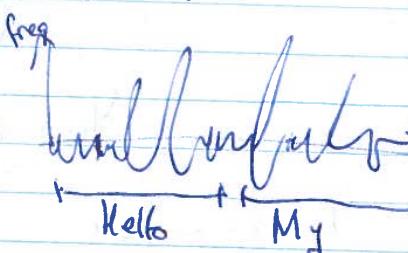
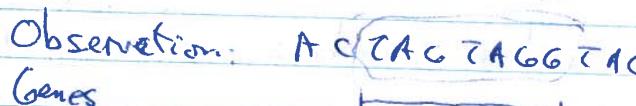
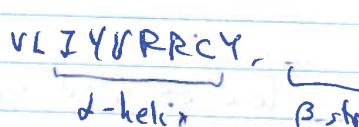


Hidden Markov Models (HMM)

Goal: "Decode" a sequence of observations

Speech recognition : Observation: 
freq

Biological sequence
Gene annotation : Observation: 
Genes

Protein: 
alpha-helix beta-str

Toy example: You visit a city with different neighborhoods.

Set of states = $S = \{F, E, C\} = S = \{S_1, S_2, \dots, S_n\}$

French English Chinese

You walk randomly in city.

Every minute, someone greets you.

Set of symbols = $\Sigma = \{b, h, n, a\} = \{\text{bonjour}, \text{hello}, \text{chinois}, \text{namaste}\}$

You record seq. of greetings: $X = x_1, x_2, \dots, x_L$

where $x_i \in \Sigma$

$X = b b q a h h h a b \dots$

Problem: Given: $X = x_1, \dots, x_n$

Find: Path $P = p_1, \dots, p_n$ that is most likely given X

We need to know:

① Emission Probabilities: $P_e[x_i = \alpha | p_i = \beta]$

States	Symbols			
	b	h	n	a
F	0.9	0.05	0.05	0
E	0.1	0.5	0.1	0.3
C	0.3	0.3	0.3	0.1

$\rightarrow \text{sum to 1}$

② Transition probabilities

$$\Pr_i [P_{i+1} = \gamma \mid P_i = \beta]$$

~~ε S~~

		Destination		
		E	C	F
Source	E	0.9	0	0.1
	C	0.1	0.4	0.5
F	0.3	0.1	0.6	

Path: E E E E E E F F F F E E F F C C C

③ Initial State probability

$$\Pr [P_1 = \gamma]$$

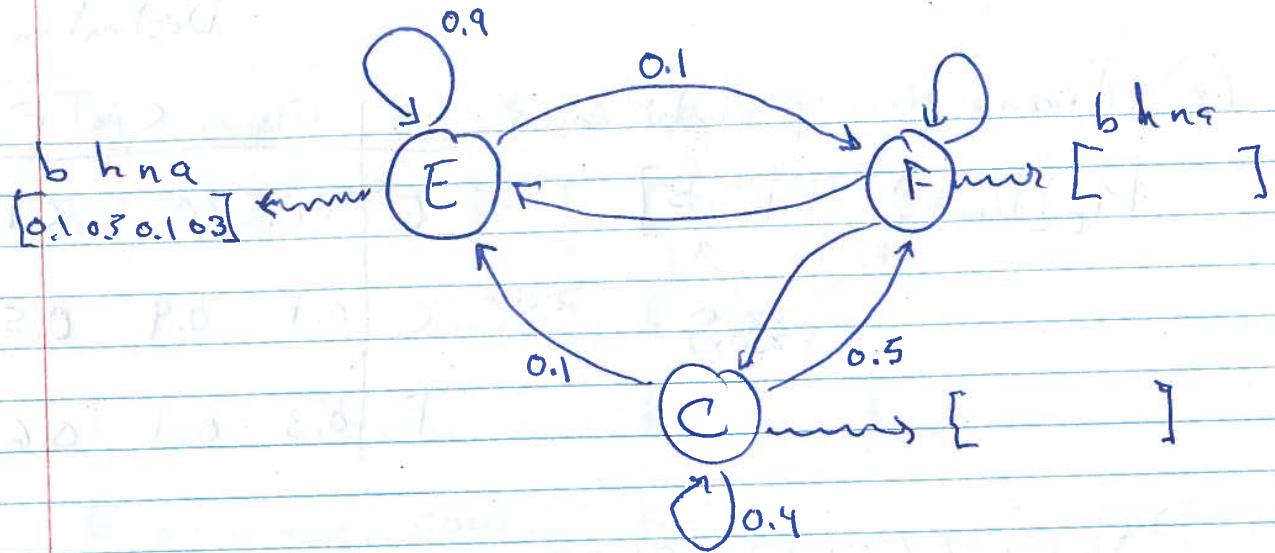
$\uparrow_{\epsilon S}$

		E	C	F
		0.5	0	0.5

Assumptions: - Markovian assumption (1st order)

Prob of going to state γ only depends on where you are now, not on how you got there

- Observations are independent from each other
(given the states)



HMM as generative model: Generate a random sequence

① Pick P_i randomly from the initial state prob. dist

Repeat:

- 2.1 Emit an observation from current state

- 2.2 Transition to next state acc. to transition prob. matrix

Questions (Assume we know S , Σ , emission prob, transition prob, initial prob)

① Maximum likelihood path:

Given: $X = x_1 \dots x_L$

Find: Path $P = P_1 \dots P_L$ that is most likely to have generated X , i.e. $\Pr[P_1 \dots P_L | x_1 \dots x_L]$ is maximized

Answer: Viterbi algo.

② Posterior decoding

Given: $X = x_1 \dots x_L$, time i , state β

Calculate: $\Pr[P_i = \beta | X]$

Answer: Forward-Backward algo

③ Estimation Problem: Assume we know S, Σ

but not emission prob

not transition prob

not initial



Given: $X = x_1 \dots x_L$

Find: Emission prob
Transit prob
Initial prob } such that $\Pr[X | E, T, I]$ is max

Answer: Baum-Welch algorithm

Viterbi Algorithm

Given: $X = x_1 \dots x_L$ (seq. of L observations)

$S, \Sigma, E, T, I = \text{HMM}$

states

alphabet

emission

transition

prob matrix prob matrix

initial state prob. vector

Find: $P = P_1 \dots P_L$, where $P_i \in S$

such that $\Pr[P_1 \dots P_L | X = x_1 \dots x_L]$ is maximized

Question 1: How to calculate $\Pr[P_1 \dots P_L | X=x_1 \dots x_L]$?
for a given path $P = P_1 \dots P_L$

$$\Pr[P_1 \dots P_L | X=x_1 \dots x_L] = \Pr[P_1 \dots P_L \wedge X=x_1 \dots x_L]$$

Independent of
Path P

$$\Pr[X=x_1 \dots x_L]$$

$$\begin{aligned} \Pr[P_1 \dots P_L \wedge X=x_1 \dots x_L] &= \Pr[X=x_1 \dots x_L | P_1 \dots P_L] \cdot \Pr[P_1 \dots P_L] \\ &= \left(\prod_{i=1}^L \Pr[x_i | P_i] \right) \cdot \Pr(P_1) \cdot \prod_{i=1}^{L-1} \Pr[P_{i+1} | P_i] \end{aligned}$$

Question 2: How to find P s.t. $\Pr[P_1 \dots P_L | X=x_1 \dots x_L]$
is max?

\Rightarrow Viterbi algorithm

Define $V(\beta, i) =$ Prob. of the most likely path of length i , given observation $x_1 \dots x_i$, assuming path ends in state β

$\in \Sigma^{\{1 \dots L\}}$

$$= \max_{\substack{P_1 \dots P_i \\ \text{where } P_i = \beta}} \{ \Pr[P_1 \dots P_i, x_1 \dots x_i] \}$$

$$V = \begin{array}{c|cccc|c|c} & x_1 & x_2 & \dots & x_i & & x_L \\ \hline S_1 & \textcircled{0} & \textcircled{0} & \textcircled{0} & \textcircled{Q} & \textcircled{0} & \\ S_2 & \textcircled{0} & \textcircled{0} & \textcircled{0} & \textcircled{Q} & \textcircled{0} & \\ \beta \rightarrow & \textcircled{0} & \textcircled{0} & \textcircled{0} & \textcircled{Q} & \textcircled{0} & \\ S_n & & & & & & \end{array}$$

Max

$$V(\beta, i) = \max_{\gamma \in S} \{ V(\gamma, i-1) \cdot P_r[\beta | \gamma] \} \cdot P_e[x_i | \beta]$$

Fill V table column by column, left-to-right

$\Rightarrow V(\beta, 1) = P_r(\beta) \cdot P_e[x_1 | \beta]$ (Initialization)

~~Pr~~ ~~P_e~~

$$\max_{P_1 \dots P_L} \{ \Pr[P_1 \dots P_L, x_1 \dots x_L] \} = \max_{\beta \in S} \{ V(\beta, L) \}$$

To recover path P that maximized $\Pr[P_1 \dots P_L, x_1 \dots x_L]$

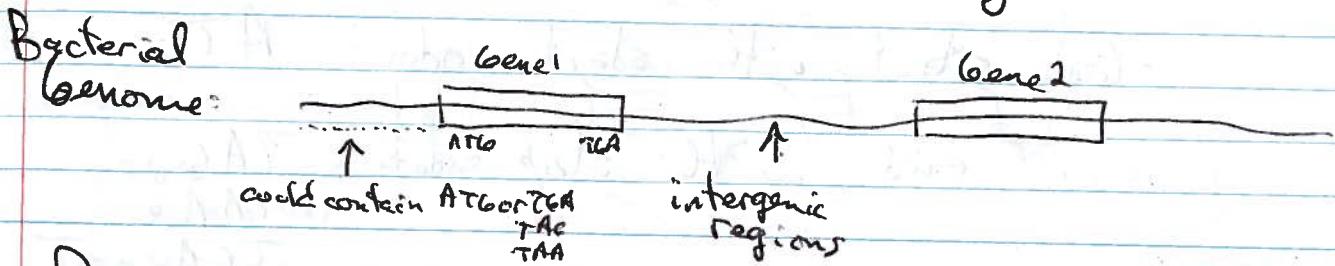
- Trace back dyn-prog algo from $\max \{ V(\beta, L) \}$

