

CS 176

---

---

---

---

---



# Exact String Matching

Def) String: ordered set of characters from alphabet  $\Sigma$ .

Def) Substring:  $P$  is a substring of  $T$  iff  $\exists 1 \leq k \leq |T|, P[i] = T[k+i]$ .

Def) Prefix: Substring that starts at the beginning of  $T$

Def) Suffix: // ends at the end of  $T$

Exact Matching: Given  $P$  and  $T$ , find all occurrences of  $P$  in  $T$ .

↪ Naively, compare every position until it fails  $\rightarrow O(MN)$  time  
 $(m := |P|, n := |T|)$

Idea: Take advantage of repeating patterns in  $P$  and  $T$

↪ Z-Algorithm: every character is compared a constant # of times

Def) Z function:  $\forall 2 \leq i \leq |S|, Z_i(S) :=$  longest substring from  $i$  matching prefix

ex)  $S = T T A T X T T T A T$

$Z(S) = / | \emptyset | \emptyset 2 4 | \emptyset |$

↪ How can this help compute the exact string matching?

$\Rightarrow$  Define a string  $S := P \$ T$  where  $\$$  is a character not in  $\Sigma$ . Compute  $Z(S)$ , and return any occurrences of  $P$  in  $T$  based on  $Z(S)$ .

Observation)  $\forall i > 1$ , we know that:

- $Z_i(S) \leq m$ .
- If  $i \leq m+1$ , then  $Z_i(S) < m$ .
- For  $i > m+1$ ,  $Z_i(S) = m$  iff  $P$  occurs in  $T$  at position  $i-(m+1)$

Def) Z-box:  $Z\text{-box } i := S[i \dots i + Z_i(S) - 1]$  (max prefix match)

↳ The rightmost Z-box's edge in any position reveals information!

Complete Z-Algorithm:  $i :=$  current index,  $(l, r) :=$  indices of rightmost Z-box

Initialization:

- $(l, r) \leftarrow (0, 0)$ ,  $i \leftarrow 2$
- Compute  $Z_2(S)$  directly
- If  $Z_2(S) > 0$ , update  $(l, r)$  accordingly

Induction:  $Z_2(S) \dots Z_{i-1}(S)$  is computed,  $(l, r)$  is based on  $(i-1)$ .

- Case 1) Uninformative:  $r < i$ , and nothing was revealed
- Case 2) Informative:  $i \leq r$ , then we know some characters!

↳ Be careful, however, for Z-values spanning outside the Z-box!

eg)  $S = \boxed{B \ B \ A \ B} \ B \ \boxed{B \ B \ A \ B} \ A \times A$

$$Z(S) = \boxed{\begin{matrix} / & 1 & 0 & 2 & 2 \\ \dots & & & & 4 \end{matrix}}$$

Concretely, define  $\beta := r - i + l$ , the length of the remaining Z-box.

Define  $j := i - l + 1$ , the matching position in the prefix Z-box.

↪ We will use  $Z_j(S)$  to compute  $Z_i(S)$ .

⇒ Case 1)  $Z_j(S) < \beta$ :  $Z_i(S) \leftarrow Z_j(S)$  ( $Z_i(j)$  is contained in Z-box)

Case 2)  $Z_j(S) > \beta$ :  $Z_i(S) \leftarrow \beta$  (outside the Z-box differs)

Case 3)  $Z_j(S) = \beta$ : direct compare from  $(r+1)$ , update  $(l, r)$

Time Complexity  $\approx$  # of character comparisons

$$= \underbrace{\# \text{ mismatches}}_{\hookrightarrow O(|S|) \text{ for every termination}} + \underbrace{\# \text{ matches}}_{\hookrightarrow \text{only happens after the rightmost Z-box} \rightarrow O(|S|)} \quad (\text{every match proceeds } r)$$

$$\Rightarrow \text{Time}(Z\text{-Algorithm}) = O(|S|) = O(m+n)$$

Space Complexity: we only need the immediate PI values of T!

Limitations: Rigid (no variations), Complexity blows up after 1-to-1

↪ Many-to-1 takes  $O(kn + km)$  time for  $k$  queries of len  $n$

Next time: How to get rid of the factor  $k$ ? (use preprocessing!)

## Exact String Matching II: Suffix Trees

Idea: use preprocessing to obtain a flexible data structure

↳ after preprocessing, queries will be constant time

Intuition: All substrings are a prefix of some suffix.

Naïve Algo: enumerate all suffixes of T, test if P is their prefix

↳ Inefficient, quadratic time

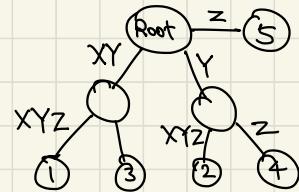
Ideally, if any two suffixes share a prefix, store only once!

Def) Suffix Tree: Rooted directed tree with n leaves, each corresponding to a suffix.

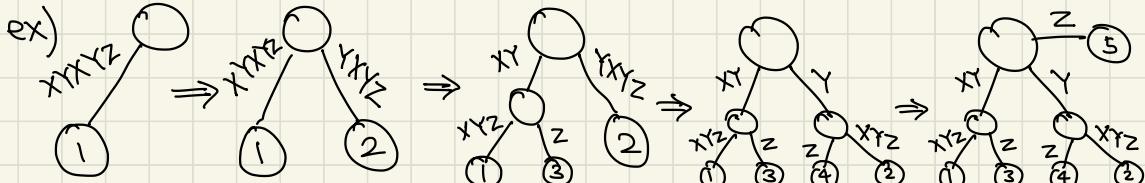
- Every edge is labeled by a nonempty substring of T
- A path to a leaf spells out a unique suffix
- Outgoing edges from the same node cannot start with the same character

ex)  $T = XYZXYZ$

- Each non-root interior vertex has at least 2 child nodes



Construction: Add suffixes from start to end, one by one.



- Observation) 1) Tree  $\rightarrow$  suffix is trivially easy.
- 2) Paths to leaves is the complete list of suffixes.
- 3) Paths are "compact"; if suffixes share prefixes, they share a parent.
- \* Not all strings have a valid suffix, specifically if ∃ indices  $i, j$  s.t.  
suffix  $i$  is a prefix of suffix  $j$

Space Complexity:  $O(|T|)$  ( $O(\# \text{ of leaves in tree})$ )

Time Complexity:  $O(|T|^2)$  (repeated character except last)

- \* There is a  $O(|T|)$  construction; refer Gusfield.
- $\Rightarrow$  The suffix tree helps solve exact matching and LCA & LCS.

Exact Matching: Traverse the tree until desired prefix is found.

All leaves of its subtree are suffixes that have a matching prefix.

Time Complexity:  $\underbrace{\text{Time traversing } P}_{\simeq \# \text{ char comparisons} = O(|P|)} + \underbrace{\text{Time enumerating occurrences}}_{\simeq \# \text{ of occurrences in } T}$

$\Rightarrow$  Overall  $O(|P| + k)$  where  $k := \# \text{ of occurrences of } P \text{ in } T$

$\hookrightarrow$  For many-to-1 searches, suffix trees runs in  $O(|T| + \sum_{i=1}^k |P_i| + k)$

Longest Common Ancestor: Find the lowest common ancestor!

Runtime: Naively, traverse the tree for  $O(|T|)$  time.

↳ There is a constant time retrieval structure for LCA in trees.

⇒ Overall  $O(1T + q)$  where  $q := \# \text{ of queries}$ .

Longest Common Substring: Define new string  $Z := S_1 \# S_2 \$$  to prevent the LCS spanning both  $S_1$  and  $S_2$ .

↳ we must "trim down" labels that spans  $S_1$  and  $S_2$  and denote from which string it originated from!

⇒ find the "deepest" vertex that is a joint ancestor of both  $S_1$  &  $S_2$

↳ recursively check if childrens have different origins (bottom-up)

⇒ Time  $O(|S_1| + |S_2|)$ , can be generalized to  $n$  strings

\* Complexity for suffix trees can be deceiving due to large constants

### Exact Matching III: Suffix Arrays

"Do we really need to store the entire suffix tree?"

Def) Lexicographical Ordering: ordering of alphabet in  $\Sigma$

↳ Assume that delimiters (\$) precede all alphabet.

Intuition: Let  $\Omega = \{i_1, \dots, i_k\}$ , the set of matching suffixes with  $P$

Claim: any suffix  $j$  not in  $\Omega$  must be lexicographically larger than or smaller than all suffixes in  $\Omega$ .

Proof of Claim: If  $j$  is between some  $i_x$  and  $i_p$ , it must be in  $\Omega_{j,x}$ .  
→ All possible sets of prefix matches will be clustered!

Def) Suffix Array:  $A_T[i] :=$  start position of lexicographically smallest suffix  
How to build the suffix array?

- 1) Merge sort:  $n \log(n)$  comparisons,  $O(n)$  per comp  $\rightarrow O(n^2 \log n)$  time
- 2) Use Suffix Tree: DFS the tree in lexicographical order!  $\rightarrow O(n)$  time, space

Algorithm) Exact matching with Suffix Array:

- build the suffix array.
- binary search the array to find the L and R window.
- return all indices between L and R.

Runtime:  $\underbrace{T(\text{binary search})}_{\hookrightarrow O(\log(|T|) \cdot |P|)} + \underbrace{T(\text{enumerating between L \& R})}_{\hookrightarrow O(k)} \rightarrow \text{total } O(\log(|T|) \cdot |P| + k)$

Speedup: If  $A_T[L]$  and  $A_T[R]$  share a prefix, we can ignore it.

Even better: Gusfield proposes a  $O(|P|)$  comparison scheme!

↳ Practically, runtime looks like  $O(\log(|T|) + |P| + k)$ .

\* Can we construct suffix arrays without suffix trees?

↳ Yes, KS Algorithm [2006] is linear, using mergesort & radixsort!

Algorithm) Radix Sort: set of  $n$  strings of length  $\leq k \rightarrow O(nk)$  time

↳ Sort by the first character, put ties into bins, recurse on them.

⇒ Naively, radix sort gives  $O(n^2)$  for suffix arrays ( $\because k=n$ ).

Intuition) What if we can sort only the first  $k$  chars of suffixes?

↳ If  $k$  is a constant, radix sort will run in  $O(k|\pi|) = O(|\pi|)$ .

What if there are ties? → consult results from  $(i+1)$  vs  $(j+1)$  or

$(i+2)$  vs  $(j+2)$  or ... or  $(i+k)$  vs  $(j+k)$ !

Naive Attempt: sort 3-suffixes via radix sort. Tiebreak by looking at the lexicographic order of the next 3 subsequent suffixes. In fact, replace all suffixes with lexicographic order of 3 subsequent suffixes.

Runtime:  $O(n^2)$  since  $n$  iteration takes  $O(n)$  time.

Intuition) Divide & Conquer: create 3 copies of the string, starting from positions 1, 2, and 3. Order two out of three using recursive radix sort, then efficiently sort & merge with the remaining third.

\* Details for KS is omitted in this note.

Suffix Array as a Permutation:  $\exists$  an inverse Suffix Array  $A_s^{-1}$ !  
↳ KS Algorithm actually wants to construct  $A_s^{-1}$  instead of  $A_s$

## Burrows-Wheeler Transform

Motivation) Efficient compression with invertability

Naïve compression: count the # of repeated characters

↳ works well with low-entropy strings, not so well with higher ones.

Lexicographical sorting: efficient, but not invertible.

