \pi P$ to solve for π .

3.2 MCMC ALGORITHMS

Theorem Let (X_0, \dots, X_n) be irreducible and aperiodic, MC with state space S and trans. mat. P . Our estimate $\hat{\pi}^{(n)}$:

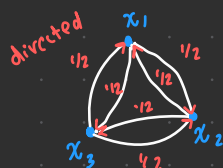
$$\hat{\pi}^{(n)} \rightarrow \pi \text{ as } n \rightarrow \infty.$$

We essentially will observe the true underlying distribution as we sample a large number of samples.

Definition: Let MC, S state set, P trans. mat., A distribution π on S is reversible if $\forall i, j: \pi_i P_{ij} = \pi_j P_{ji}$

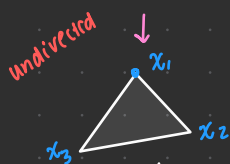
LECTURE 12 MCMC Continued

RANDOM WALKS ON GRAPHS



These are our transition probabilities \rightarrow so let's move it to a graph!

Graph $G = (V, E) \rightarrow$ Vertices and Edges



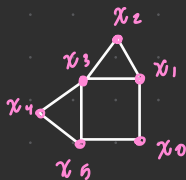
A random walk on a graph is a Markov chain with state set $V = \{v_1, \dots, v_k\}$ and the following transition mechanism:

- If at v_i , it moves at time $t+1$ to one of neighbors w. equal IP.
- This is then a function of $\text{degree}(v_i) =$ the number of neighbors.

Assume = IP to transition to any arbitrary neighbor

$$P_{ij} = \begin{cases} 1/d_i, & \text{if } i, j \text{ are neighbors} \\ 0 & \text{otherwise.} \end{cases}$$

Extended Example



Then the stationary distribution $\pi = \left[\frac{d_0}{d}, \frac{d_1}{d}, \dots, \frac{d_5}{d} \right]$ for $d = \sum d_i$

Intuitively, the nodes with the most neighbors will be visited the most number of times.

MARKOV CHAIN MONTE CARLO

Given a probability distribution $\pi = [\pi_1, \pi_2, \dots, \pi_K] \rightarrow \sum \pi_i = 1$ over a state space S .

How do we simulate a random object over the distribution π ?

the "simple solution":

Let $x = \Psi(u)$ where u is uniform $[0, 1]$ and the function Ψ :

$$\Psi(x) = \begin{cases} s_1 & x \in [0, \varphi(s_1)] \\ s_2 & x \in (\varphi(s_1), \varphi(s_1) + \varphi(s_2)] \\ \vdots & \\ s_K & x \in (\sum_{i=1}^{K-1} \varphi(s_i), \sum_{i=1}^K \varphi(s_i)] \end{cases} \rightarrow \text{a piecewise function that takes the "INVERSE CDF" method}$$

"INVERSE CDF METHOD": an algorithm

learn more about this in APMA 1690!

- Input: a given CDF $F(t)$ + a sample size n (the number of random variables)
- Output: A series of "iid" random variables.

1) Generate (pseudo) iid random variables u_1, \dots, u_n from $\text{Unif}(0, 1)$

2) for $i=1, 2, \dots$ $X_i = \inf \{t \in \mathbb{R} \mid F(t) \geq u_i\}$

Proof of Thm: (prev page)

Assume F^{-1} exists, then $G = F^{-1}$ (that's why we choose smallest t in Thm).

$$F(t) = \mathbb{P}\{w \in \Omega \mid X(w) \leq t\} = \mathbb{P}\{w \in \Omega \mid F^{-1}(u(w)) \leq t\} = \mathbb{P}\{w \in \Omega \mid u(w) \leq F(t)\}$$

$$= F(t) \checkmark \text{ this proof is ONLY possible with the assumption } F^{-1} = G.$$

this becomes impractical! Especially when K very large.

Metropolis's Algorithm

We want to simulate a given π over a set of states. The graph that represents the states must:

- the graph must be connected and irreducible
- each vertex should not be the endpoint of too many edges.

We then define

IP of choosing given = IP to transition to an arbitrary edge

$$P_{i,j} = \begin{cases} \frac{1}{d_i} \min \left\{ \frac{\pi_j d_i}{\pi_i d_j}, 1 \right\} & \text{if } i \neq j \\ 0 & \text{if not neighbors} \\ 1 - \sum_{e \sim i} \frac{1}{d_i} \min \left\{ \dots \right\} & \text{if } i=j \end{cases}$$

↓ debiasing term

Now we will prove that the