Complexity

Time = $O(n^2 \cdot L)$

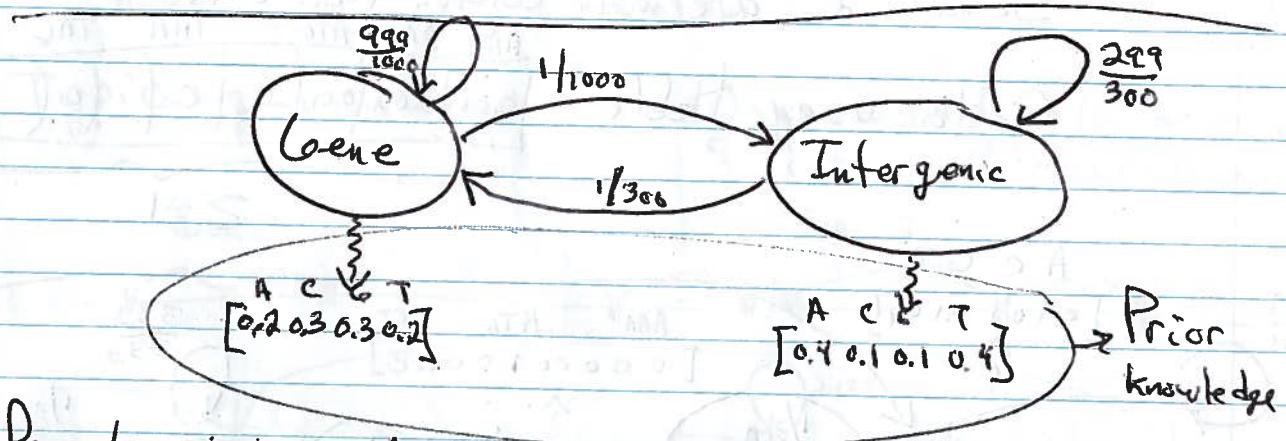
Space = $O(n \cdot L)$

HMMs for Gene Finding



Problem: Given: Genome sequence X

Find: Start/End position of each gene in X



Prior knowledge: Avg. gene length ~ 1000 bp
Avg. intergenic length ~ 300 bp

Genome $X = \text{ATAGATAAACAGGCCTCGTGGTC|ATAT}$

Viterbi Path = I I I I I I | 66 66 66 66 66 66 | I I I I

Integ Gene Interg

Gene Properties

- Gene start with start codon: ATG

- end with stop codons: TAG
TAA
TGA

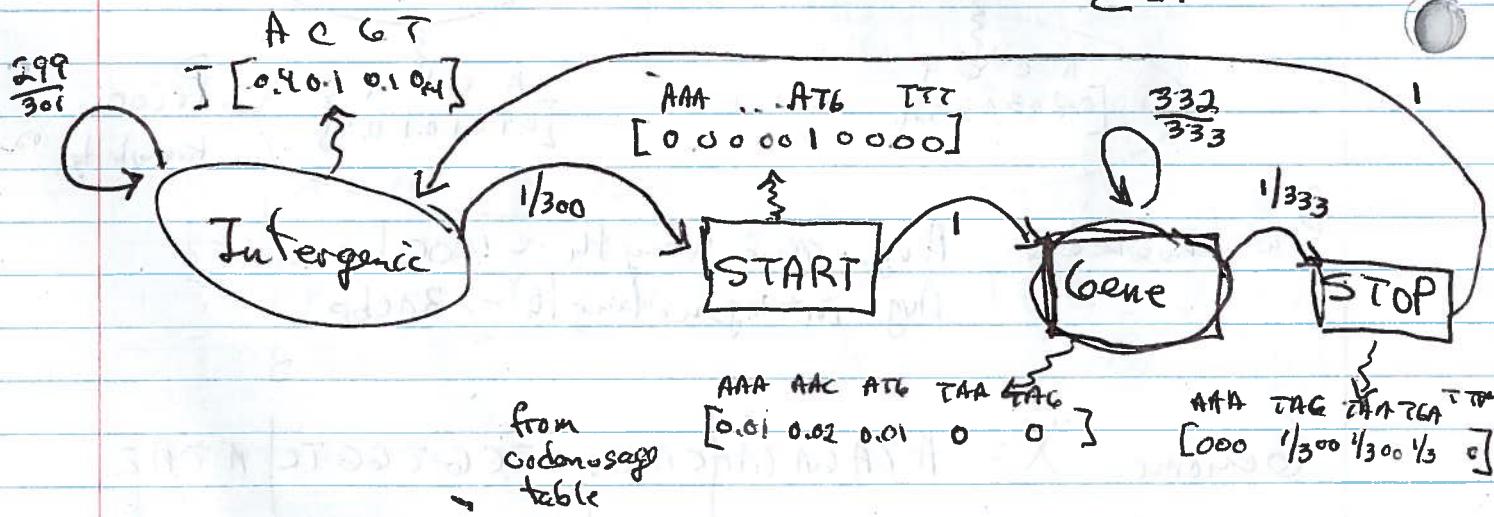
- Gene is made of codons (triplets of nuc.)

- Some codons are more common than others

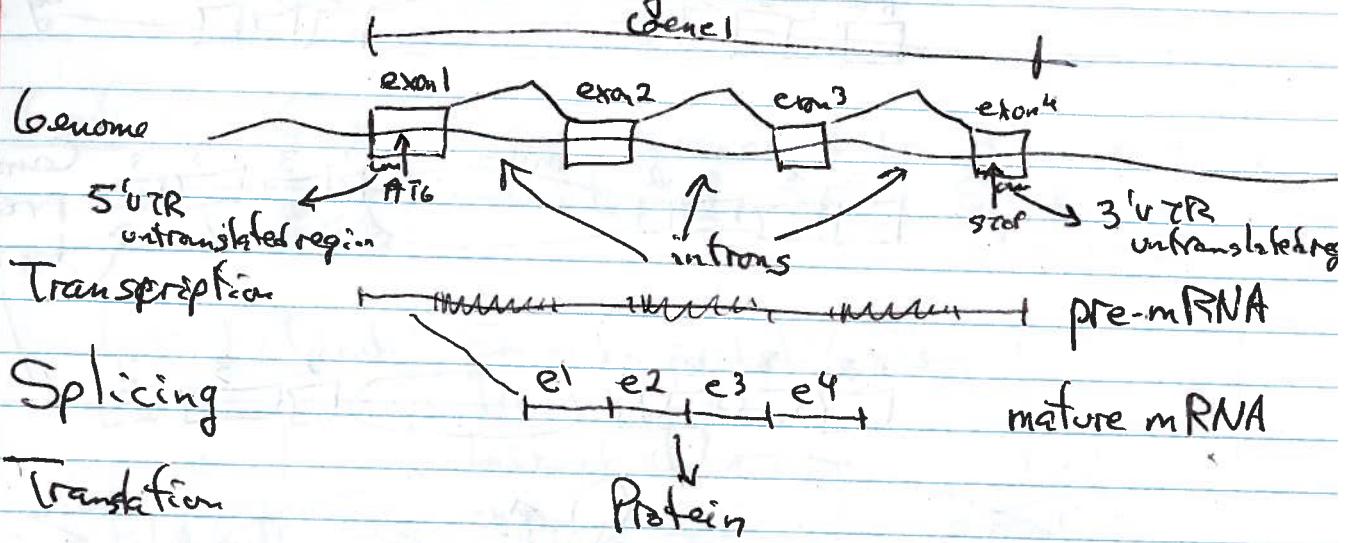
Codon usage table

AAA	AAC	ATG	TAA	TAG	TGA
0.01	0.02	0.01	0	0	0

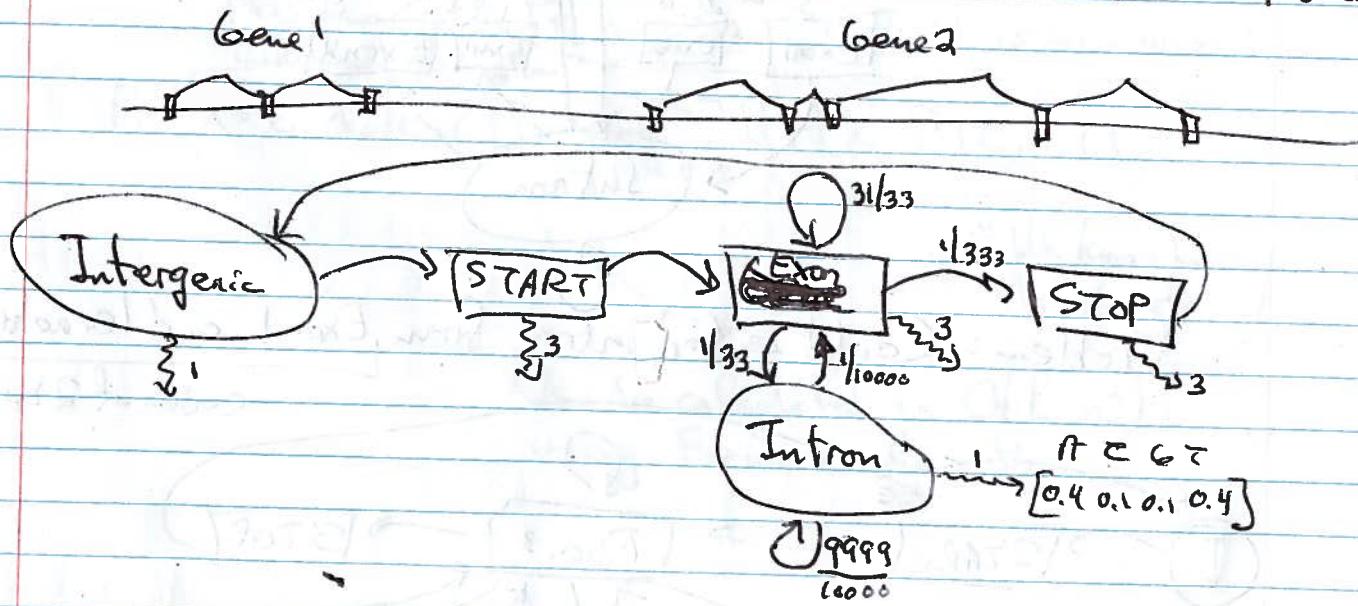
$\sum = 1$

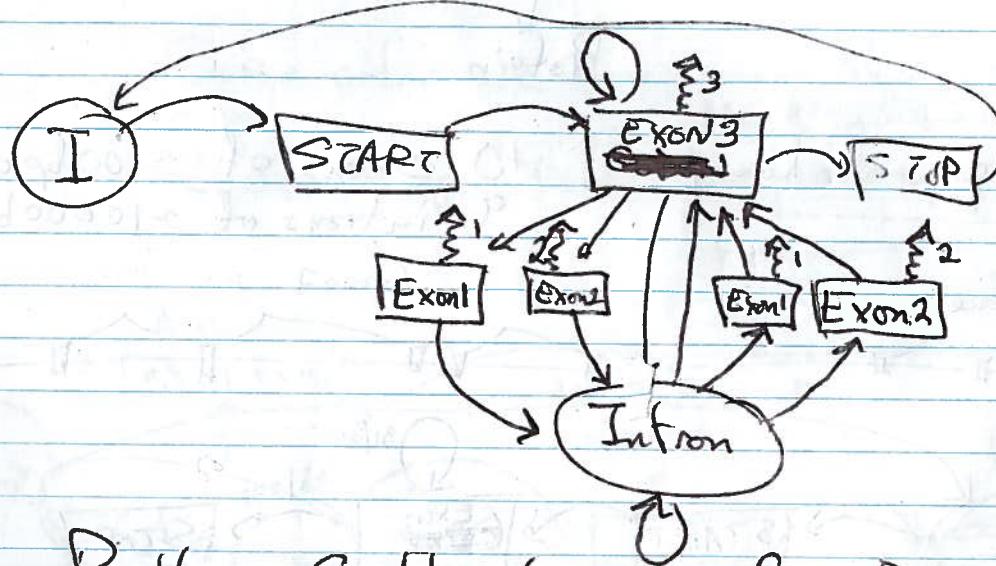
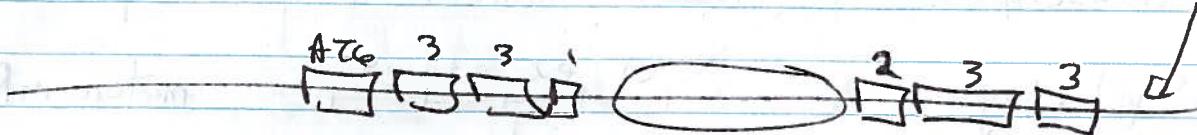
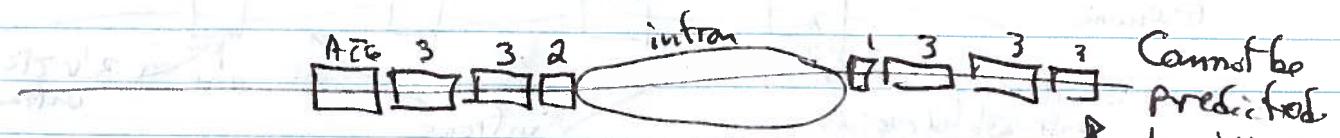
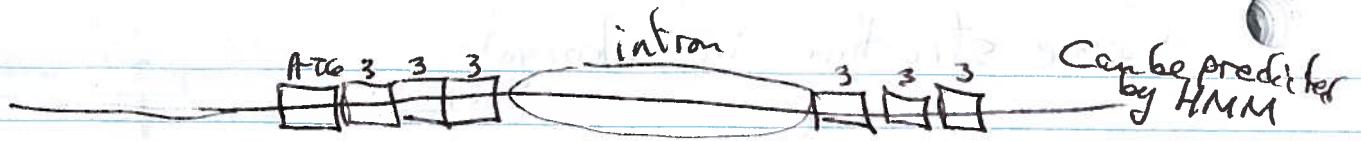