BWT: list of characters that precedes each row in a SA.

Intuition: Strings often have patterns, i.e. some sequences are repeated

↳ Thus, BWT is likely to contain long repeated characters!

e.g. 'h' is likely to be preceded with 't', so repeated 't's are expected

Construction of BWT: Consider the cyclic rotations of  $S$ . Sorting these yields a suffix array with nice properties containing BWT.

Def) BWT:  $\pi(S) :=$  rotations of  $S$ ,  $\pi^s(S) :=$  lex. sorted list of  $\pi(S)$

$F, L :=$  first & last columns of  $\pi^s(S)$ .

which can be reconstructed from  $L$  as with  $S$ .

↳  $F$  is simply the ordered list of characters, and  $L$  is the BWT.  
( $\$$  pointer to  $\$$ )

\* Using KS algorithm, BWT can be constructed in  $O(n)$ .

Claim: we can unambiguously reconstruct S using only L (and F).

↳ In other words, BWT is invertible.

Key Obs: F and L are order-preserving!

Proof of Claim: all ties are broken by their next permutations, which are exactly the strings that end with that character! Then, since the SA is already sorted, appearance in L is order-preserved.

ex) A\$BARBAR, < ARAS\$BARB  $\Leftrightarrow$  \$BARBARAA < RA\$BARBAA

Reconstruction from L: Obviously  $F[1]$  is \$, and the last character in S is also \$. Then, we can chase down its preceding character,  $L[i]$ . Next, find the first occurrence of  $L[i]$  in F, and let its index be j. Chase down the next one,  $L[j]$ .

Continue finding the n-th occurrence of  $L[i]$  in F, j.  $L[j]$  is the next preceding character, and update  $i \leftarrow j$ .

Efficiency: to find the m-th copy of character k, keep M, the first index for each character in F & C, the # of occurrences of  $L[r]$  in string  $L[1:r-1]$ .  $\Rightarrow M[L[r]] + C[r]$  is the index of r in F!

\* There also exists a "forward decoding" that is equivalent.

↳ Intuitively, we chase down characters in F instead of L!

## Querying BWT

The Perfect Match Problem: Is BWT any useful for this?

↳ yes, BWT compactly stores the information of the suffix array!

Idea: Scan the query string backwards and update pointers

$S_p$  &  $E_p$  that points to the first & last matching rotations in  $\pi^s(T)$ .

Update Rules: For a query of length 1, consult M

↳  $S_p(c) \leftarrow M[c]$ ,  $E_p(c) \leftarrow M[\text{next}(c)] - 1$ .

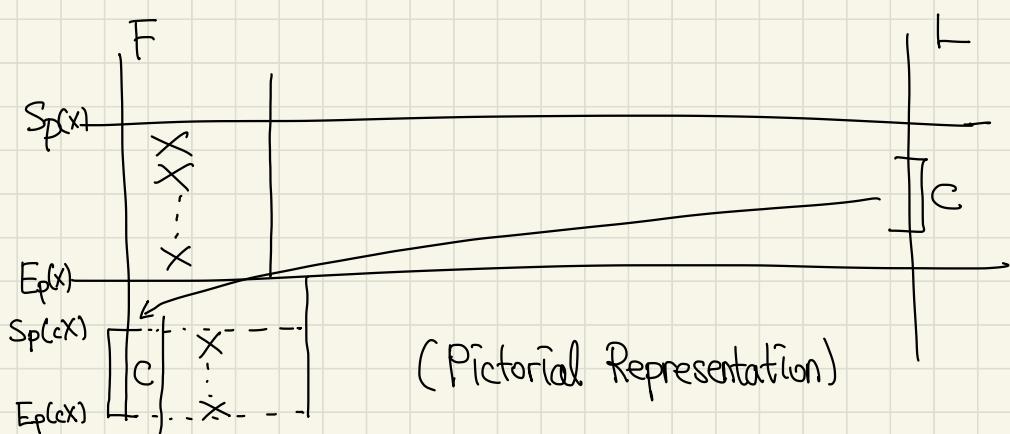
Then, look at  $L[S_p]$  and  $L[E_p]$  for their matching occurrence

In F again using M and OCC, where  $OCC[c, i] := \#$  of occurrences of character c in the first i positions in L.

↳  $S_p(cX) \leftarrow M[c] + OCC[c, S_p(X)-1]$

$E_p(cX) \leftarrow M[c] + OCC[c, E_p(X)] - 1$

Why does this work? OCC increases on the interval  $[S_p(X), E_p(X)]$  as long as  $L[i] = c$ , the next character we are chasing!



→ This solves half of the problem: the # of occurrences of  $P$  in  $T$ .

↪ If  $P$  is exhausted, output  $(E_p - S_p + 1)$ , if  $E_p < S_p$  output NULL.

Claim: If  $X$  occurs in  $T$ , then  $cX$  occurs in  $T$  if  $S_p(cX) \leq E_p(cX)$ .

Proof:  $OCC[c, E_p(X)] = OCC[c, S_p(X) - 1]$  ( $c$  never occurs b/wn  $[S_p, E_p]$ )

$$\Rightarrow M[c] + OCC[c, E_p(X)] - 1 \leq M[c] + OCC[c, S_p(X) - 1]$$

$$\Rightarrow E_p(cX) < S_p(cX).$$

Runtime:  $O(|P|)$ , but still need to find locations of matching

(also  $OCC$  is very demanding in space)

Locating Matches: Simply return  $\{A_T[S_p], \dots, A_T[E_p]\}$

... but this defeats the purpose of not storing  $A_T$ !

Idea: Only store some values of  $A_T$ , then recover as necessary.

Sparse Storage of  $A_T$ : evenly store a small fraction ( $\sim 2\%$ ) of  $A_T$ .

- ↳ If  $A_T[i]$  is not found, consult  $j = M[C] + OCC[C, i-1]$ , which is the index of  $cX$  if  $i$  corresponds to  $X_c$ . Then, if  $A_T[j]$  is found, return  $A_T[j]+1$  for  $A_T[i]$ . Else, recurse on  $i \leftarrow j$ .
- ⇒ This introduces a time vs space tradeoff on how much to store!

Next Idea: Approximately matching strings!

## Inexact String Matching

Why? Sequencing Errors & Genetic Variation

Given two strings: overall level of similarity, which parts are similar?

Pairwise Global Alignment ⇒ Inserting gaps to enforce matching

Events: Substitution ( $\begin{smallmatrix} A \\ c \end{smallmatrix}$ ), Insertion ( $\underline{C}$ ), Deletion ( $\overline{C}$ ) ... Descendant Ancestor

↳ Between siblings, insertion/deletion is called indels

Scoring Alignment: Total Score =  $\sum$ (score of aligned pair) +  $\sum$ (score of indel)

↳ amino acids  $C_1 \leftrightarrow C_2$  will have a higher substitution score if that substitution happened frequently in evolution (assumption)

Score :=  $\log\left(\frac{\Pr[C_1 \rightarrow C_2 | R]}{\Pr[C_1 \rightarrow C_2 | U]}\right)$  where  $R :=$  related,  $U :=$  unrelated

=  $\log\left(\frac{p_{C_1, C_2}}{q_{C_1} \cdot q_{C_2}}\right)$  where  $p_{C_1, C_2}$  := frequencies of  $C_1 \& C_2$  both in a same index,

$q_{C_1}, q_{C_2}$  := marginal probabilities of  $C_1$  and  $C_2$

=  $\sigma(C_1, C_2) \propto$  amount of expected  $C_1 \rightarrow C_2$  mutations

$\Rightarrow \text{Score} = \sum_{\text{aligned}(C_1, C_2)} \sigma(C_1, C_2) + \sum(\text{indel score})$

Indel scoring: Linear  $\rightarrow$  treat every indel independently: ( $W_s \times \# \text{ of indels}$ )

Affine  $\rightarrow$  treat consecutive indels as a single event:

$(\sum_{\text{gaps}} (W_g + W_s \cdot (\# \text{ of indels in the gap}))$  where  $(W_g / \text{gap penalty}) \gg |W_s|$ .

$\Rightarrow \text{Score} = \sum_{\text{aligned}(C_1, C_2)} \sigma(C_1, C_2) + (\# \text{ gaps}) \times W_s = \sum_{C_1, C_2} \sigma(C_1, C_2)$  where  $\sigma(C_1, -) = W_s$ .

Global alignment: Given strings  $x, y$  and  $\sigma, W_s$ , find maximal alignment.

$\hookrightarrow$  Use DP through defining appropriate subproblems!

Claim: If  $V_{ij}$  is the value of optimal alignment for  $x[:i], y[:j]$ , then

$V_{ij}$  can be derived trivially from  $\{V_{i-1,j}, V_{i,j-1}, V_{i-1,j-1}\}$ .  $\rightarrow$  find  $V_{nm}$ !

Algorithm) Needleman-Wunsch:  $x, y, \sigma, W_s \rightarrow V(n, m)$  ( $|x|=n, |y|=m$ )

Base case:  $V(0,0) = 0$ , for  $0 < i \leq n$ :  $V(i, 0) = \sum_{k \in [i]} \sigma(x[k], -) = i \cdot W_s$

$V(0, j) = \sum_{k \in [j]} \sigma(-, y[k]) = j \cdot W_s$ .

Recursion:  $V(i, j)$  can end with 1)  $(x[i], y[j])$  2)  $(x[i], -)$  3)  $(-, y[j])$

Then, we need information from 1)  $V(i-1, j-1)$  2)  $V(i-1, j)$  3)  $V(i, j-1)$  to calculate each possible  $V_{ij}$  score.

Take the max of:  $\{V(i-1, j-1) + \sigma(x_i, y_j), V(i-1, j) + \sigma(x_i, -), V(i, j-1) + \sigma(-, y_j)\}^{=ws}$

→ Fill the  $n \times m$  matrix, then back track pointers (path not unique!)

Runtime: if  $n = O(m)$ , total  $O(n^2)$  time.

Space: if  $n = O(m)$ , naively  $O(n^2)$ , can reuse one row for  $O(m)$  space.

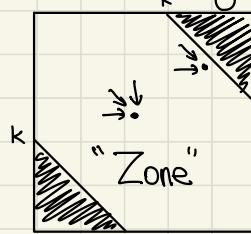
↳ can also do a diagonal zig-zag fill out, which is parallelizable!

## Variations on String Matching I

1) K-bounded Alignment: allow at most  $k$  mismatches (sub & indels)

Assume  $\sigma(n, m) = \begin{cases} 0 & \text{if } n=m \\ 1 & \text{if } n \neq m \text{ OR Gap.} \end{cases}$  Obviously,  $|n-m| \leq k$ .

Idea: only fill in the relevant portions of the matrix  $V$ .



Border cells ignore cells outside the zone.

However, a path in the zone can still have more than  $k$  alignment. → If  $V(n, m) \leq k$ , trace back.

Otherwise, there is no  $k$ -bounded alignment (iff relation!)

$$\text{Area(Zone)} \approx nm - 2 \cdot \frac{(n-k)(m-k)}{2} = k(m+n) - k^2 \rightarrow O(m+n) = \underline{O(n) \text{ time!}}$$

2) Affine Gap Penalty: treat consecutive indels as single events

$$\rightarrow \text{Score} = \sum_{(c_1, c_2)} \delta(c_1, c_2) + \sum_{x: \text{gap}} W_{\text{gap}} + |X| \cdot W_s.$$

↪ NIN algorithm no longer applies since we introduce dependencies!  
(We might want to go back and fix previous indices now)

Naïve Algorithm: Consider all possible consecutive gaps  $\rightarrow O(n^3)$   
matrices

Idea: keep track of three possibilities for preceding solutions

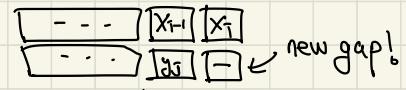
$M(i, j)$  := best alignment of  $X[1:i], Y[1:j]$  ending with  $x_i$  &  $y_j$  match

$I_x(i, j)$  := // ending with  $x_i$  & gap

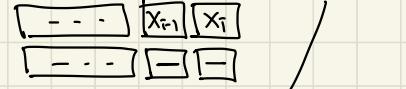
$I_y(i, j)$  := // ending with gap &  $y_j$

$$V(i, j) \leftarrow \max \{ M(i, j), I_x(i, j), I_y(i, j) \}.$$

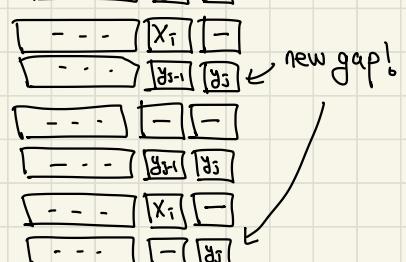
$$\Rightarrow I_x(i, j) = \max \begin{cases} M(i-1, j) + W_g + W_s \\ I_x(i-1, j) + W_s \\ I_y(i-1, j) + W_g + W_s \end{cases}$$



$$\Rightarrow I_y(i, j) = \max \begin{cases} M(i, j-1) + W_g + W_s \\ I_x(i, j-1) + W_g + W_s \\ I_y(i, j-1) + W_s \end{cases}$$



$$\Rightarrow M(i, j) = \max \begin{cases} M(i-1, j-1) + \delta(x_i, y_j) \end{cases}$$



$$\begin{cases} I_x(i-1, j-1) + \sigma(x_i, y_j) \Leftrightarrow V(i-1, j-1) + \sigma(x_i, y_j) \\ I_y(i-1, j-1) + \sigma(x_i, y_j) \end{cases}$$

Base Cases:  $M(0,0) = I_x(0,0) = I_y(0,0) = \emptyset$ .

$I_x(i,0) = Wg + i \cdot Ws$ ,  $I_y(i,0)$  &  $M(i,0)$  are not defined. ( $0 < i \leq n$ )

$I_y(0,j) = Wg + j \cdot Ws$ ,  $I_x(0,j)$  &  $M(0,j)$  are not defined. ( $0 < j \leq m$ )

## Variations on String Matching II

1) Local Alignment: alignment of substrings of  $X$  and  $Y$

↪ find the highest scoring alignment, considering all substrings of  $X$  &  $Y$

\* makes sense when  $X$  and  $Y$  are very distant or have very different length

Pairwise Local Alignment: outputs substrings  $\alpha$  and  $\beta$ , naïve solution

is to enumerate all  $O(n^2m^2)$  pairs of substrings  $\rightarrow O(n^3m^3)$  runtime

Intuition: {all suffixes of all prefixes} = {all substrings}

↪  $V_{ij} :=$  best alignment score between a suffix of  $X[i:-]$  and  $Y[1:j]$ .

$\Rightarrow$  best local alignment is  $\max_{i,j} \{V_{ij}\}$ . [Smith-Waterman]

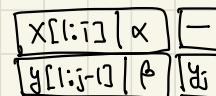
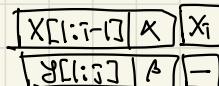
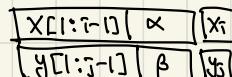
Recursion:

$$V(i,j) := \max$$

$$V(i-1, j-1) + \sigma(x_i, y_j)$$

$$V(i-1, j) + W_s \xrightarrow{\sigma(x_i, -)}$$

$$V(i, j-1) + W_s \xrightarrow{\sigma(-, y_j)}$$



$\emptyset$  (is selected when all above are negative)

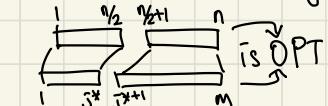
Base Case:  $V(0,0) = V(i,0) = V(0,j) = \emptyset$ .

\* After taking  $\max_{i,j} \{V_{ij}\}$ , back track until we hit a  $\emptyset$  (restart point).

2) Linear Space Global Alignment: How to recover with only the score?

Idea: Record "midway" points, then recursively continue. [Hirschberg]

↪ How to find where it passes at  $i = \frac{n}{2}$ ?



Intuition:  $\exists j^*, V(X[1:n], Y[1:m]) = V(X[1:\frac{n}{2}], Y[1:j^*]) + V(X[\frac{n}{2}:n], Y[j^*:m])$

↪ in other words,  $j^* = \arg\max_j \{ V_1(j) + V_2(j) \}$ .

→ Compute NW with  $X[1:\frac{n}{2}]$  and  $Y \rightarrow$  last row gives  $V_1(j) \downarrow$

For the bottom half, find NW of  $X[n:\frac{n}{2}+1]$  and  $Y^R \rightarrow V_2(m-j) \downarrow$   
 ↪ Reversed!

$\Rightarrow j^* = \arg\max_j \{ V_1(j) + V_2(j) \}$ , path passes through point  $(\frac{n}{2}, j^*)$ !

↪ For the rest of the path, recursively find smaller midpoints until M is small.

Space:  $O(mn) = O(n)$  (no saving entire matrix)

Time: First iteration  $\rightarrow O(nm)$



Second iteration  $\rightarrow O(\frac{nm}{2})$



Third iteration  $\rightarrow O(\frac{nm}{4})$



... i-th iteration  $\rightarrow O(\frac{nm}{2^{i-1}}) \Rightarrow$  Overall  $O(nm \sum_{i=0}^{\log(n)} \frac{1}{2^{i-1}}) = O(nm)$ , //

$\Rightarrow$  We trade some constant multiplicative factor in time to reduce space requirement from quadratic to linear!

## Phylogeny

Def) Phylogeny: a weighted tree capturing ancestral relationships

↳ leaves represent extant species, branching is a species divergence

↳ a rooted tree will encode a "flow of time" for mutations

Why Trees?  $\rightarrow$  we only consider "splitting" and not "joining" (not always...)

Applications: 1) Taxonomy 2) Ancestral Relations 3) Human Evolution

4) Evolution within an Individual (cancer mutations)

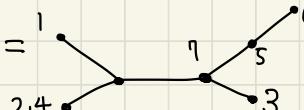
$\rightarrow$  max likelihood, max parsimony, ...

Two Algorithms: 1) Distance-Based 2) Character-Based

Distance-Based: Given a pairwise dissimilarity table, label a tree s.t. path lengths correspond to the dissimilarity

Parsimony-Based: Find a tree that require the minimum # of changes

Def)  $X$ -Tree:  $T := (T, \varphi)$  where  $\varphi(X)$  contains all vertices with degree  $\leq 2$ , and the label set  $X$  ( $\varphi$  is not one-to-one nor onto)

ex)  $T(T, \varphi) =$   ,  $X = [7]$  \* degree 2 vertex only is introduced to put a label there!

Def) Phylogenetic Tree:  $X$ -Tree where  $\varphi$  is a bijection of  $X \leftrightarrow T$ .

Theorem) [Schröder]  $B(n) := \{ \text{all inequivalent unrooted binary phylogenetic trees with } n \text{ leaves} \}$ , then  $|B(n)|$  is super exponential.

Def) Dissimilarity Map: function  $\delta: X \times X \rightarrow \mathbb{R}$  s.t.  $\delta(x, x) = 0$ ,  $\delta(x, y) = \delta(y, x)$ .  
 \* DM is pseudometric if it satisfies the triangle inequality, and it is metric if  $\delta(x, y) = 0 \iff x = y$ .

↪ How can  $\delta$  be "represented" by a tree, and how do we find it?

⇒ Make  $T$  with edge weight  $w(e)$  s.t.  $d_{(T, w)}(u, v) := \sum_{e \in \text{path}(u, v)} w(e)$ .

Def) Tree Metric:  $\delta$  on  $X$  s.t.  $\exists T, w: E \rightarrow \mathbb{R}_{>0}$  s.t.  $\delta(x, y) = d_{(T, w)}(\varphi(x), \varphi(y))$

↳ What is the necessary & sufficient condition?

Def) 4 Point Condition:  $\forall a, b, c, d \in X, \delta(a, b) + \delta(c, d) \leq$

$$\max \{ \delta(a, d) + \delta(b, c), \delta(a, c) + \delta(b, d) \}$$

\* 4PC implies TI,  $\delta$ 's nonnegativity, and  $\delta$ 's pseudometricness.

Theorem) (Tree-Metric) For a  $\delta \geq 0$  on  $X$ , 4PC  $\iff$  Tree-Metric

↗ two are equal and  $\geq$  than the third

Def) Ultrametric:  $\forall a, b, c \in X, \delta(a, b) \leq \max \{ \delta(a, c), \delta(b, c) \}$ .

\* Ultrametric  $\Rightarrow$  Tree-Metric, and it arises naturally in phylogenetic  $X$ -trees.

Def) Equidistant Representation: A phylogenetic  $X$ -tree  $T(T, \varphi)$

rooted at  $\rho$  and an edge weight  $w: E(T) \rightarrow \mathbb{R}$  s.t.:

$$1) d_{(T,w)}(\rho, \varphi(x)) = d_{(T,w)}(\rho, \varphi(y)) \quad \forall x, y \in X,$$

$$2) w(e) \geq 0 \quad \forall e \in E^o(T) \text{ where } E^o(T) := \text{set of interior edges of } T.$$

↳ Useful to think that current species have roughly the same # of mutations (distance) when compared with a common ancestor (root node).

Theorem) If  $\delta$  is ultrametric on  $X$ , then  $\exists$  an equidistant representation  $(T, w)$  such that  $\delta = d_{(T,w)}$ .



Def) Gromov Product: For a  $\delta$  on  $X$ , fix an element  $r \in X$ . The Gromov product of  $x$  and  $y$  in  $X \setminus \{r\}$  is:

$$\delta_r(x, y) = \frac{1}{2} [\delta(x, y) - \delta(x, r) - \delta(r, y)] \text{ if } x \neq y, 0 \text{ otherwise.}$$

Theorem) A non-negative  $\delta$  on  $X$  satisfies FPC  $\Leftrightarrow \delta_r$  is an ultrametric on  $X \setminus \{r\}$   $\forall r \in X$ . (Proof in HW, key to proving Tree Metric Theorem!)

Phylogeny Reconstruction Problem: Given  $\delta$  on  $X$ , find a phylogenetic  $X$ -tree and weight function  $w: E \rightarrow \mathbb{R}_>$  that minimizes:

$$L(d_{(T,w)}, \delta) := \sum_{a,b \in X} g_{ab} [\delta(a, b) - d_{(T,w)}(a, b)]^2. \quad (g_{ab} \text{ are scaling factors})$$

$\hookrightarrow$  If  $T(T, \psi)$  is given, optimizing  $w$  is in P. However, if  $T$  is unknown, this problem is NP-Complete.

## Distance-Based Phylogeny

"Given a dissimilarity matrix, output a tree that is as consistent as possible."