Gene structure in eukaryote

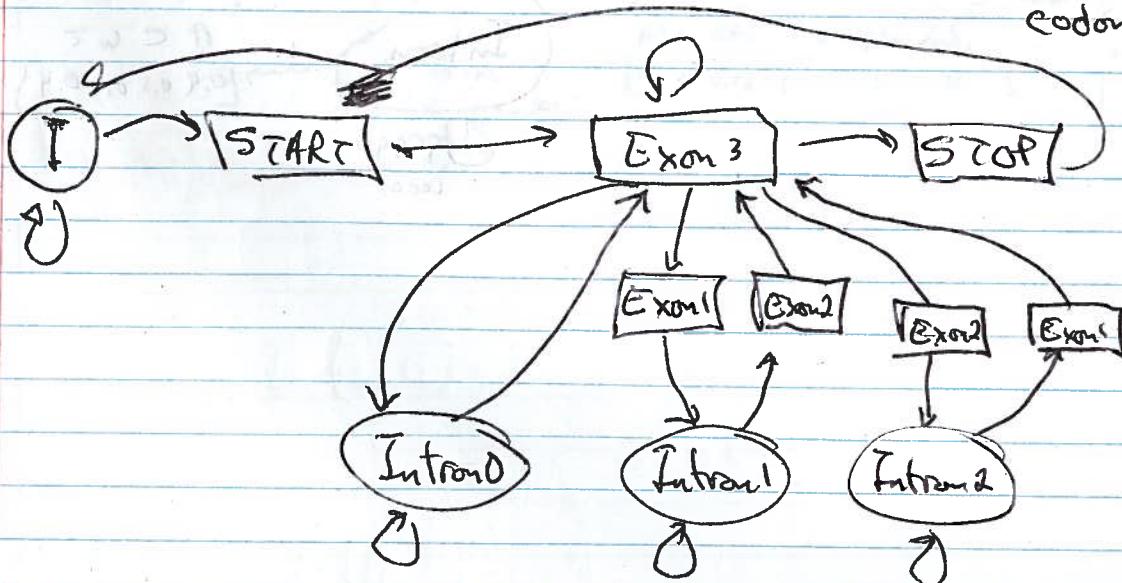


Typical gene in human : 10 exons of ~ 100 bp on average
9 introns of ~ 10000 bp on average





Problem: Could enter intron from Exon1 and leave from Exon1
codon of 2 bp



Training HMMs

Goal:

- Given:
 - Long sequence of observations $X = x_1 \dots x_L$
 - Structure of an HMM: Set of states S
 - Set of possible trans. T

~~fill in~~

Find:

- Emission probabilities : E
- Transition prob. : T
- Initial state prob : I

$$\Pr[A] = \sum_b \Pr[A, B=b]$$

such that E, T, I best represent the observed seq. X

Maximum Likelihood estimation $\leftarrow \Pr[X = x_1 \dots x_L | E, T, I]$ is maximized

$$\Pr[X = x_1 \dots x_L | E, T, I] = \sum_{\text{path } P = p_1 \dots p_L} \Pr[X, P | E, T, I]$$

We know how to calculate

Can be calculated in $O(L \cdot n^2)$ using Forward algorithm

$$(0, e)_{AB} = (0, e)_{AB} = 0.20, 0.10$$

$$(1, e)_{AB} = 0.10$$

$$P = (1, e) \rightarrow \text{what?}$$

~~class notes~~

Simplified version:

Given: $X = x_1 \dots x_L$

$P = p_1 \dots p_L \leftarrow \text{annotation of } X$

Find: E, T, I s.t. $\Pr[X, P | E, T, I]$ is max

Let $N_e(s, s') = \# \text{of times transition happened}$

b/w s and s'

= # of positions i s.t. $p_i = s \wedge p_{i+1} = s'$

Then, choose $T(s, s') = \frac{N_e(s, s')}{\sum_{x \in S} N_e(s, x)}$

$$\sum_{x \in S} N_e(s, x)$$

total # of
times in states

Example: $X = \overbrace{\text{ACA}}^1 \text{CAC} \overbrace{\text{ATGAC}}^2 \overbrace{\text{CTG}}^3$
 $P = \underbrace{\text{GGG}}_1 \overbrace{\text{TTT}}^2 \text{TGGG} \overbrace{\text{GG}}^3 \text{II}$

$$T(6, 1) = \frac{2}{7} \quad T(6, 6) = \frac{5}{7}$$

For emissions $\rightarrow N_e(s, a) = \# \text{times } a \text{ was emitted from } s$
 $= \# \text{positions } i \text{ s.t. } p_i = s \wedge x_i = a$

Choose $E(s, a) = \frac{N_e(s, a)}{\sum_b N_e(s, b)}$

Example: $E(6, A) = \frac{4}{7}$

Back to problem where P is not given

① Viterbi training

- Ⓐ Choose "reasonable" E, T, I (either from prior knowledge or randomly)

Repeat until convergence

- Ⓑ Use Viterbi algo to find best path for $X = x_1 \dots x_k$ assuming E, T, I \Rightarrow Produces path P

- Ⓒ Re-estimate E, T, I from P, X as done previously

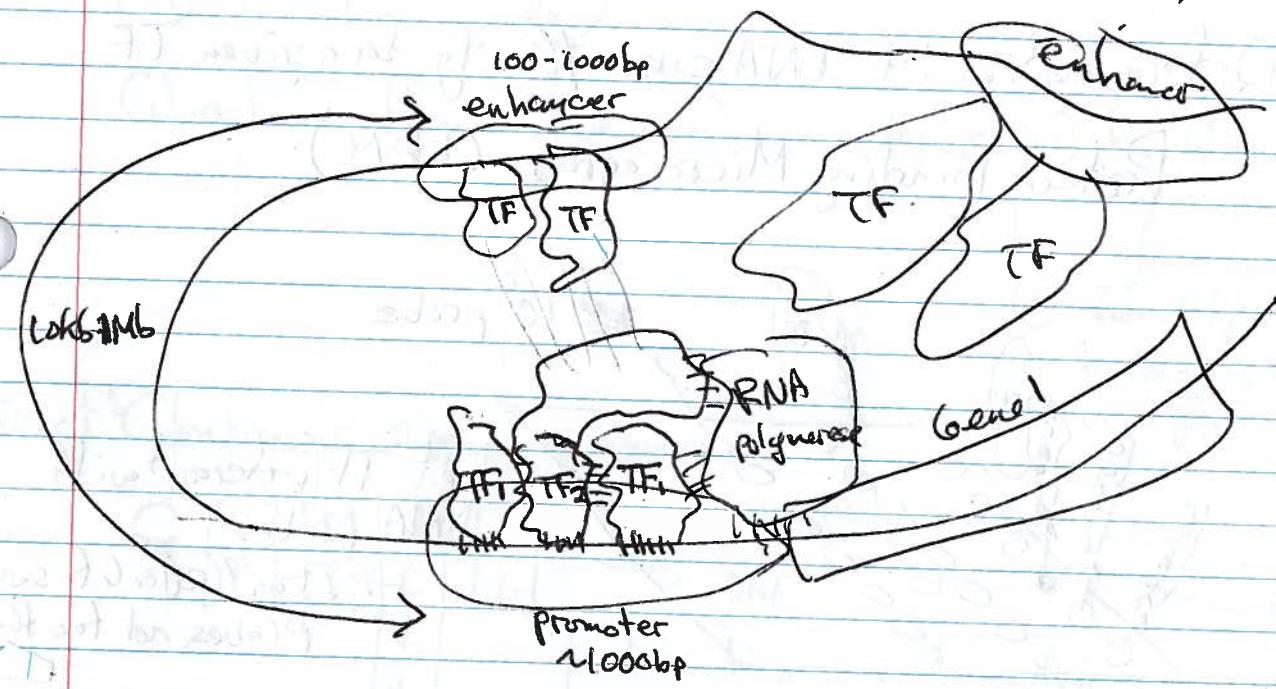
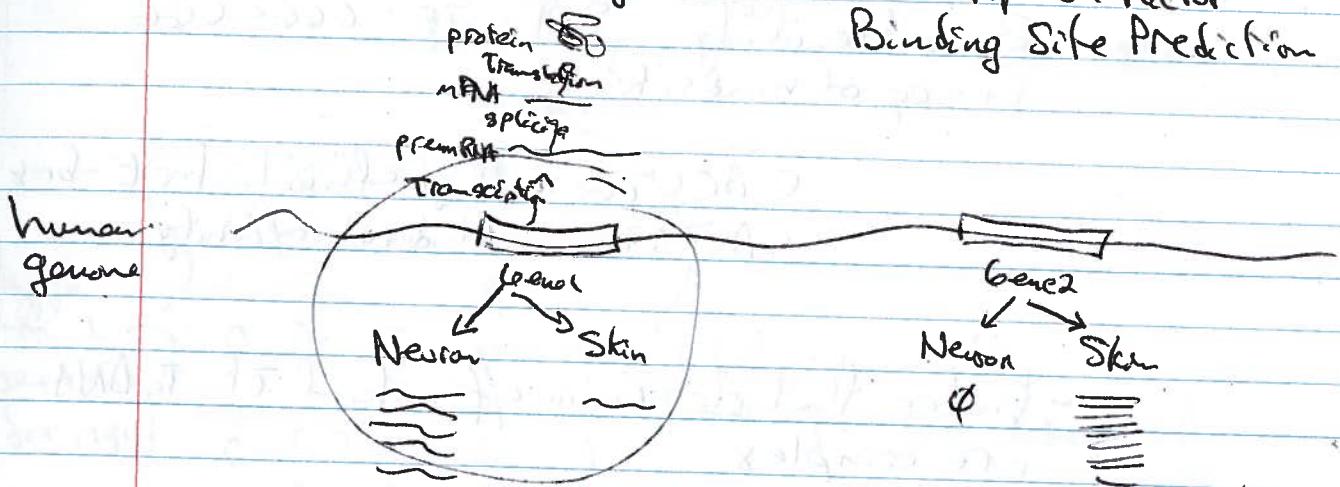
Problem: Algo. frequently gets stuck in local optima

② Baum-Welsch algorithm

Bases re-estimation of E, T, I on all paths (weighted by their probability)

\Rightarrow Reduces the prob of getting stuck in local optima

Gene Regulation + Transcription Factor Binding Site Prediction



Transcription factors:

- Proteins that
 - Bind specific DNA seq.
 - Alter expression of nearby gene(s)
 - activation
 - repression
- In humans: 2000 different TFs
 - ↳ each binds to different DNA sequences
- Transcription Factor Binding sites

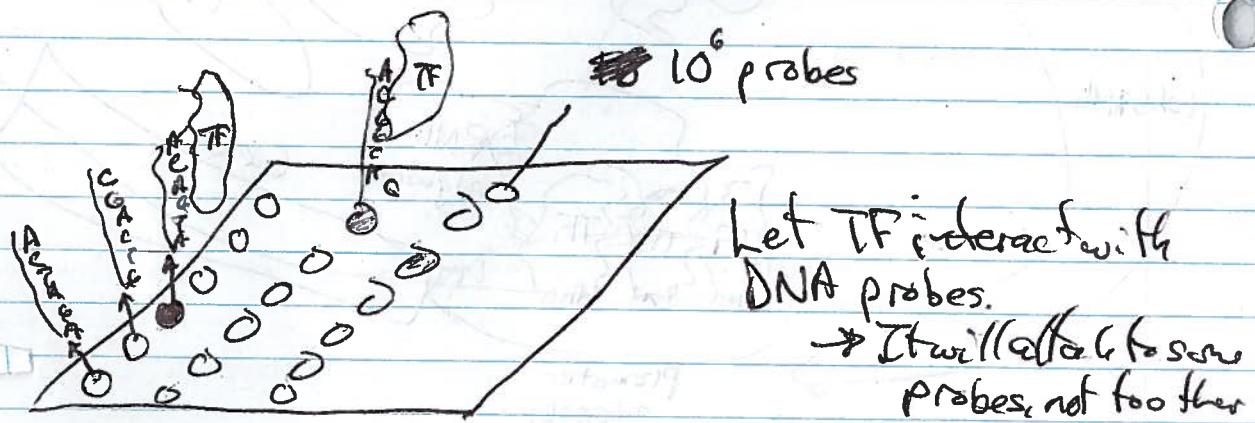
- TFBS:
- 6-15 bp long
 - Some flexibility in seq. of bind site
- E-box TF bind: CACGTC
- SPI TF: GGGC-GGG

CACGTC: High affinity for E-box
 CA~~T~~GTC: Moderate affinity

- Rules that determine affinity of TF to DNA are complex

Determination of DNA seq. affinity for a given TF

Protein Binding Micro-array (PBM)



Take picture of PBM \Rightarrow infer affinity of TF to each probe

\Rightarrow ACAGTA : intensity 1000
 AC~~G~~TA : 250

standard deviation:
 mean: AUG standard deviation:
 standard deviation:
 extra point!