Two greedy algorithms: 1) UPGMA (ultrametric) 2) Neighbor Joining (tree metric)

Vertex Representation of Ultrametric: Let  $h: V(T) \rightarrow \mathbb{R}$  be a height function s.t.  $h(u) \geq h(v)$  if  $u$  is an ancestor of  $v$ , and  $d_{(T,h)} := \begin{cases} 2h(\text{LCA}_T(\varphi(x), \varphi(y))) & \text{if } x \neq y \\ 0 & \text{if } x = y \end{cases}$ .

Algo) UPGMA: Unweighted Pair Group Method with Arithmetic Mean

↪ Input:  $S$  on  $X$ , Output: Phylogenetic X-tree  $T$  & height function  $h$ .

\* Always returns an equidistant (ultrametric) tree.

Intuition: Find the "closest" pair and merge them, recursively. Then, we need to update  $S$  based on the new supervertex (collapsing).

Algorithm:  $\forall i \in X$ , assign cluster  $C_i := \{i\}$ . Initialize  $\Delta(C_i, C_j) := S(i, j)$ .

Define a leaf node  $\emptyset T$  and assign a height of  $0$ . Iterate:

Find  $\min\{\Delta(C_i, C_j)\}$ . Merge them into  $C_m := C_i \cup C_j$ . Assign a height of  $\Delta(C_i, C_j)/2$  for a new node  $C_m$ . Update dissimilarity by

$$\Delta(C_i \cup C_j, C_k) := \frac{|C_i|}{|C_i| + |C_j|} \Delta(C_i, C_k) + \frac{|C_j|}{|C_i| + |C_j|} \Delta(C_j, C_k).$$

$$(* \Delta(C_i, C_j) := \frac{1}{|C_i| \times |C_j|} \sum_{x \in C_i, y \in C_j} S(x, y) \text{ originally})$$

Runtime:  $O(n^3)$  ( $n$  iterations,  $n^2$  to find the min-pair at each step)

↪ Can reduce down to  $O(n^2)$  with optimization.

But what if  $S$  is not ultrametric, i.e. one edge of a sibling is longer?

↪ UPGMA might merge the wrong nodes!

Algo) Neighbor Joining: Return the correct tree if  $S$  is tree metric.

Intuition: Find clusters to merge like UPGMA, but the update rule is different!

Pair Selection:  $\min_{(i,j)} \{Q(i,j) := d(i,j) - (r_i + r_j)\}$  where  $r_i := \frac{1}{|L|-2} \sum_{k \in L} d(i,k)$ .

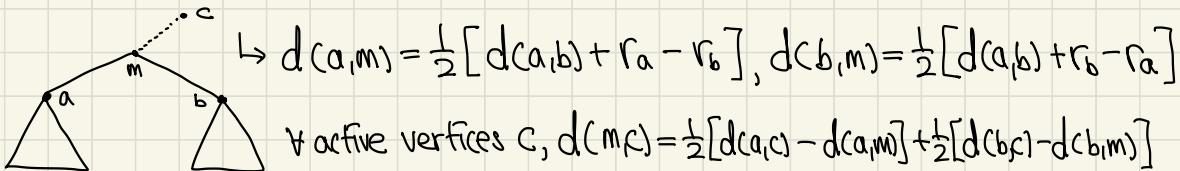
$\hookrightarrow Q(i,j) \approx$  distance btwn  $(i,j)$  minus mean distance to all other active leaves.

Theorem) If  $S$  is tree metric, then  $\arg\min_{(i,j)} \{Q(i,j)\}$  is guaranteed to be neighbors.

$\hookrightarrow$  Justifies Neighbor Joining, proof in HW. (intuitively, confidence boosting)

Update: Leverage the tree metric property to derive the rule:

$$d(c,m) = (d(c,b) + d(c,a) - d(a,b)) / 2 \text{ when } m \text{ is a merge of } a \text{ and } b.$$



$\hookrightarrow$  intuitively, take the average of the two estimators

Runtime:  $O(n^3)$ .

## Character-Based Phylogeny

A character on  $X$  captures attributes of the sample. ex) DNA seq, alignment

Def) Character: function  $X: X' \rightarrow S$  where  $X' \subseteq X$  and  $S$  is a set of states.

$\hookrightarrow X'$  is defined in this way to represent missing data.

	$x_1$	$x_2$	...
1	A	A	
2	T	C	
3	*	C	

A character  $X$  is said to be:

1) Full if  $X' = X$ . 2)  $k$ -state if  $|X(X')| = k$  (# of state outcomes)

3) Trivial if there exists at most one state with multiplicity  $\geq 2$ .

\* Nontrivial characters imply information about tree topology!

Def) Extension: for  $T(T, \psi)$ , an extension  $\bar{X}: V \rightarrow S$  of  $X$  on  $T$  has  
 $\bar{X} \circ \psi(a) = X(a) \quad \forall a \in X$ . (i.e. interior vertices also gets a state  $\in S$ )

Parsimony: Given a collection  $C$  of characters on  $X$ , provide a tree that gives the simplest explanation, i.e. minimizes the number of mutations.

Def) Changing Number: given  $T \& \bar{X}$ ,  $ch(\bar{X}, T) := |\{(u, v) \in E(T) | \bar{X}(u) \neq \bar{X}(v)\}|$ .

Def) Parsimony Score:  $\pi(X, T) = \min_{\bar{X} \in \text{ext}(T, X)} \{ch(\bar{X}, T)\}$ .

Small Parsimony Problem: given a rooted Phy. X-tree  $T(T, \psi)$  & a full character  $X$  on  $T$ , find an extension  $\bar{X}$  of  $X$  minimizing  $ch(\bar{X}, T)$ .

Algorithm) Fitch's Algorithm: if leaf  $v$ ,  $\psi(v) := \{X(\psi^{-1}(v))\}$ . if interior vertex  $v$  with children  $u, w$ , let  $\psi(v) := \begin{cases} \psi(u) \cap \psi(w) & \text{if } \psi(u) \cap \psi(w) \neq \emptyset \\ \psi(u) \cup \psi(w) & \text{if } \psi(u) \cap \psi(w) = \emptyset \end{cases}$ . Then,  $\pi(X, T)$  is the number of union operations called. To obtain an

optimal extension, pick an arbitrary state in  $\psi(\text{root})$ , then for a parent  $p$  with state  $s_p$ , the child  $c$ 's state will be  $s_c \leftarrow \begin{cases} s_p & \text{if } s_p \in \psi(c) \\ s_{\pi(\psi(c))} & \text{otherwise.} \end{cases}$

Runtime:  $O(kn)$  where  $k = |S|$ ,  $n = |X|$  ( $\forall v \in V$ , takes  $O(k)$  time to find  $\psi(v)$ ).

Weighted Small Parsimony:  $\forall i, j \in S$ , a mutation  $i \rightarrow j$  costs  $m(i, j)$ .

Algorithm) Sankoff's Algorithm:  $\forall v \in V$ ,  $A_v :=$  an array indexed by  $S$  such that

$A_v[S] :=$  OPT cost for the subtree rooted at  $v$  with state  $s \in S$ . 

Base Case:  $\forall$  leaf  $l$ ,  $A_l[S] = \begin{cases} 0 & \text{if } s = \text{state}(l) \\ \infty & \text{otherwise} \end{cases}$   $\hookrightarrow$  penalize as much as possible.

Recursion:  $\forall v$  with children  $u, w$ ,  $A_v[i] := \min_{j \in S} [A_u[j] + m(i, j)] + \min_{k \in S} [A_w[k] + m(i, k)]$

$\hookrightarrow$  keep a pointer for  $\arg \min_{j \in S}$  and  $\arg \min_{k \in S}$ . Then,  $\Pi(X, T) = \min_{S \in S} \{A_{\text{root}}[S]\}$ .

To obtain the optimal extension, find  $s_c = \arg \min_{j \in S} \{A_c[j] + m(s_p, j)\}$ .

Runtime:  $O(k^2 n)$ .

Key Idea: There is a fundamental equivalence between  $X$ -trees and a certain type of collection of binary characters on  $X$ .

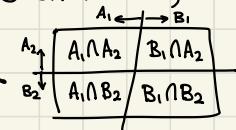
Def)  $X$ -Splits: a partition of  $X$  into two nonempty subsets ( $A \sqcup B$ ), that is, 1)  $A \cup B = X$ , 2)  $A \cap B = \emptyset$ , 3)  $A \neq \emptyset, B \neq \emptyset$ .

$\hookrightarrow$  Removing an edge  $e$  introduces an  $X$ -split (a cut).

Claim: For a given X-tree, no two different edges induce the same X-split.

Proof: If so, we will have an unlabeled degree 2 vertex, which is invalid.

Def) Compatibility: two X-splits  $A_1|B_1$  and  $A_2|B_2$  if at least one of  $A_1 \cap A_2, A_1 \cap B_2, B_1 \cap A_2, B_1 \cap B_2$  are empty. Otherwise, they are incompatible.  $\rightarrow$  all four "quadrants" are nonempty.



Def) 4-Gamete Test: If all four combinations  $\{00, 01, 10, 11\}$  appear in the two columns corresponding to a given pair of binary characters, they are incompatible.

Theorem) Splits-Equivalents [Buneman, 1971]: Let  $C$  be a collection of X-splits. Then, there exists an X-tree  $T(T, \varphi)$  s.t.  $\sum(T) = C$  iff  
 $\hookleftarrow$  = full binary characters  
the splits in  $C$  are pairwise compatible. Furthermore, the X-tree is unique up to isomorphism. ( $\sum(T) = \{X\text{-split}(e) \mid e \in E(T)\}$ )

Def) Convexity: a character  $X: X' \rightarrow S$  on  $T(T, \varphi)$  s.t.  $\exists$  an extension  $\bar{X}: V \rightarrow S$  of  $X$  s.t.  $\forall s \in S$ , the subgraph  $T := \{v \in V \mid \bar{X}(v) = s\}$  is connected, i.e. we can label inner vertices such that each subgraph consisted of the same state is connected.

Remark: a binary character  $X$  on  $T$  if  $\exists$  a bipartition of  $X$  induced by  $X$  is an  $X$ -split induced by an edge of  $T$ . (and that edge is unique)

→ Biologically, homoplasy-free (no recurrent or back mutations)

Def) Compatible: a collection  $C$  of characters on  $X$  s.t.  $\exists$  an  $X$ -tree on which all of the characters in  $C$  are convex. (generalized to  $k$ -states)

Def) Perfect Phylogeny Problem: Determine whether  $C$  is compatible, and if so, find the  $X$ -tree on which they are convex.

↪ For the binary case,  $C$  is compatible iff pairwise compatible!

Naive Algorithm: If  $|X|=n$ ,  $|C|=m$ ,  $O(nm^2)$  for checking all pairs.

Gusfield's  $O(nm)$  Algorithm: Assume that root = all-zero string.

- 1) View each column as an  $n$ -bit binary number, radix sort in decreasing order.
- 2) At row  $r$ , let  $k_r :=$  string of labels of columns with 1s in row  $r + \$$  at the end.
- 3) Build a "keyword tree" for strings  $k_1, k_2, \dots, k_n$ , akin to suffix trees.
- 4) Check whether each column label appears exactly once. If so, the tree is correct upto permutation of labels on the same edge.

Why does this work? → For each character  $X_i$ , let  $D_i := \{x \in X \mid X_i(x) = i\}$ .

If  $\exists$  a perfect phylogeny solution, then after radix sorting, if character  $X_i$  appears to the left of  $X_j$ , then either  $O_j \subseteq O_i$  or  $O_j \cap O_i = \emptyset$ .