Example: ~~Binding sites~~

DNA sequences bound by TF ~~myc~~ myc

not independent

high affinity sequences for myc

C	C	T	A	A
C	T	A	C	G
C	T	T	A	G
C	T	T	G	G
C	C	T	T	A
C	C	T	C	A
C	C	T	T	A

Representative sample of all possible sites TF

Question: How can we use

- in order to
- ① Build model for that TF's affinity
 - ② Identify new binding sites in a genome

Strict Consensus Sequence approach

C [C] [A] [A] [A]
[T] [T] [C] [G]
6
7

Match consensus

Position weight matrices

	1	2	3	4	5
A	0	0	1/7	2/7	4/7
C	1	4/7	0	2/7	0
G	0	0	0	4/7	3/7
T	0	3/7	6/7	2/7	0

or

C [C] [T] [A] [A]
[T] [C] [C] [G]
6
7

Candidate sequence
score

C T A G C

$$1 \times \frac{3}{7} \times \frac{4}{7} \times \frac{1}{7} \times \frac{3}{7} = \frac{9}{7^4} = 0.00...$$

How to identify matches for given PWM in given sequence

$S = \underline{\text{A C T G T C A C T T}}$

For each starting position i in S

Calculate score of $S[i, i+1, \dots, i+5-i]$ on PWM

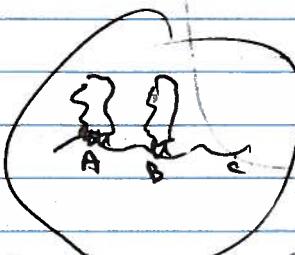
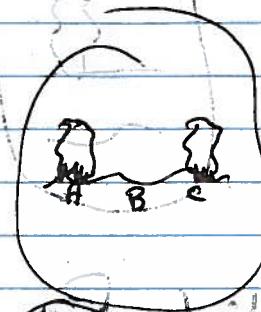
If score > Threshold ; then predict Binding
other no binding

Binded DNA

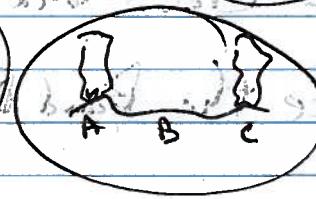
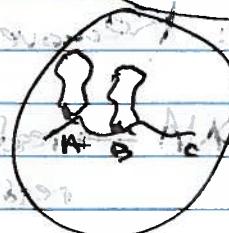
ChIP-Sq + Motif discovery

Goal of ChIP-Sq: Identify regions of genome bound by a given TF in a given type of cells.

ChIP-Sq: Chromatin Immunoprecipitation followed by Sq.



(10 Million cells)



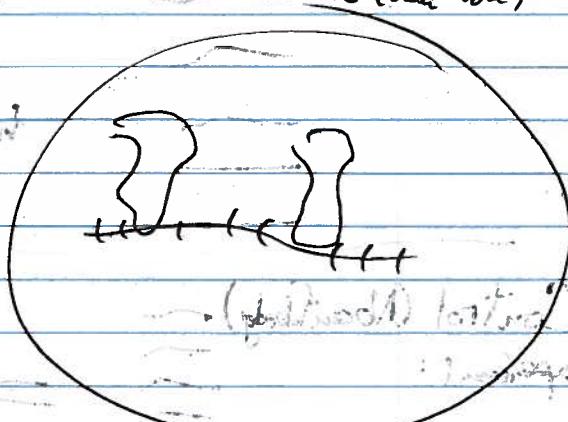
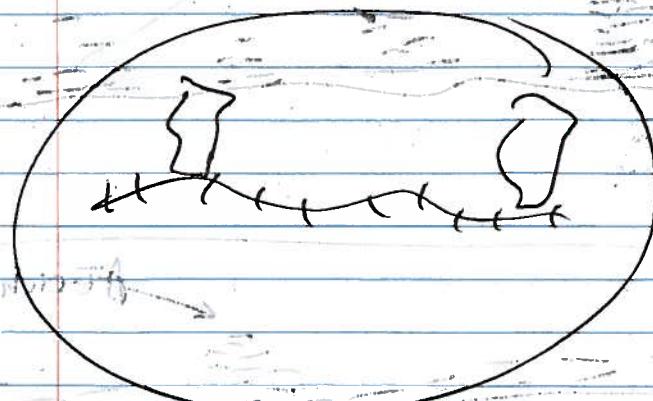
(Typically,

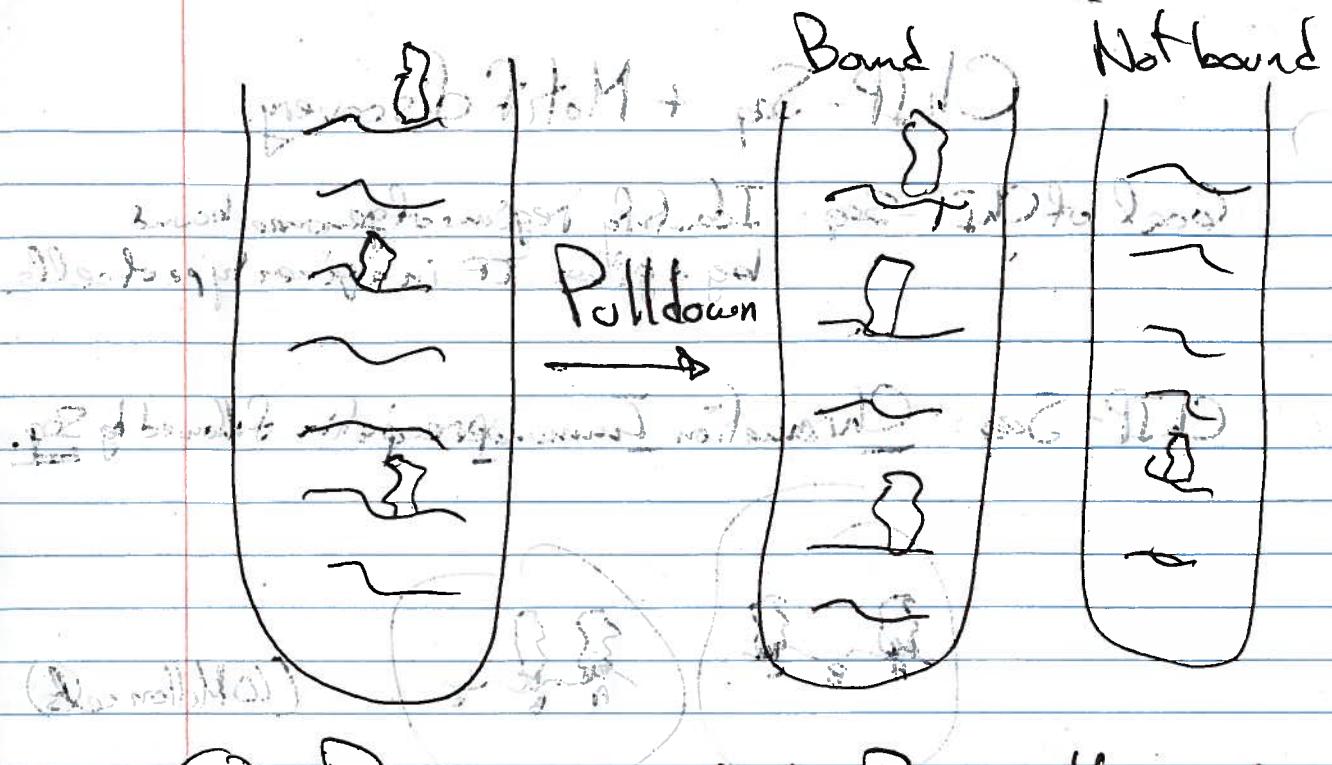
100 - 100,000 sites
occupied by TF)

① Cross-linking proteins to DNA: Strengthens bonds b/w protein and DNA

② Extract DNA with proteins attached to it

③ Fragment DNA in pieces of 120bp (~~sonication~~
sonication)





(5) Reverse cross-link: Remove the TF from DNA

(6) Sequence the Band DNA → read1: ACGTGCGTC
read2: TGCTGCAG..

read 10⁷: TAATGAT..

(7) Read Mapping

Local alignment b/w

Reads

human genome

Control (No antibody) ~
experiment:

peak: 100-500bp

Normalize: $\frac{\text{readcount}_{\text{exp}}}{\text{readcount}_{\text{ctrl}}}$

Stretch

100-100,000

ChIP-seq outcome: List of genomic regions (00-500 bp)
that are believed to be bound by TF

Problem: Doesn't tell us which 6-15 bp regions bind

⇒ Can't build consensus ~~or~~ or PWM for TF

Motif Discovery Problem

Given: Set of sequences S_1, S_2, \dots, S_n → lengths 60-500
believed to contain binding site for TF

Finds Consensus sequence for TF
OR
PWM for TF

Consensus sequence: Regular expression

$$w = \{A\} \{C\} \{A\} \{T\} \{C\} \{A\}$$

How should we score a candidate cons. seq $w = w_1 \dots w_k$
Criteria:

- w should occur in each of $S_1 \dots S_n$
 - too strict: Some S_i might be false positives
 - $\{\underline{A}\} \{\underline{C}\} \dots \{\underline{T}\}$ matches everywhere

Motif enrichment approach

Let $M_w = \# \text{of matches for } w \text{ in } S_1, S_2, \dots, S_n$

\rightarrow ~~subset of sequences from set S~~

$E_w = \text{Expected } \# \text{of matches of } w \text{ in}$
~~a set of random sequences } R_1, R_2, \dots, R_m~~
 where R_i has ~~same~~ length S_i

\hookrightarrow each nucleotide is chosen indep. with $P_A = 0.3$

$$w \in \{A, C, G, T\}^k$$

$$P_C = 0.1$$

$$P_G = 0.2$$

$$P_T = 0.3$$

How to compute E_w ? $w = w_1, w_2, \dots, w_k$
 where $w_i \in \{A, C, G, T\}$

$$w \text{ has } k \text{ matches} \rightarrow E_w = \sum_{i=1}^k P_{w_i}$$

$\Pr[w \text{ has a match starting at position } p \text{ in random seq}] = ?$
 \rightarrow ~~as if random N bp was broken into $k+1$ blocks~~

$$R : \dots \quad \boxed{\dots} \quad \dots$$

$$\hookrightarrow = \prod_{i=1}^k \Pr[w_i \text{ has match at position } p+i-1]$$

$$P_{\text{match}} = \prod_{i=1}^k \left(\sum_{a \in w_i} P_a \right)$$

$$E_w = (\#\text{positions eligible for match}) \cdot P_{\text{match}}(w)$$

$$= \sum_{i=1}^n (\text{length of } S_i - k + 1) \cdot P_{\text{match}}(w)$$

Finally, from N_w , E_w find Z_w

2. We want to find w where $|N_w - E_w|$ are the most different, i.e. $N_w \gg E_w$

Z-score approach

$$Z_w = \frac{N_w - E_w}{\sqrt{E_w}}$$

Complete algo. $\{w\}$ store st wth

For each possible consensus seq. w

Calculate N_w, E_w, Z_w

Report word w with highest Z_w

[Implementation to understand (w) \Rightarrow $\{w\}$ \Rightarrow Z_w]

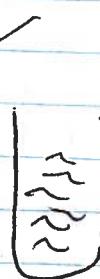
$$(w) \xrightarrow{\text{WSP}} \{w\} \xrightarrow{\text{WSP}} (w) \Rightarrow$$

$$(w) \xrightarrow{\text{WSP}} \{w\} \cdot (\text{down and right orientation}) \xrightarrow{\text{WSP}} (w) \Rightarrow$$

$$(w) \xrightarrow{\text{WSP}} \{w\} \cdot (1 + 1 - 12 \text{ rotation}) \xrightarrow{\text{WSP}} (w) \Rightarrow$$

DNA sequencing + Genome sequencing

Goal:



File: >seq1

ACGTGCTA

read1

>seq2

TGATCGATG...

read2

:

Illumina Sequencing: See video.