For the unknown root case, take majority string as the root, and swap  $1 \rightarrow 0$ .

$\hookrightarrow$  If data is consistent, the majority string will appear somewhere in the tree!

For  $k \geq 3$ , we check for compatibility via convexity (pairwise not sufficient)

$\hookrightarrow$  Necessary & Sufficient condition is a chordal completion of  $\text{Int}(C)$ , the partition intersection graph  $((x_a, B_{ai}), (x_b, B_{bj})) \in E(\text{Int}(C))$  if  $B_{ai} \cap B_{bj} \neq \emptyset$ .

$\Rightarrow$  full characters run in  $O(2^k n m^2)$  time, any partial is NP-C.

## Probability

Probability Space:  $(\Omega, \mathcal{F})$ , sample space and  $\sigma$ -field (or  $\sigma$ -algebra).

Probability Measure:  $P$ , a function  $\mathcal{F} \rightarrow [0, 1]$  with such conditions:

1)  $P[\emptyset] = 0$ . 2)  $P[\Omega] = 1$ . (normalization)

3)  $A_1, A_2, \dots \in \mathcal{F}$ .  $A_i \cap A_j = \emptyset \forall i \neq j$ . Then,  $P[\bigcup_{i=1}^{\infty} A_i] = \sum_{i=1}^{\infty} P[A_i]$  (additive)

Random Variable:  $X$ , a function  $\Omega \rightarrow \mathbb{R}$  s.t. " $X \leq x$ " :=  $\{w \in \Omega | X(w) \leq x\} \in \mathcal{F}$

Cumulative Distribution Function:  $F_X(x) := P[X \leq x]$ .

Convergence in Distribution: seq.  $X_1, X_2, \dots$  and  $X$ .  $X_n \xrightarrow{d} X$  as  $n \rightarrow \infty$  if

$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_x(x)$ ,  $\forall x \in C(F_x)$  i.e. points in  $\mathbb{R}$  where  $F_x$  is continuous.

Central Limit Theorem:  $X_1, X_2, \dots$  i.i.d with finite mean  $\mu$ , finite variance  $\sigma^2$ .

Let  $S_n := X_1 + \dots + X_n$ .  $\frac{\sqrt{n}}{\sigma} \left( \frac{S_n}{n} - \mu \right) \xrightarrow{d} N(0, 1)$  as  $n \rightarrow \infty$ .

Conditional Probability: Let  $B$  be some event s.t.  $P[B] > 0$ .  $E_1, E_2$  are events.

We want to define a new prob. space  $(B, \mathcal{F}_B, P_B)$  s.t.  $P_B$  is consistent.

$$P[B \cap E_1] / P[B \cap E_2] = P_B[B \cap E_1] / P_B[B \cap E_2] \text{ (relative size)}$$

$$\Rightarrow P_B = cP \text{ where } c \text{ is a constant } \frac{1}{P[B]}. \Rightarrow P[A|B] := P_B[B \cap A] = \frac{P[B \cap A]}{P[B]}.$$

Independence:  $A \perp\!\!\!\perp B \Leftrightarrow P[A|B] = P[A]$ ,  $P[B|A] = P[B]$ .

$$\Rightarrow \text{If } A \perp\!\!\!\perp B, P[A \cap B] = P[A] P[B].$$

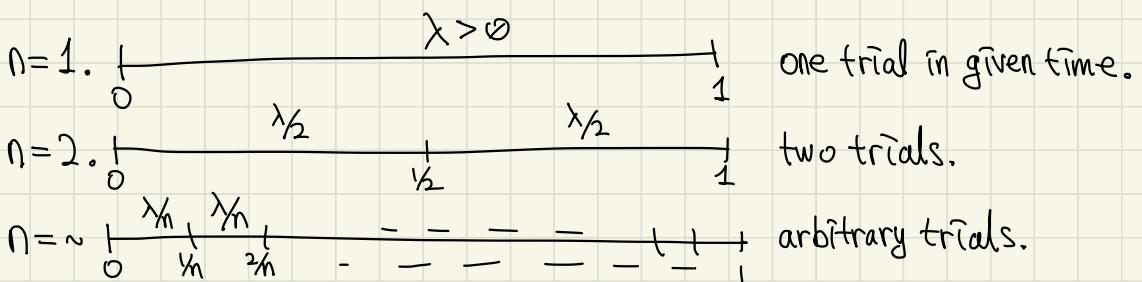
Theorem) Law of Total Probability: Let  $B_1, B_2, \dots \in \mathcal{F}$  be a partition of  $\Omega$  and  $A \in \mathcal{F}$ . Then,  $\sum_{i=1}^{\infty} P[A|B_i] P[B_i] = \sum_{i=1}^{\infty} P[A \cap B_i] = P[\bigcup_{i=1}^{\infty} (A \cap B_i)] = P[A]$  by additivity of  $P$ .

Stochastic Process:  $I \sim \text{Ber}(p)$  where  $P[I=0] = p$ ,  $P[I=1] = (1-p)$ .

$I_1, I_2, \dots$  i.i.d.  $\text{Ber}(p)$ .  $S_n := \# \text{ of } 1\text{s upto } n \text{ trials} = \sum_{i=1}^n I_i \sim \text{Bin}(n, p)$ .

$W_i := \text{waiting time of next } 1 \sim \text{Geo}(p)$ .  $P[W_i=k] = (1-p)^{k-1} p$ ,  $k \in \mathbb{N}^+$ .

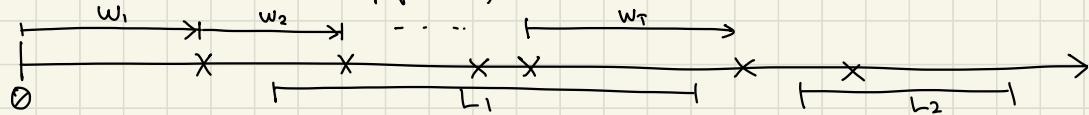
↪ How do we think about continuous processes?



$Y_n := \# \text{ of successes} \sim Bi(n, \lambda_n)$ .  $Y_1, Y_2, \dots, Y_n \xrightarrow{d} Po(\lambda)$  as  $n \rightarrow \infty$ .

Similarly,  $G_n \sim Geo(\frac{\lambda}{n})$ , then  $X_n := \frac{G_n}{n} \xrightarrow{d} Exp(\lambda)$  as  $n \rightarrow \infty$ .

Poisson Point Process:  $PPP(\lambda), \lambda > 0$ .



Each  $W_i \sim Exp(\lambda)$ .  $N_i := \# \text{ of hits in interval } L_i \sim Po(\lambda \cdot |L_i|)$ .

↳ This will be useful in describing mutation events!

## Markov Chains

(finite)  
(state space)

Def) Discrete Time Markov Chain:  $\{X_n | n=0, 1, \dots\}, X_n: \Omega \rightarrow S$ .

$$P[X_n=j | X_0=i_0, \dots, X_{n-1}=i_{n-1}] = P[X_n=j | X_{n-1}=i_{n-1}] \quad \forall n \in \mathbb{N}, \forall i_0, \dots, i_{n-1}, j \in S.$$

The one-step  $P[X_{n+1}=j | X_n=i] = p_{ij}$   $\forall n \in \mathbb{N}$  is homogeneous, and the elements can be represented as a matrix  $P[i, j] = p_{ij}$ .

Some Notations:  $\vec{u}_n = (P(X_n=i))_{i \in S}$ . How to relate this to  $\vec{u}_{n-1}$ ?

$$P(X_n=j) = \sum_{i \in S} P(X_n=j | X_{n-1}=i) \cdot P(X_{n-1}=i) = \sum_{i \in S} p_{ij} \cdot P(X_{n-1}=i)$$

$\hookrightarrow \vec{U}_n = \vec{U}_{n-1} P.$   $\rightarrow \vec{U}_n = \vec{U}_{n-2} P^2$ , and so on...  $\rightarrow \vec{U}_n = \vec{U}_0 P^n.$

Denote  $P^n[i, j] = P_{ij}^{(n)}.$

$$\text{Chapman-Kolmogorov: } P_{ij}^{(n+m)} = \sum_{k \in S} P_{ik}^{(n)} \cdot P_{kj}^{(m)} \text{ (total probability again)}$$

Def) Accessible:  $j$  is accessible from  $i$ , ( $i \rightarrow j$ ), if  $\exists n \in \mathbb{N}$  s.t.  $p_{ij}^{(n)} > 0.$

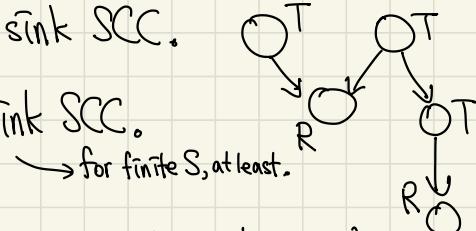
Def) Communicate:  $(i \leftrightarrow j)$  iff  $(i \rightarrow j) \wedge (j \rightarrow i).$

Def) Strongly Connected Component: Set of states s.t.  $\forall i, j \in S, (i \leftrightarrow j).$

$\hookrightarrow$  We can decompose states into SCCs & topologically sort them into DAGs.

Def) Transient Component: All states not in a sink SCC.

Def) Recurrent Component: All states in a sink SCC.



Def) First Passage Time:  $f_{ij}^{(n)} := P(X_1=j, \dots, X_{n-1}=j, X_n=j | X_0=i).$

Def) Return Probability:  $f_{ii} := \sum_{n=1}^{\infty} f_{ii}^{(n)}$  (every  $f_{ii}^{(n)}$  are disjoint!)

Obs) Transient  $\Rightarrow f_{ii} < 1.$  Recurrent  $\Rightarrow f_{ii} = 1.$  (actual definitions of these)

Def) Irreducibility: MC s.t.  $\forall i, j \in S, (i \leftrightarrow j)$  (one SCC for entire graph!)

Claim) Suppose  $j \in S$  is transient. Then,  $\lim_{n \rightarrow \infty} P_{ij}^{(n)} = 0 \quad \forall i \in S.$  (without proof here)

Decompose  $S = A \sqcup B$  where  $A$  is absorbing ( $\in$  recurrent SCC),  $B$  is transient.

For  $i \in A, j \in B$ , what is  $h_{ij} = P(\text{Enter } A \text{ through } j \mid X_0 = i)$ .

First Step Analysis:  $h_{ij} = \sum_{k \in S} P(\varepsilon_j \mid X_i = k, X_0 = i) \cdot \overbrace{P(X_i = k \mid X_0 = i)}^{p_{ik}}$

$$= \sum_{k \in S} P(\varepsilon_j \mid X_i = k) \cdot p_{ik}. P(\varepsilon_j \mid X_i = k) = \begin{cases} \delta_{jk}, & k \in A \\ h_{kj}, & k \in B. \end{cases} \Rightarrow h_{ij} = p_{ij} + \sum_{k \in B} h_{kj} p_{ik}.$$

We can decompose  $P$  into  $\begin{matrix} B & Q \\ A & R \end{matrix}$ . Then,  $H = R + QH \Rightarrow H = (I - Q)^{-1}R$ .

→ Of course,  $(I - Q)$  needs to be invertible, which it always is!