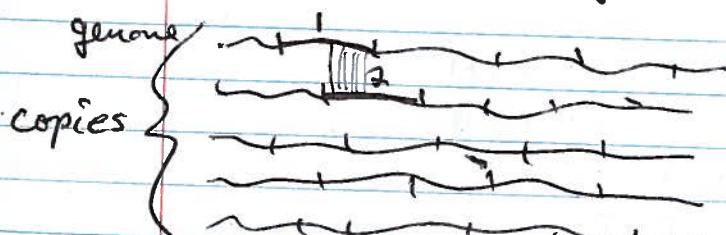
Limitation: - Length of a read is limited ($\leq 300\text{ bp}$) Illumina

Genome Sequencing + Assembly

Goal: Get entire DNA seq. of a genome ^{long}

Problem: Seq. machines produce reads that short

Shotgun sequencing

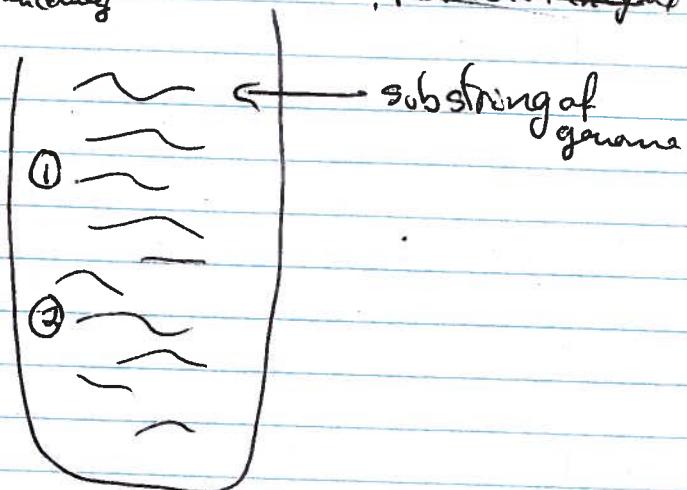


③ Sequence DNA fragments

ACTAGC
TAGCC...
:

④ Assemble reads into genome

- ① Generate many copies of genome
- ② Cut seq. into small pieces randomly \rightarrow sonication, restriction enzymes



True genome: ACTAGCTTTAGCCTT
(Unknown)

Read 1	CTTCTT
Read 2	ACTA <u>GC</u>
Read 3	<u>AGCC</u> TT
Read 4	CTTA <u>GC</u>
Read 5	T <u>AGC</u> TT
Read 6	TT <u>CTT</u> A

Problem: Shortest Superstring Problem

Given: Set of reads $R_1 \dots R_n$, of length $L = 6$

Find: Sequence σ such that

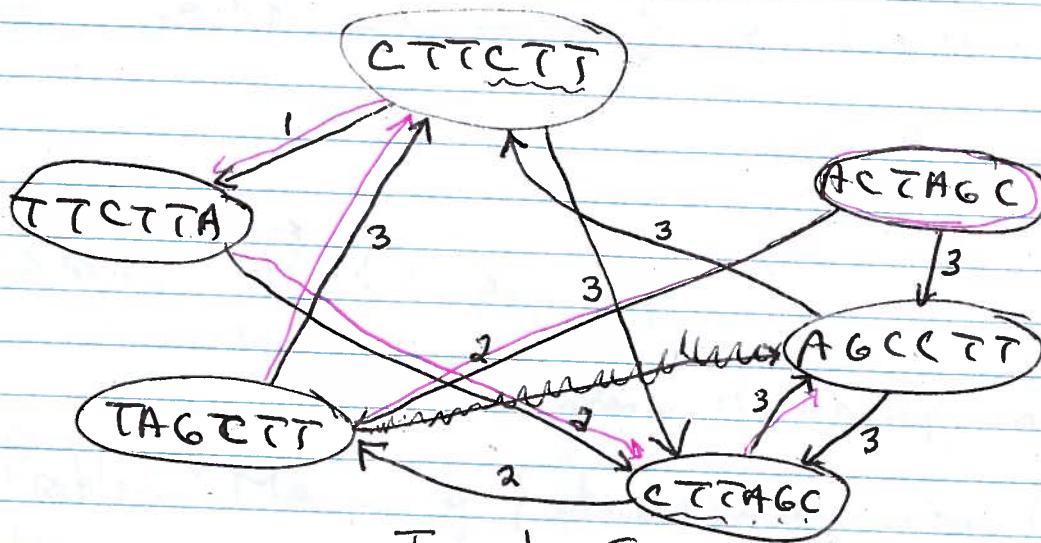
- 1) Each R_i is a substring of σ
- 2) σ is as short as possible

Assumption: No sequencing error.

Overlap-Layout-Consensus Approach

Build Graph: V : set of reads

E : overlaps between reads of at least $k = 3$ bases



Traveling Salesperson Problem

Goal: Find shortest Hamiltonian Path in G

Smallest total weight

Path that visits each vertex exactly once

NP Complete Problems

ACCTAGC
TA GCTT
CTTCTT

→ TTCTTA

CTTAGC

AGCCCTT

$$\text{weight: } 2 + 3 + 1 + 2 + 3 \\ = 11$$

Prefixed

ACCTAGC TTCTTA AGCCCTT

Closure

Gene Expression + Class comparison

Final Exam: Dec 14th ^{6pm} open book, covers all topics

Goal: Capture and compare "state" of different cells

cells with
positive action
upon treatment

negative
outcome

$$\vec{P} = (p_1, p_2, \dots, p_{20000})$$

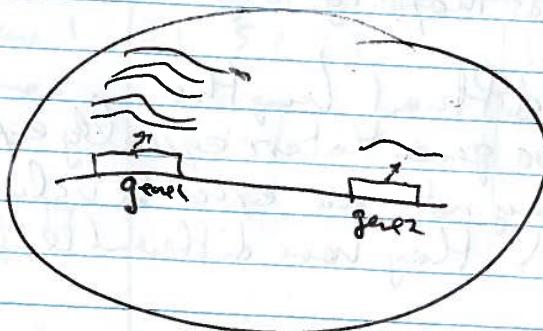
where p_i = abundance of protein p_i in cells.

Problem: Measuring protein abundance is hard (mass spectrometry)

Alternative: Measure mRNA abundance

$$\vec{G} = (g_1, g_2, \dots, g_{20000})$$

where g_i = abundance of mRNA from gene i in cells



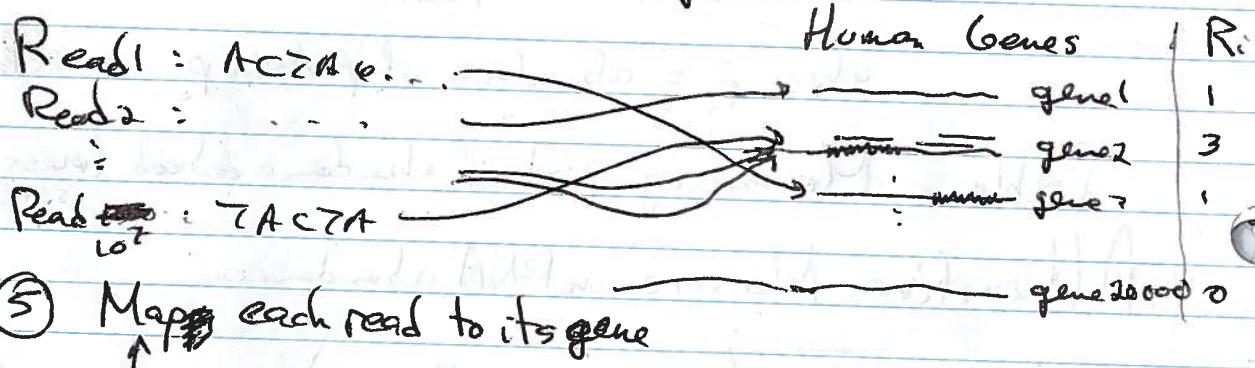
$$\begin{aligned} g_1 &= 5 \\ g_2 &= 1 \end{aligned}$$

Note: $g_i \neq p_i$ because mRNA degradation
translation is regulated

RNA-sequencing (RNA-seq)

Goal: Measure gene expression levels $\vec{G} = (g_1, \dots, g_{20000})$

- ① Extract RNA from cells
- ② Reverse transcribed RNA to cDNA
↑
complementary
- ③ Fragment DNA in ~200bp pieces
- ④ Sequence each cDNA fragment



- ⑤ Map each read to its gene

Find the gene to which the read aligns

- ⑥ Count $R_i = \#$ of reads mapping to gene i

Problem: Genes have different lengths

→ Two genes that are equally expressed may not have equal R values if they have different length

- ⑦ Normalize:

$$FPKM(g_i) = \frac{R_i}{\text{length}(g_i) \cdot (\text{Total RNA reads in Millions})}$$

Fraction per kilobase per million reads

Class Comparison Problem

Given: Normalized gene expression data from two sets of samples

A : (control) with $N_A = 20$ samples

$$\vec{A}_i = (A_1(i), A_2(i), \dots, A_{20000}(i))$$

$A_i(i)$ ↑ FPKM of gene i in sample 1

$$\vec{A}_2 = ()$$

$$\vec{A}_{N_A} = ()$$

B = (treatment) with $N_B = 25$ samples

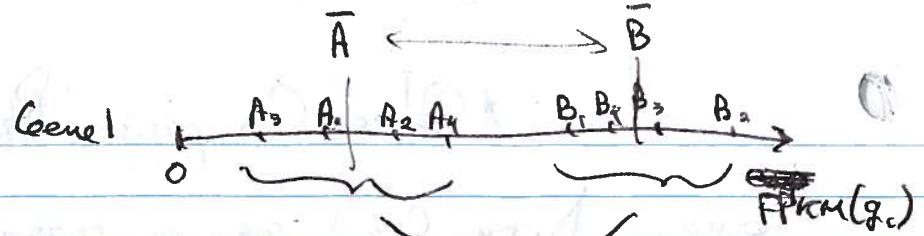
$$\vec{B}_1 = ()$$

$$\vec{B}_{N_B} = ()$$

	20			25			t-stat	p-value
	A_1	A_2	\dots	A_{N_A}	B_1	\dots	B_{N_B}	
Gene 1	5.1	5.7	3.7		1.7	1.3	1.9	2.1
Gene 2								1.3
:								
Gene 17								
Gene 2000								0.35
							best gene	$0.0001 = 4 \times 10^{-4}$

Goal: Find genes that are "differentially expressed" b/w A, B

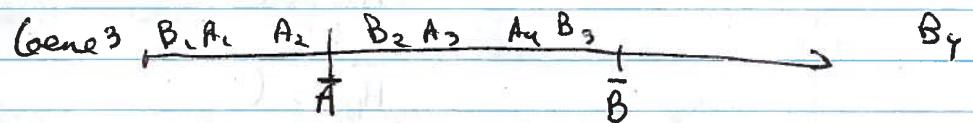
Example



Gene2 B₁ A₁ A₂ A₃ B₂ B₃ A₂ A₃

no significant diff,

$d(g_i)$ = diff. obs. $= \bar{A}(g_i) - \bar{B}(g_i)$



Student t-test (performed separately for each gene)

H_0 : Expression values from sample A and B come from the same normal distribution

$$\begin{aligned} \mu_A &= \mu_B \\ \sigma_A &= \sigma_B \end{aligned} \quad \downarrow \text{Assumption}$$

$$H_1: \mu_A \neq \mu_B$$

① Calculate $t(g_i) = \frac{\bar{A}(g_i) - \bar{B}(g_i)}{\sqrt{\frac{s_A^2}{N_A} + \frac{s_B^2}{N_B}}}$

s_A^2 = Variance in A

② Calculate p-value for $t(g_i)$ = Prob that two random samples drawn from same dist. would have t-statistic $\geq t(g_i)$

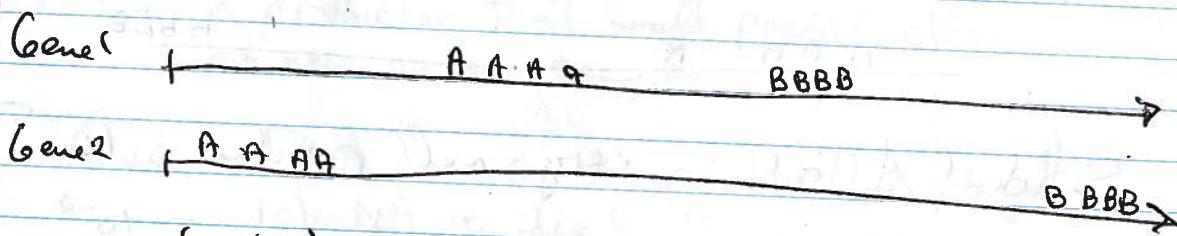
Under H_0 : t follow a Student (m)

degrees of freedom

$$m = \frac{\left(\frac{S_A^2}{N_A} + \frac{S_B^2}{N_B} \right)^2}{\frac{\left(\frac{S_A^2}{N_A} \right)^2}{N_A - 1} + \frac{\left(\frac{S_B^2}{N_B} \right)^2}{N_B - 1}}$$

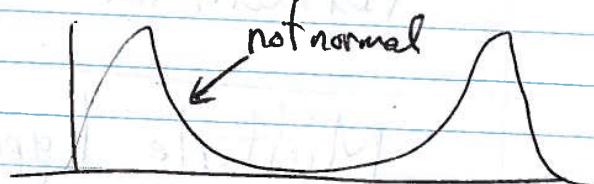
If p-value(g_i) ≤ 0.05 , then call g_i "diff. expressed"
 > 0.05 , then g_i is not diff. exp.

Issue: Violation of assumptions: - Normality of data



$$p\text{-value}(g_2) < p\text{-value}(g_1)$$

Suppose distribution under H_0



Permutation test: Estimate p-value without assumption about underlying distrib. of data

$$(1 - (-g)) \cdot \ln(-g) = (g) \ln(g) \quad \text{and this is monotonic}$$

For each gene i

Real	A ₁	A ₂	A ₃	A ₄	A ₅
Shuffled	1	1	1	1	1
	B ₁	A ₁	A ₂	B ₂	B ₃

① Calculate $t(g_i)$

② Repeat $K=1000$ times

2.1 Randomly reshuffle class to obtain \tilde{A}, \tilde{B}

2.2 Calculate $\tilde{t}(g_i)$ from

2.3 If $\tilde{t}(g_i) \geq t(g_i)$ then

③ Report $p\text{-value}(g_i) = \frac{\text{success}}{K}$

AAA A B BBB

Student t-test : very small p-value = 1e

Permutation : $p\text{-value} = \frac{1}{(8)} = 10^{-1}$

Multiple hypothesis testing

We've done 20,000 tests

If all genes come from H_0 , then best value we would expect to observe would be: $1/20,000$

Bonferroni correction : corrected $P\text{-value}(g_i) = p\text{-value}(g_i) \cdot \frac{N}{1}$

Class Prediction Problem

	A ₁	A ₂	...	A _{N_A}	B ₁	...	B _{N_B}	X
Gene 1								
2								
:								
Gene 20000								

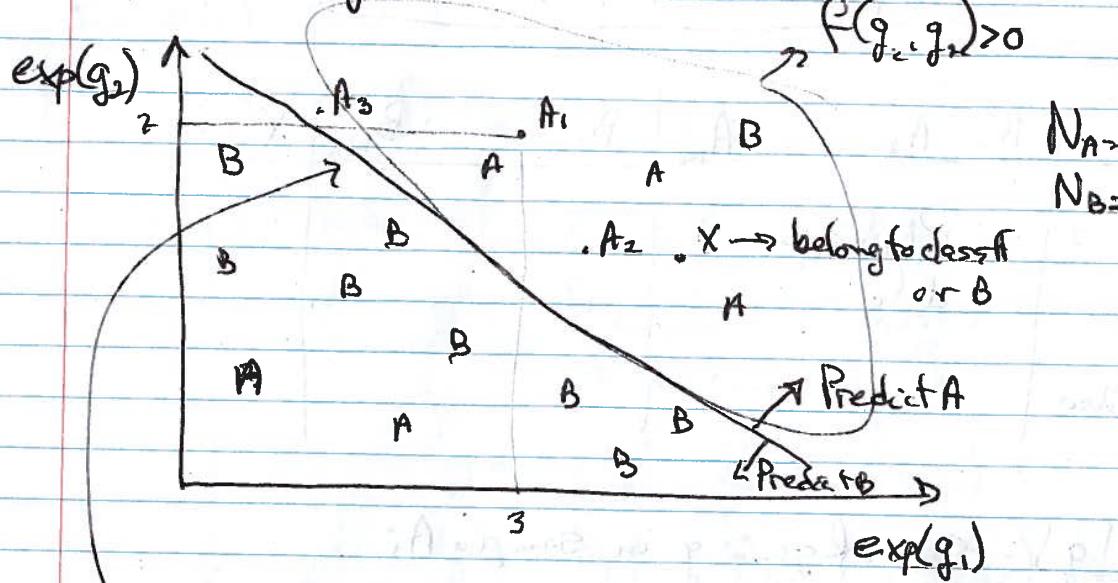
$A_i[g]$ = exp. of gene g in sample A_i

$B_j[g]$ = ...

Goal: From RNA-seq data from A, B, train a predictor that would predict class of new, unseen samples

Given: RNA-seq data X , predict if X belongs to class A or class B.

Assume #genes = 2



$$N_A = 8 \\ N_B = 9$$

Linear classifier: linear function of $\exp(g_1), \exp(g_2)$

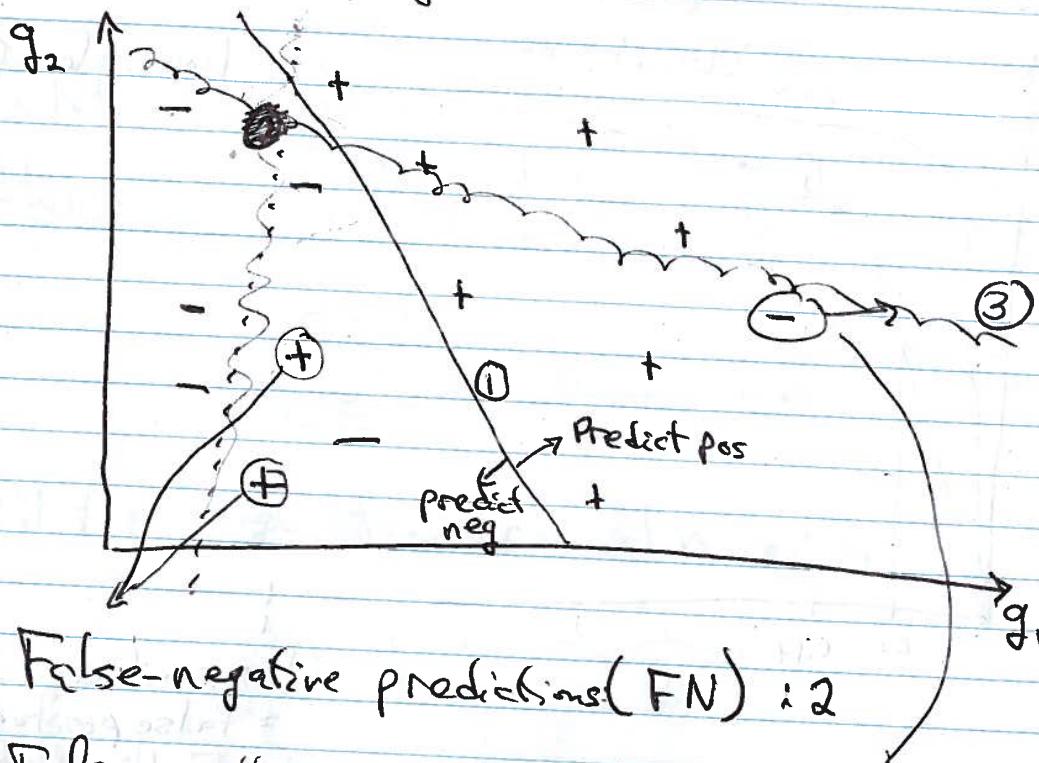
$$g_0 = 10 - 2g_1$$

$$f(g_1, g_2) > 0 \rightarrow \text{predict A} \\ < 0 \rightarrow \text{predict B}$$

$$f(g_1, g_2) = 4g_1 + 2g_2 - 10$$

$$2g_1 + g_2 - 10$$

Assessing a classifier's accuracy



False-negative predictions (FN) : 2

False-positive prediction (FP) : 1

True-positive (TP) : 7

True-negative (TN) : 5

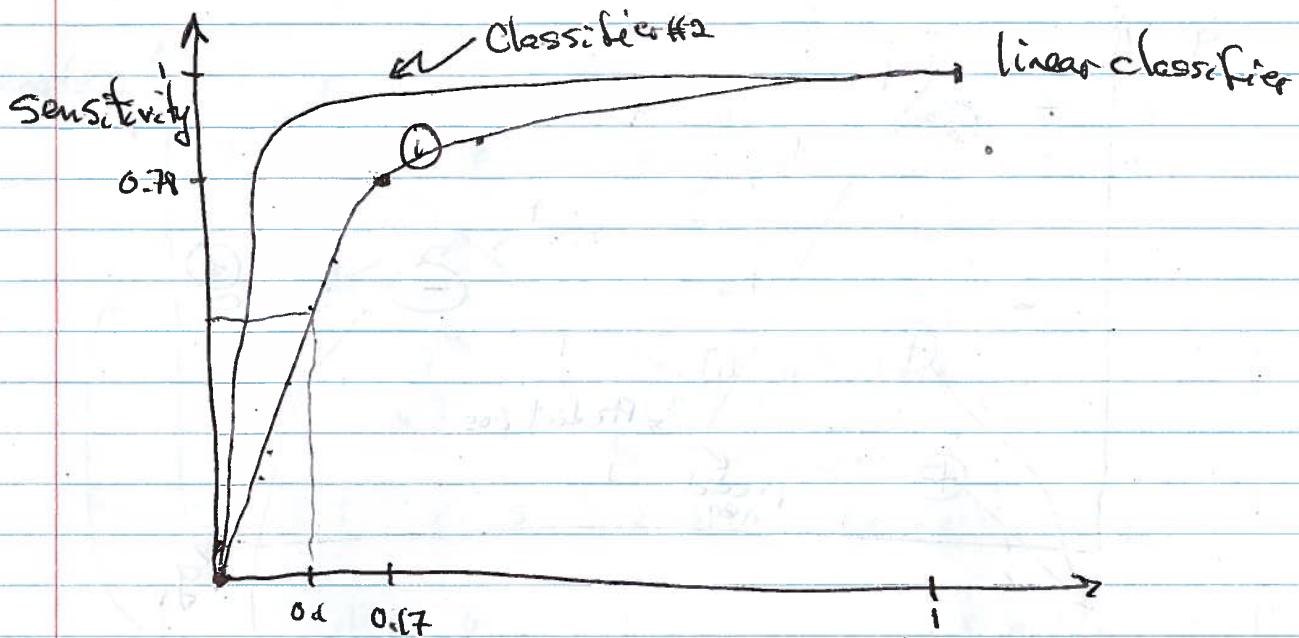
$$\text{① } \left\{ \text{Sensitivity: } \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{7}{7+2} = \frac{7}{9} \approx 78\% \right.$$

$$\left. \text{Specificity: } \frac{\text{TN}}{\text{TN} + \text{FP}} = \frac{5}{5+1} = \frac{5}{6} = 83\% \right.$$

② Sensitivity: 100%
Specificity: 50%

③ Sensitivity: $4/9 \approx 55\%$
Specificity: 100%

Receiving-Operating Curve



1-specificity

= false positive rate

= fraction of neg. examples
that are predicted pos.

$$AT + F = ST$$

$$\therefore S + F = 1 - ST$$

$$P_{EB} = \frac{S}{S+F} = \frac{S}{1-ST}$$

See next slide
for interface?