Lemma) Suppose a square matrix  $M$  satisfies  $\lim_{n \rightarrow \infty} M^n = \emptyset$ . Then,  $(I - M)^{-1}$  exists, and it is given by  $\sum_{n=0}^{\infty} M^n$ . (proof omitted)

↪ Then, because  $\lim_{n \rightarrow \infty} Q^n = \emptyset$ ,  $\underbrace{(I - Q)^{-1}}_{\text{Fundamental Matrix}}$  exists. //

Hitting Time:  $T_A = \left| \{ n \in \mathbb{N}_0 \mid X_n \in A \} \right| = \sum_{n=0}^{\infty} \sum_{j \in B} 1 \{ X_n = j \} \rightarrow E[T_A \mid X_0 = i]$   
 $= \sum_{j \in B} E \left[ \sum_{n=0}^{\infty} 1 \{ X_n = j \} \mid X_0 = i \right].$  (exact derivation in notes)

Recall that  $\vec{u}_n = \vec{u}_{n-1} P$ . If  $\pi = \pi P$ ,  $\pi$  is a stationary distribution.

That is,  $\pi$  is a left eigenvector with  $\lambda = 1$  and nonnegative elements.

If MC is finite,  $\exists$  a SD. If MC is irreducible,  $\pi$  is unique.

Def) Period of State: For  $i \in S$ ,  $d(i) = \gcd(\{ n \in \mathbb{N} \mid P_{ii}^{(n)} > 0 \})$ .

↪ MC is aperiodic if  $d(i) = 1$ .

Consider a MC that is finite, irreducible, & aperiodic.  $\lim_{n \rightarrow \infty} [P^n]_{ij} = \pi_j$ .

Def) Reversible MC:  $\pi_i P_{ij} = \pi_j P_{ji} \forall i, j \in S$ . (aka detailed balance)

↪ In fact,  $\pi_i [P^n]_{ij} = \pi_j [P^n]_{ji}$  holds by induction.

A simple model: Latent root  $Z \sim \pi$ ,  $n_1$  &  $n_2$  mutations for X & Y branches.

$$\begin{array}{c}
 \text{Z} \sim \pi \\
 \xrightarrow{\quad} \\
 \left. \begin{array}{c}
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \end{array} \right\} E \\
 \left. \begin{array}{c}
 \xrightarrow{n_1 \text{ mut}} \\
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \end{array} \right\} X=i \\
 \left. \begin{array}{c}
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \xrightarrow{\quad} \\
 \end{array} \right\} Y=j
 \end{array}
 \quad
 \begin{aligned}
 P(X=i \wedge Y=j | E) &= \sum_{k \in S} P(X=i \wedge Y=j | E, Z=k) P(Z=k | E) \\
 &= \sum_{k \in S} P(X=i | E, Z=k) P(Y=j | E, Z=k) \pi_k. \\
 &= \sum_{k \in S} [P^{n_1}]_{ki} [P^{n_2}]_{kj} \pi_k. \text{ Assuming reversibility,} \\
 \rightarrow \sum_{k \in S} [P^{n_1}]_{ki} \pi_i [P^{n_2}]_{jk} &= \pi_j [P^{n_1+n_2}]_{ji} = \pi_i [P^{n_1+n_2}]_{ij}.
 \end{aligned}$$

## Continuous Time Markov Chains

$$\text{PPP}(\mu), \mu > 0. \quad x_0 \xrightarrow{\quad} \cdots \xrightarrow{\quad} \xrightarrow{\quad} \{X_t | t \in \mathbb{R}_{\geq 0}\}.$$

Let  $N(t) := \# \text{ of mutations on the edge } \sim Po(\mu t)$ .

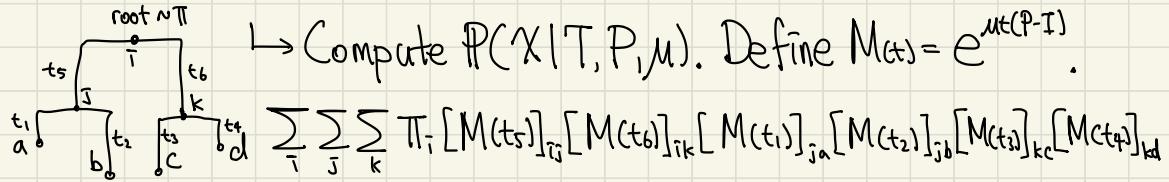
$$\hookrightarrow P(N(t)=k) = e^{-\mu t} \cdot \frac{(\mu t)^k}{k!} \quad \forall k = 0, 1, 2, \dots$$

$$\begin{aligned}
 P(X_t=j | X_0=i) &= \sum_{k=0}^{\infty} P(X_t=j | X_0=i, N(t)=k) \cdot P(N(t)=k | X_0=i) \\
 &= \sum_{k=0}^{\infty} [P^k]_{ij} \cdot e^{-\mu t} \frac{(\mu t)^k}{k!} = e^{-\mu t} \sum_{k=0}^{\infty} \frac{[\mu t P]^k}{k!} = e^{-\mu t} \cdot [e^{\mu t P}]_{ij}.
 \end{aligned}$$

Consider  $e^{-\mu t I} = \begin{cases} e^{-\mu t}, & i=j \\ 0, & i \neq j \end{cases} \rightarrow e^{-\mu t} [e^{\mu t P}]_{ij} = [e^{-\mu t I} e^{\mu t P}]_{ij}$ .

\* If  $[A, B] = 0$  ( $AB - BA = 0$ ) (commutes), then  $e^A e^B = e^{A+B}$ .  $[e^{\mu t(P-I)}]_{ij}$

A Tree Example: full character  $X$ , edge weighted phylo.  $X$ -tree  $T$ ,  $P$ ,  $M$ .



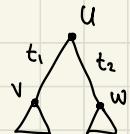
→ This gives  $O(|S|^3)$ , but with varying  $|X|=n$ ,  $O(|S|^{n+1}) \rightarrow$  intractable!

⇒ With DP, we can achieve  $O(n|S|^2)$ . Similar to Sankoff's Algorithm!

Consider a subtree rooted at  $v$  with children restricted to  $X$ .

→  $f_v(i) = \text{Prob. of observing children with } X \text{ given } v \text{ is state } i$ .

Initialization: For leaf nodes  $l$ ,  $f_l[i] = 1 \{ i = X_l \}$ .

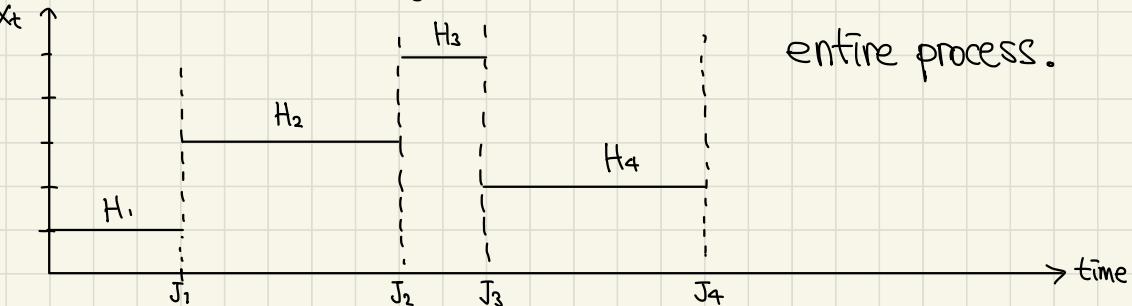


Iteration:  $f_u[i] = \left[ \sum_{j \in S} [M(t_1)]_{ij} f_v[j] \right] \times \left[ \sum_{k \in S} [M(t_2)]_{ik} f_w[k] \right]$  (cond.prob & LTP)

↪ Runtime:  $O(|S|^2)$  for every node,  $O(n)$  nodes  $\rightarrow O(n|S|^2)$ !

Termination:  $P(X|T, P, M) = \sum_{k \in S} \pi_k f_{\text{root}}[k]$ . (Felsenstein's Algorithm)

General CTMC:  $Q$ , a generator.  $Q = (q_{ij}) \forall i, j \in S$  describes the entire process.



$P(X_{u+t} = j | X_u = i) = P_{ij}(t)$ .  $\Rightarrow P(t) = (P_{ij}(t))$  satisfies  $\frac{dP(t)}{dt} = QP(t) = P(t)Q$ .

→  $P(t) = e^{tQ}$  is the general ODE solution.  $\Rightarrow = e^{Mt(P-I)}$

Properties of Q:  $\sum_{j \in S} q_{ij} = 0$ .  $q_{ii} = \sum_{j \neq i} q_{ij}$ ,  $q_{ij} | i \neq j > 0$ ,  $q_{ji} = -q_{ij}$ .

$H_i := \inf\{U | X(t+u) \neq i\} \sim \text{Exp}(q_i)$ . normalized  
✓

Embedded Jump Chain:  $\{Y_n = X_{J_n} | n \in \mathbb{N}_0\}$ .  $P(Y_{n+1} = j | Y_n = i) = \frac{q_{ij}}{q_{ii}}$ .

$P(Y_{n+1} = i | Y_n = i) = \begin{cases} 1 & \text{if } q_{ii} = 0. \xrightarrow{\text{degenerate case}} \\ 0 & \text{if } q_{ii} \neq 0. \end{cases}$  Finally,  $H_n \perp\!\!\!\perp Y_n$ . +o(h)

For  $0 < h \ll 1$ , we assume that  $p_{ij}(h) = q_{ij} \cdot h + o(h)$ ,  $p_{ii}(h) = 1 - q_{ii}h$ .

↳ This yields  $\frac{dP}{dt} = PQ = QP$  (KFE, KBE).

ex) PPP counter: states are enumerated, incremented at every arrival.

Q looks like  $\begin{bmatrix} -\lambda & \lambda & 0 \\ \lambda & -\lambda & \lambda \\ 0 & \lambda & \ddots \end{bmatrix}$  where diagonals are  $-\lambda$  and directly upper entries are  $\lambda$ .  $\rightarrow P'_0(t) = [PQ]_{00} = \sum_{k \in S} P_{0k} Q_{k0} = P_{00}(t) \cdot (-\lambda)$ .

Assume initial condition is  $P_{00}(0) = 1$ .  $\rightarrow P_{00}(t) = e^{-\lambda t}$ .  $\sim P_0(\lambda t)$ !

$P'_{0j}(t) = [PQ]_{0j} = \sum_{k \in S} P_{0k} Q_{kj} = (P_{0(j-1)}(t) - P_{0j}(t)) \lambda \rightarrow P_{0j}(t) = \frac{(\lambda t)^j}{j!} e^{-\lambda t}$ .

## Hidden Markov Models

Motivating Example: Occasionally Dishonest Casino that has fair/loaded dice.

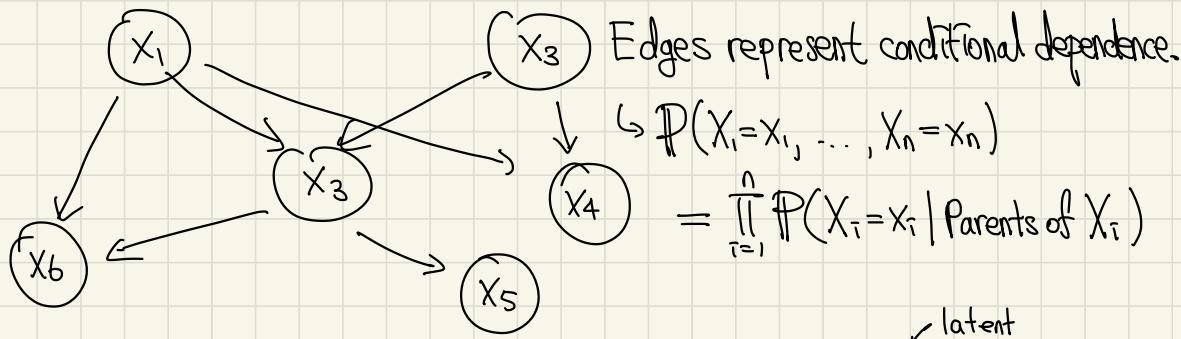
loaded probability  $x$ , initial fair dice Prob  $\pi$ , change prob  $a_{FL}, a_{LF}$ .