Look at slide 2
for interface?

RNA secondary structure prediction

DNA

Stable:

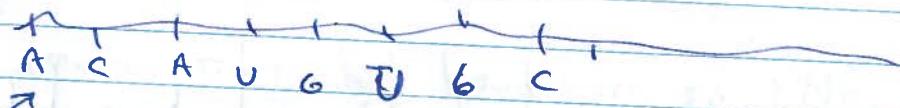


$$A = T$$

$\swarrow \searrow$ hydrogen bonds

$$C = G$$

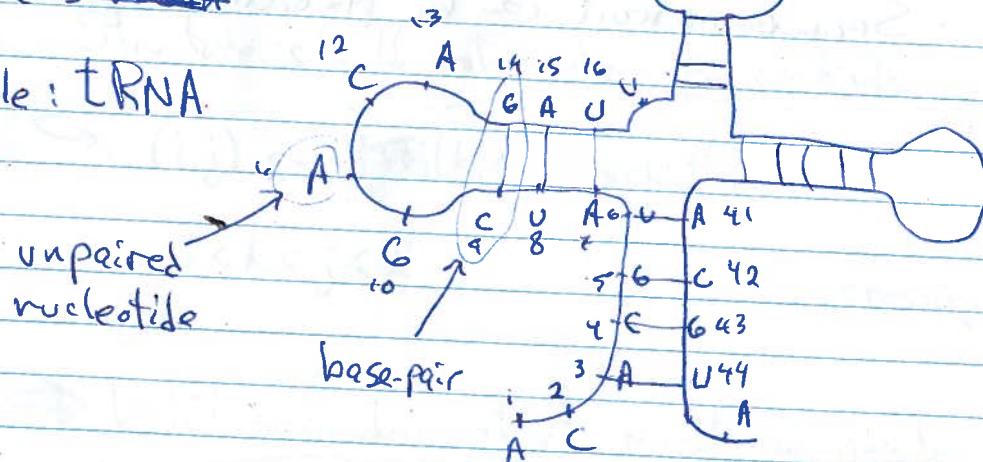
RNA : single-stranded chain of nucleotide



Unfolded structure is unstable \Rightarrow Sequence folds to maximize its stability

Folded sequence:

Example: tRNA



Function of RNA molecules depends on their structure

depends on sequence

Key idea: A sequence typically folds in its most stable structure

Secondary structure of sequence $S = S_1 S_2 \dots S_L$

is defined as list of pairs of positions that form base pair

Sec. struct. $\{ (3, 44), (4, 43), (5, 42), (6, 41), (7, 46), (8, 18), \dots \}$

Tertiary structure \equiv 3D structure: (x, y, z) coordinates of each atom in the sequence

Idea: • We want tertiary structure, but it is hard to predict

- Secondary struct. is sufficient to let us predict function, and also useful for predicting tertiary struct.
- Secondary struct. can be predicted computationally

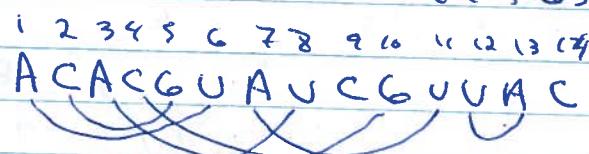
Secondary Structure prediction problem

Version 1: Given: Sequence $S_1 \dots S_n$

Find: Secondary structure for $S_1 \dots S_n$
that is the most stable

of base pairings present
in the sec. struct.

Example:



$$\{(1,6), (2,5), (3,8), \dots\}$$

Problem: Ignores rules about bendability of RNA

Rules

① if $(i,j) \in \text{Struct}$, then $|i-j| \geq 3$

② Structure should not contain pseudoknots

$\hookrightarrow (i,j)$ and ~~(k,l)~~ such that $i < k < j < l$

A diagram showing two crossing arcs representing base pairs. One arc starts at position i and ends at j, passing over another arc that starts at k and ends at l. This configuration is labeled as a "crossing base pairs".

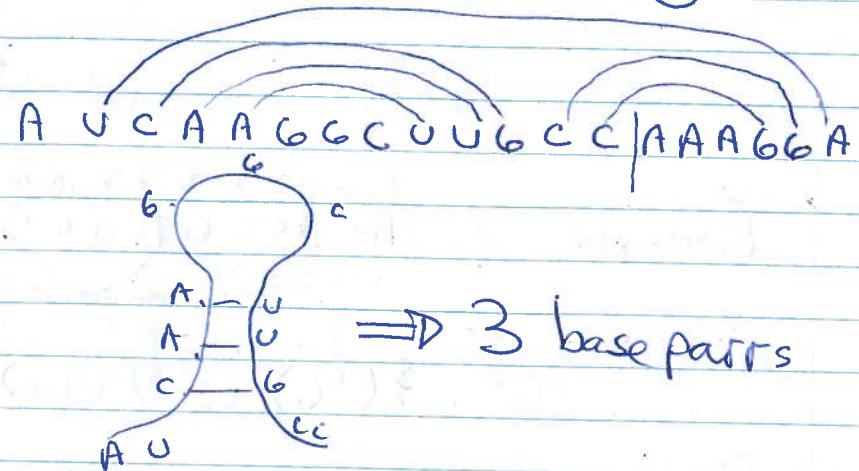
\Rightarrow Valid secondary struct must be nested



Version 2: Given: RNA seq S_{true} , S_L

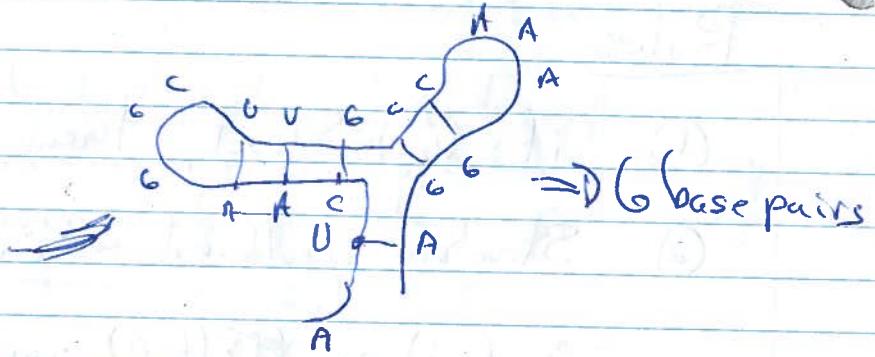
Find: Sec. struct. that maximizes total # of base pairs, subject to ~~the~~ rules ① and ②

Example:



Solution:
(short)

longer:

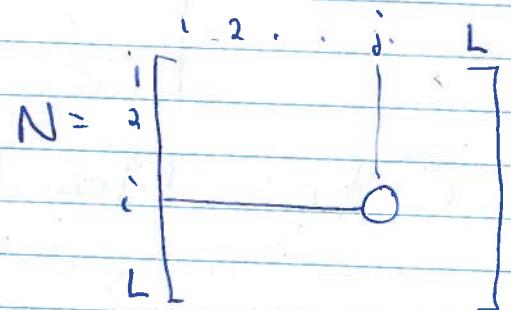


Nussinov Algorithm: Dynamic prog. algo.

Define $N(i, j)$ = Maximum # of base pairs that can be formed for $S_i \dots S_j$

We want $N(1, L)$

$S: 1 \ 2 \ 3 \ i \ \dots \ j \ L$



How to calculate $N(i, j)$?

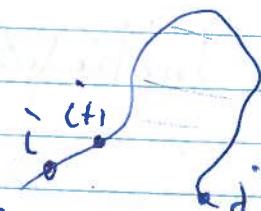
Case 1: S_i is paired with S_j

$$1 + N(i+1, j-1)$$



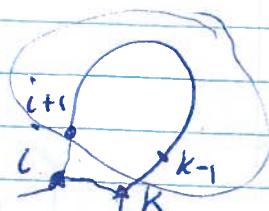
Case 2: S_i is not paired with anything

$$0 + N(i+1, j)$$

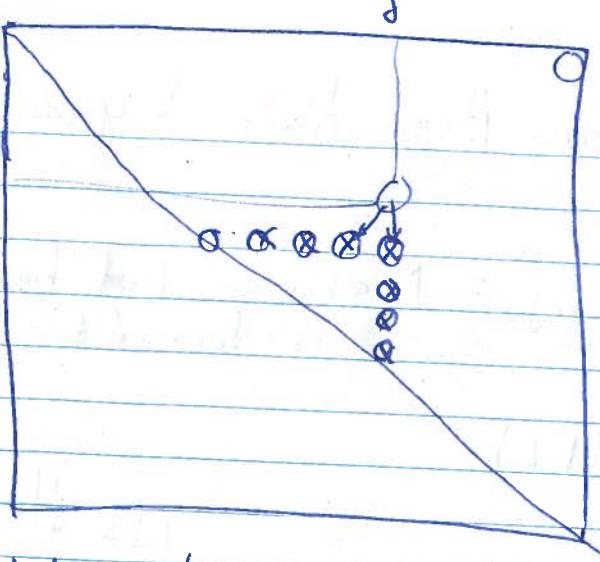


Case 3: S_i is paired with some nucleotide S_k , where $k < j$

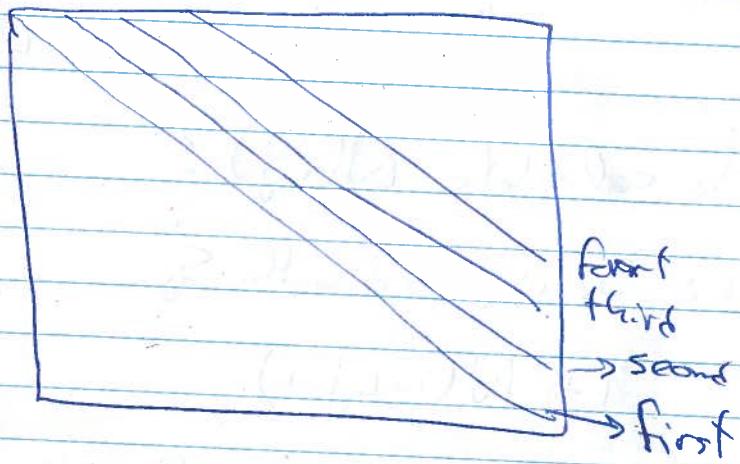
$$1 + N(i+1, k-1) + N(k+1, j)$$



$$N(i, j) = \max \begin{cases} 1 + N(i+1, j-1) & \leftarrow \text{if } j \geq i+3 \\ \quad \text{and } S_i \text{ and } S_j \text{ are complementary} \\ 0 + N(i+1, j) \\ 1 + \max \left\{ N(i+1, k-1) + N(k+1, j) \right\} \\ \quad i+3 \leq k < j, \text{ and } S_i \text{ is complementary to } S_k \end{cases}$$