↳ After observing a sequence of rolls, what can we say about parameters?

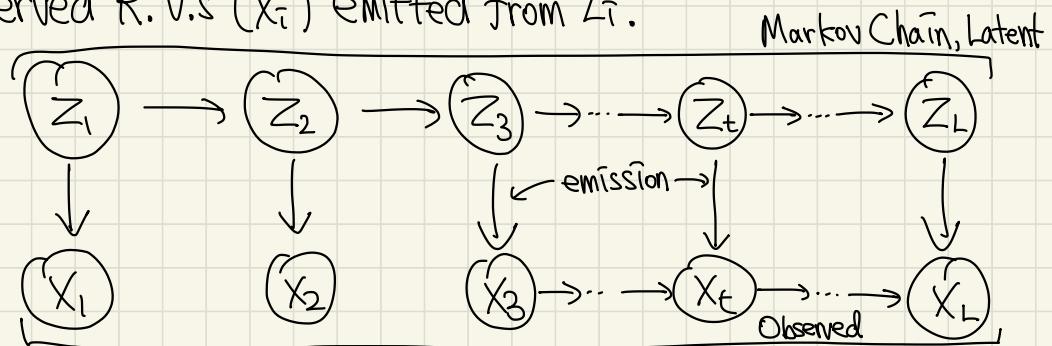
Concretely, what is  $P(Z_i = \text{fair} | \text{Data})$ ? (If we don't have the model, we can estimate it first)

Applications in Comp Bio: CpG finding, Coalescent fHMM → finding time from most recent common ancestor

Graphical Model : A DAG of R.V.s ( $X_i$ ).



HMMs are special graphical models governed by a MC( $Z_i$ ) and observed R.V.s ( $X_i$ ) emitted from  $Z_i$ .



→ This way,  $X_i$  will be conditionally independent of all other  $X_j$  given  $Z_i$ !

Denote the state space of  $Z_i \in S$ ,  $X_i \in A$  (some alphabet).

The HMM is specified by tuple  $(A, S, \Theta)$  where  $\Theta :=$  some parameters.

We don't know ground truth of  $\Theta$ , so we estimate  $\hat{\Theta}^* := \underset{\Theta}{\operatorname{argmax}} P_{\Theta}(X_{1:L} | Z_{1:L})$ .  

$$X_{1:L} = X_1, \dots, X_L$$

→ How to calculate  $P_{\Theta}(X_{1:L})$ ? → enumerating manually,  $\sum_{Z_{1:L}} P_{\Theta}(X_{1:L} | Z_{1:L}) P(Z_{1:L})$ .

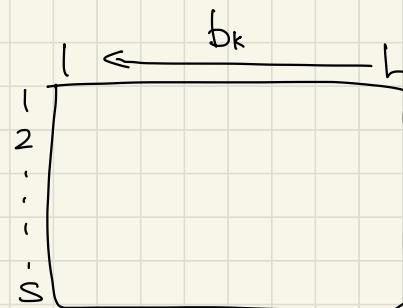
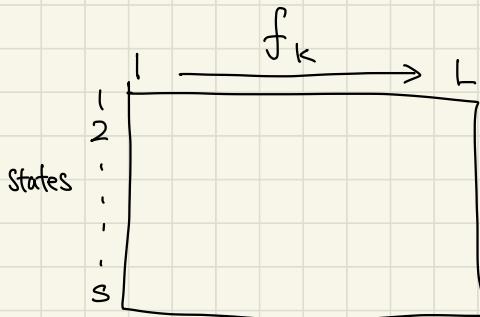
⇒ This sum involves  $O(|S|^L)$  terms, intractable!

$$\hookrightarrow \prod_{t=1}^L P_{\Theta}(X_t | Z_t) \text{ (emission)}$$

Inference on  $Z_t$ :  $P(Z_t=k | X_{1:L}) = \frac{P(X_{1:L}, Z_t=k)}{P(X_{1:L})}$ .

$$\underbrace{P(X_{1:t}, Z_t=k) P(X_{t+1:L} | X_{1:t}, Z_t=k)}_{=: f_k(t), \text{ forward probability}} \rightarrow \underbrace{P(X_{t+1:L} | Z_t=k)}_{=: b_k(t), \text{ backward probability}}, \text{ by cond. indep.}$$

Given  $f_k(t) \& b_k(t)$ ,  $P(Z_t=k | X_{1:L}) = \frac{f_k(t) \cdot b_k(t)}{P(X_{1:L})}$ .



Observe that  $P(X_{1:L}) = \sum_k P(X_{1:L}, Z_t=k)$  by marginalization.

Then, these values are exactly  $\sum_k f_k(L)$ ! (Using backward probs,

$$P(X_{1:L}) = \sum_k \underbrace{P(X_{1:L} | Z_1=k) P(Z_1=k)}_{\pi_k} \rightarrow \underbrace{P(X_{2:L} | Z_1=k) P(X_1 | Z_1=k)}_{\substack{b_k(1) \\ (X_2:L \perp\!\!\!\perp X_1 | Z_1)}} e_k(x_1)$$

$$= \sum_k b_k(1) \cdot e_k(x_1) \cdot \pi_k.$$

How to fill in  $f_k(t) \& b_k(t)$ ?  $\rightarrow$  Use DP! For the forward table:

Base Case:  $f_k(1) = P(X_1, Z_1=k) = P(X_1 | Z_1=k) P(Z_1=k) = e_k(x_1) \cdot \pi_k$ .

General Case:  $f_k(t) = P(X_{1:t}, Z_t=k) = \sum_l P(X_{1:t}, Z_{t-1}=k, Z_t=l).$  (cond. indep.)

$$\underbrace{P(X_{1:t}, Z_{t-1}=k) P(X_t, Z_t=l | X_{1:t-1}, Z_{t-1}=k)}_{f_k(t-1)} \rightarrow \underbrace{P(X_t, Z_t=l | Z_{t-1}=k)}_{e_k(x_t)}$$

$$\Rightarrow f_k(t) = e_k(x_t) \sum_{k \in S} f_k(t-1) \alpha_{kl} \cdot // \quad \underbrace{P(X_t | Z_{t-1}=k, Z_t=l) P(Z_t=l | Z_{t-1}=k)}_{e_k(x_t)} \quad \alpha_{kl} \text{ (given by the MC)}$$

Upshot: in the parameter  $\theta$ , we need  $\{a_{k,l} | k, l \in S\}$ ,  $\{e_k(\sigma) | k \in S, \sigma \in A\}$ , and  $\{\pi_k | k \in S\}$  to specify the entire table of  $f_k(t)$ .

Backward Table:  $b_k(t) = P(x_{t+1:L} | Z_t = k)$ .

Base Case:  $b_k(1) = 1 \quad \forall k \in S$ .

General Case:  $b_k(t) = P(x_{t+1:L} | Z_t = k) = \sum_{l \in S} \underbrace{P(x_{t+1:L} | Z_t = k, Z_{t+1} = l)}_{\downarrow a_{k,l}} \cdot \overbrace{a_{k,l}}$

Observe that  $(x_{t+1:L} \perp\!\!\!\perp Z_t | Z_{t+1}) \rightarrow P(x_{t+1:L} | Z_{t+1} = l)$ . This is then

$$P(x_{t+1} | Z_{t+1} = l) P(x_{t+2:L} | Z_{t+1} = l) = e_l(x_{t+1}) \cdot b_l(t+1).$$

$$\rightarrow b_k(t) = \sum_{l \in S} e_l(x_{t+1}) b_l(t+1) a_{k,l}, // \frac{f_k(t) b_k(t)}{P(x_{1:L})}$$

Posterior Decoding: for  $t=1\dots L$ , let  $z_t^* := \underset{k \in S}{\operatorname{argmax}} P(Z_t = k | X_{1:L})$

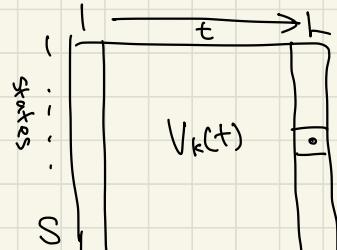
↳ This has a pathological issue that  $z_1^* \dots z_L^*$  may be an invalid sequence!

Viterbi Decoding: goal is to find  $z_{1:L}^{**} = \underset{z_1 \dots z_L}{\operatorname{argmax}} P(z_{1:L} | X_{1:L})$ .

Observe that  $P(z_{1:L} | X_{1:L}) = P(z_{1:L}, X_{1:L}) / P(X_{1:L})$ , but the denom.

never changes!  $\rightarrow z_{1:L}^{**} = \underset{z_1 \dots z_L}{\operatorname{argmax}} P(z_{1:L}, X_{1:L})$  is equivalent.

$\rightarrow$  Let  $v_k(t) := \max_{z_1 \dots z_{t-1}} P(z_{1:t-1}, Z_t = k, X_{1:t})$ . We fill in the DP table:



Base Case:  $v_k(1) = P(Z_1 = k, X_1) = P(X_1 | Z_1 = k) \pi_k$

$$= e_k(x_1) \cdot \pi_k.$$

General Case:  $v_k(t) = \max_{z_{1:t-1}} \max_{k \in S} P(z_{1:t-1}, Z_t = k, X_{1:t})$

We can rewrite the inner prob. as:  $\overbrace{P(Z_{1:t-2}, Z_{t-1}=k, X_{1:t-1})}^{\text{"almost" } V_k(t-1)} P(Z_t=l, X_t | Z_{1:t-2}, Z_{t-1}=k, X_{1:t-1})$ . By cond. ind., this is  $P(Z_{1:t-2}, Z_{t-1}=k, X_{1:t-1}) \underbrace{P(Z_t=l, X_t | Z_{t-1}=k)}_{P(X_t | Z_t=l, Z_{t-1}=k) P(Z_t=l | Z_{t-1}=k)} = e_k(x_t) \cdot a_{kl}$ .

$\rightarrow V_k(t) = e_k(x_t) \max_{k \in S} a_{kl} \cdot V_k(t-1)$ . Keep pointers of argmax to back trace!

At column  $L$ ,  $Z^{**}_L = \arg \max_{k \in S} V_k(L)$ , then back trace to find optimal path.

## Parameter Estimation

Want to find  $\Theta^* = \arg \max_{\Theta} P_{\Theta}(X_{1:L})$ , but this is not necessarily convex!

Case 1: Latent Variables are known,  $(Z_1, \dots, Z_L) \in S^L$ .

Let  $A_{kl} := \#$  of transitions from state  $k$  to  $l$  in the sequence  $Z_{1:L}$ .