Order: Main diagonal \rightarrow Corner

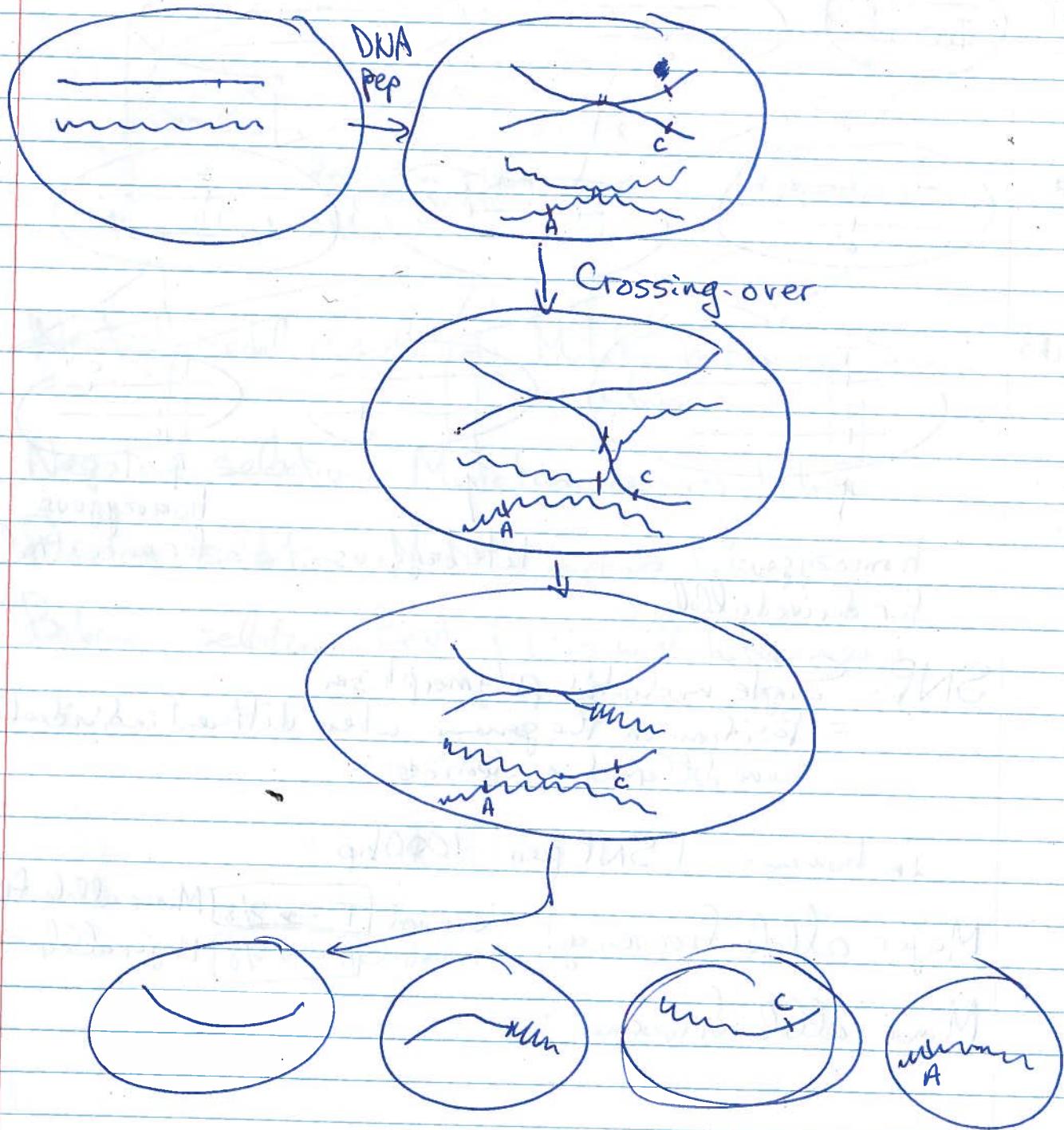


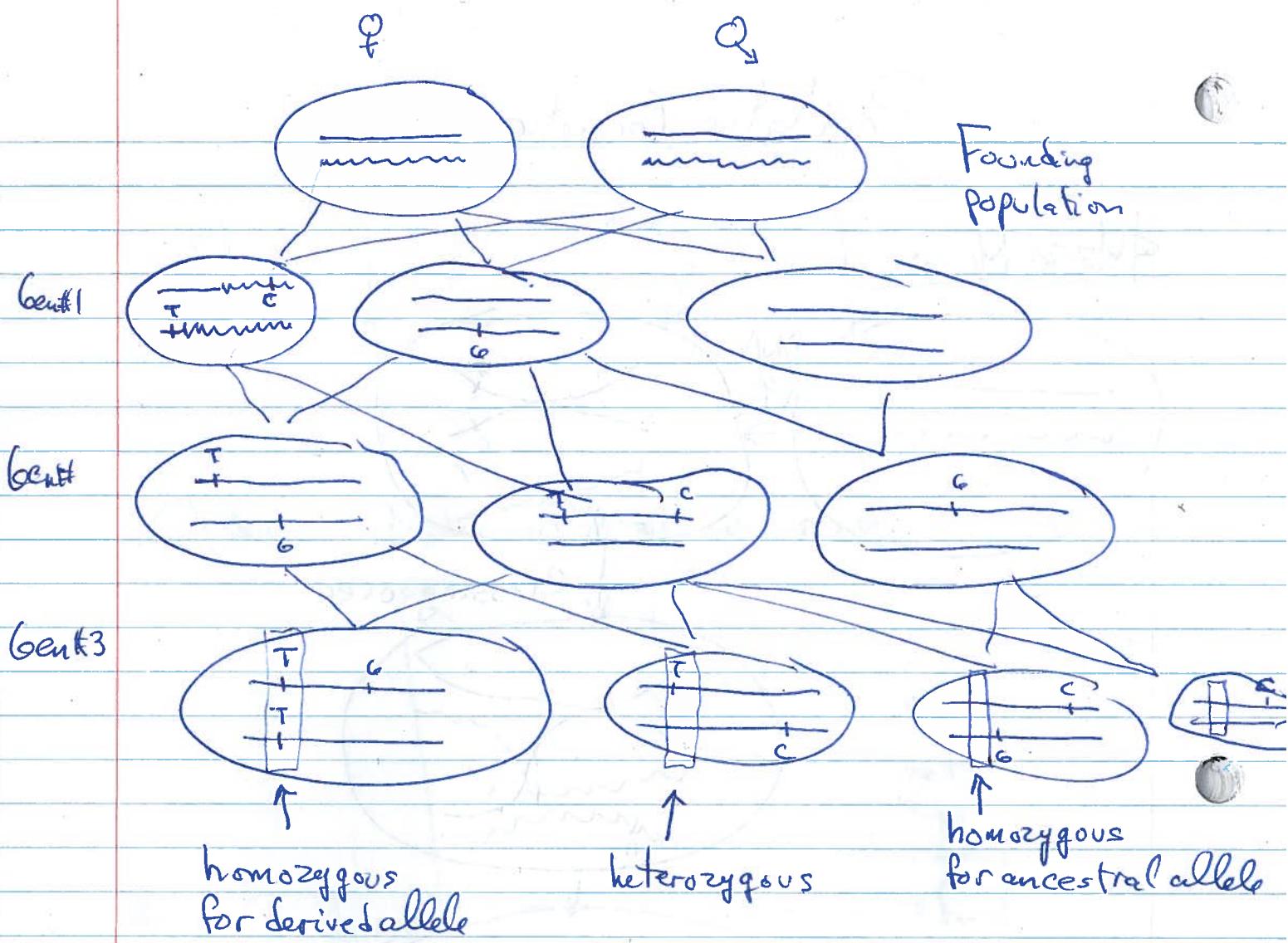
Initialization: $N(i,i) = 0 \quad \forall i \in L$
 $N(i,i+1) = 0$
 $N(i,i+2) = 0$

From there, use recurrence

Population Genetics

~~Meiosis~~





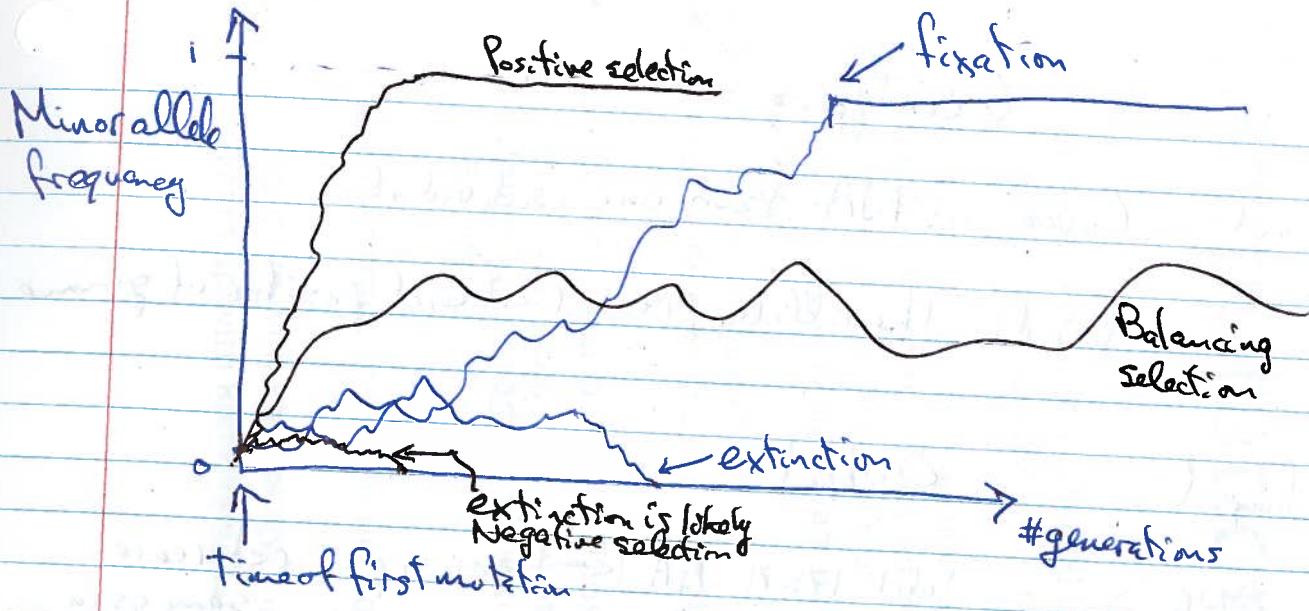
SNP: Single nucleotide polymorphism
 = Position in the genome where different individuals have different nucleotides

In human: ~ 1 SNP per 1000 bp

Major allele frequency:

derived: $T \rightarrow 3/8$ Minor allele freq: 0.375
 ancestral: $A \rightarrow 5/8$ Major allele: 0.625

Minor allele frequency:



— Neutral model of selection: Mutation has no consequences on fitness

— Negative selection: Mutation reduces fitness

Positive selection: Mutation improves fitness

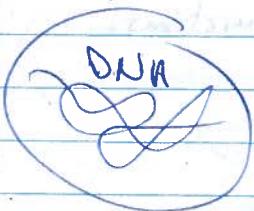
Balancing selection: Best fit is with heterozygous

Genotyping

Goal: Given: DNA from one individual

Find: The alleles present at each position of genome

Input



Output

chr1	173271	AA	← homozygous reference = same as in reference human genome
chr2	173471	AC	← heterozygous
chr7		TT	← homozygous non-reference

• ~~Genotyping arrays~~

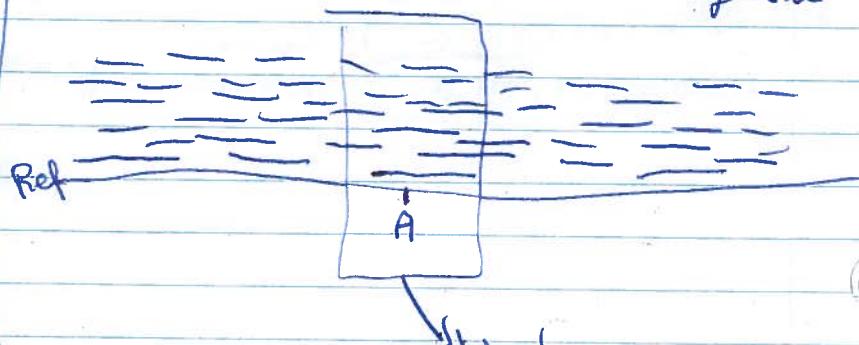
• Genotyping by sequencing

① Extract + fragment + sequence DNA

Fasta file

ACCTAAAC...
ATTACACCA

② Align (map) each read to reference human genome



Next page

C	A	A	T
A	A	A	T
A	A	A	T
C	A	A	T
C	A	A	T
C	A	A	T
A	A	A	T
I		T	T
A			T

↓ ↓ ↓

4C 1SA Homozygous
3A 1T non-referenc

Homozygous
references