Let  $E_k(\sigma) := \#$  of emissions of  $\sigma \in \mathcal{A}$  from state  $k \in S$ .

$$P_{\Theta}(X_{1:L}, Z_{1:L}) = \underbrace{P_{\Theta}(X_{1:L} | Z_{1:L})}_{\prod_{k \in S} \prod_{\sigma \in A} [e_k(\sigma)]^{E_k(\sigma)}} \cdot \underbrace{P_{\Theta}(Z_{1:L})}_{\prod_{k \in S} [\prod_{i:j} A_{ij}]^{A_{ij}}}$$

$\left[ \prod_{i:j} A_{ij}^{A_{ij}} \right]$ . Since this is monotone, we can take a log in the arg max.

$$(\text{In HW}) \text{MLE given } Z_{1:L} \text{ are: } \hat{a}_{ij} = \frac{A_{ij}}{\sum_k A_{ik}}, \hat{e}_k(\sigma) = \frac{E_k(\sigma)}{\sum_{\sigma'} E_k(\sigma')}$$

Case 2: Expectation Maximization. Observed RV  $X$ , Latent RV  $Z$ .

$$\log p(x) = \text{ELBO} + KL(q(z|x) \| p(z|x))$$

not symmetric on  $q$  &  $p$ !!  
arbitrary

For discrete RV  $X$ ,  $\mu, \nu :=$  prob. measures on the same prob. space.

Def) KL Divergence:  $\text{KL}(\mu \mid \nu) := \sum_x \mu(x) \log\left(\frac{\mu(x)}{\nu(x)}\right)$ . (assume  $\nu(x)=0 \Rightarrow \mu(x)=0$ )

$\hookrightarrow \text{KL}(\mu \mid \nu) \geq 0$ , and is  $= 0$  iff  $\mu(x) = \nu(x) \forall x$  (same dist.)

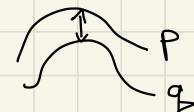
$\rightarrow$  Since  $\text{KL} \geq 0$ ,  $\log p(x) \geq \text{ELBO}$ , and we can optimize this value.

Def) Evidence Lower Bound:  $\text{ELBO} = \mathbb{E}_{z \sim q_{\theta}( \cdot | x)} [\log(p(x|z))] - \text{KL}(q_{\theta}(z|x) \parallel p(x))$ .

$$= \mathbb{E}_{z \sim q_{\theta}( \cdot | x)} [\log(p(x,z)) - \log(q_{\theta}(x,z))] + \log q_{\theta}(x).$$

$$\rightarrow \log p(x) \geq \mathbb{E}_{z \sim q_{\theta}( \cdot | x)} [\log p(x,z) - \log(q_{\theta}(x,z))] + \log q_{\theta}(x)$$

$$\rightarrow \underbrace{\log p(x) - \log q_{\theta}(x)}_{\text{ELBO}} \geq \mathbb{E}_{z \sim q_{\theta}( \cdot | x)} [\log p(x,z) - \log(q_{\theta}(x,z))].$$



Back to HMM, defined as:  $Q(\Theta' \mid \Theta_m) - Q(\Theta_m \mid \Theta_m) \geq 0$ . We maximize  $\Theta'$ , given by  $\Theta' = \underset{\Theta}{\operatorname{argmax}} Q(\Theta \mid \Theta_m)$ .

Idea: Devise an iterative algorithm that finds  $\Theta_0, \Theta_1, \dots, \Theta_m$  where

$P_{\Theta_i}(X_{1:L}) \leq P_{\Theta_{i+1}}(X_{1:L}) \quad \forall i$ , i.e. monotonicity on likelihood.  $\rightarrow$  EM!

Analysis of ELBO bound: Let  $q := P_{\Theta_m}$ , the current estimate, and  $p$

$\models P_{\Theta}$ , some arbitrary measure. The equation becomes:  $(x := (x_1, \dots, x_L))$

$$\log P_{\Theta}(x) - \log P_{\Theta_m}(x) \geq \underbrace{Q(\Theta \mid \Theta_m)}_{\text{constant!}} - \underbrace{Q(\Theta \mid \Theta_m)}_{\text{constant!}}$$

$$\mathbb{E}_{z \sim P_{\Theta}( \cdot | x)} [\log(P_{\Theta}(x,z))] = \mathbb{E} [\log(P_{\Theta}(x,z)) \mid x, \Theta_m] =: Q(\Theta \mid \Theta_m).$$

Now, if the RHS  $\geq 0$ , then  $\log P_{\Theta}(x) \geq \log P_{\Theta_m}(x)$ , so  $P_{\Theta}(x) \geq P_{\Theta_m}(x)$ !

Then, our goal is to find  $\underline{\Theta_{m+1}} = \arg\max_{\Theta} Q(\Theta | \Theta_m)$ , which satisfies the condition  $P_{\Theta_{m+1}}(x) \geq P_{\Theta_m}(x)$ .

$$\underline{\Theta_{m+1}} = \arg\max_{\Theta} E[\log P_{\Theta}(x, z) | x, \Theta_m]. \quad (\text{Baum-Welch})$$

↳ Upshot: In HMM, this has a closed form, so we can do argmax analytically!

Recall that when  $Z := (z_1, \dots, z_L)$  is known,  $\hat{a}_{ij} = \frac{A_{ij}}{\sum_k A_{ik}}$ ,  $\hat{e}_k(\sigma) = \frac{E_k(\sigma)}{\sum_{\sigma'} E_k(\sigma')}$ , from maximizing  $\log(P_{\Theta}(x, z)) = \prod_l [\pi_l]^{1\{z_l=l\}} \cdot \prod_{i,j} (a_{ij})^{A_{ij}} \cdot \prod_{\sigma, k} [e_k(\sigma)]^{E_k(\sigma)}$ .

If  $Z$  is unknown, the change is that now  $\underline{z_i, A_{ij}, E_k(\sigma)}$  depends on  $Z$ !!

↳ Taking expectation of log, we have a nice disjoint expression by variables:

$$E\left[\sum_l 1\{z_l=l\} \log(\pi_l) + \sum_{i,j} A_{ij}(Z) \log(a_{ij}) + \sum_{\sigma} \sum_k E_k(\sigma, Z) \log(e_k(\sigma))\right].$$

$$E\left[\sum_{i,j} A_{ij}(Z) \log(a_{ij})\right] = \sum_{i,j} \log(a_{ij}) E[A_{ij}(Z) | x, \Theta_m]. \text{ Solving this:}$$

$$\hat{a}_{ij}^{(cm)} = \frac{E[A_{ij}(Z) | x, \Theta_m]}{\sum_k E[A_{ik}(Z) | x, \Theta_m]}. A_{ij}(Z) \text{ is represented as } \sum_{t=1}^{L-1} 1\{z_t=i \wedge z_{t+1}=j\}.$$

$$\text{Then, } E[A_{ij}(Z) | x, \Theta_m] = \sum_{t=1}^{L-1} E[1\{z_t=i \wedge z_{t+1}=j\} | x, \Theta_m] = \sum_{t=1}^{L-1} P(z_t=i, z_{t+1}=j | x, \Theta_m)$$

$$= \sum_{t=1}^{L-1} f_i^{(cm)}(t) \cdot a_{ij}^{(cm)} \cdot e_j^{(cm)}(x_{t+1}) \cdot b_j^{(cm)}(t+1) / P_{\Theta_m}(x). \text{ (derivation in discussion)}$$

## Pseudoalignment [Guest Lecture by Sina Booshagi]

### De Bruijn Graphs: Construction & Properties

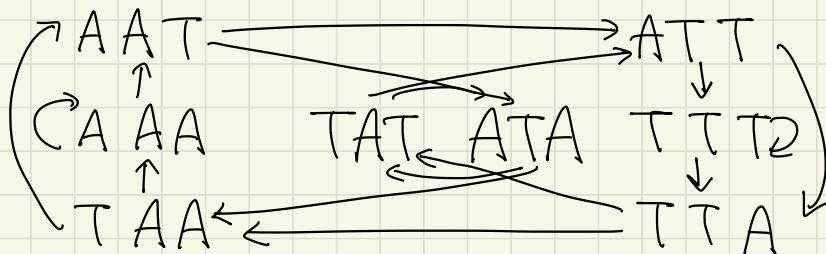
↳ Directed Graph, Alphabet  $\Sigma = \{A, C, G, T\}$ , integer  $K > 0$ .

Each node represents a char sequence of length K.

An edge is drawn between nodes sharing  $(K-1)$ -length overlapping.

i.e.,  $a \rightarrow b$  implies  $a = X \underline{S} \quad \& \quad b = \underline{S} Y$

ex)  $\Sigma = \{A, T\}$ ,  $K=2$ , then AAA  $\rightarrow$  AAT. full graph:



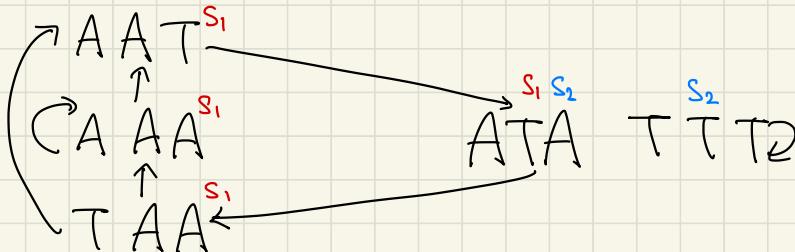
Some Properties of  $U = \text{DBG}(\Sigma, k)$ :

$(|\Sigma| \cdot n)$

1) # of nodes =  $|\Sigma|^k$ . 2) Indegree = Outdegree =  $|\Sigma|$ , # of edges =  $|\Sigma|^{k+1}$ .

Coloring DBG: Let  $S_1 := \{AAT, AAA, TAA, ATA\}$ ,  $S_2 := \{TTT, ATA\}$ .

$S_1, S_2 \subset U$ ,  $S_1 \cap S_2 = \{\text{ATA}\}$ . Let  $W := S_1 \cup S_2$ . Coloring is membership:



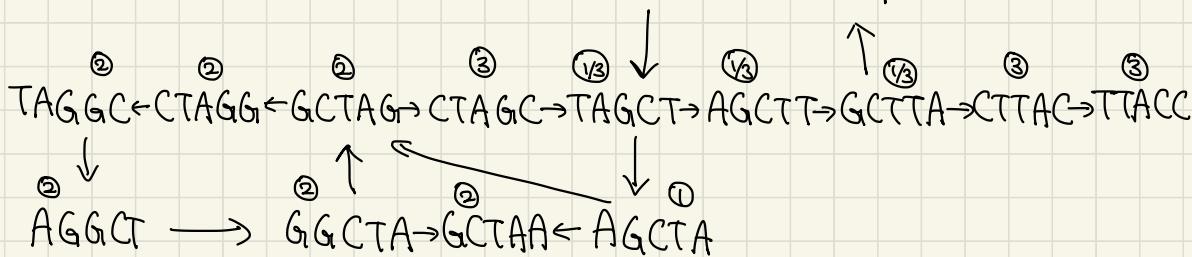
Construction of DBG: Given a set of sequences, extract the overlapping k-mers. e.g. AAGAT,  $K=4 \rightarrow \{AAGA, AGAT\}$ . Then, identify unique

k-mers & annotate them with its sequence origin. Create directed edges that link ( $K-1$ ) overlaps.

Ex) 3 length  $L=10$  sequences,  $K=5$ .  $\rightarrow \overbrace{(L-K+1)}^6$  k-mers per sequence

$\rightarrow \underbrace{AGCTT}_{①} \underbrace{AGCTA}_{②}, \underbrace{GCTAG}_{③} \underbrace{GCTAA}_{④}, \underbrace{CTAGC}_{⑤} \underbrace{TTAAC}_{⑥}$ .

$\overset{①}{TTAGC} \leftarrow \overset{①}{CTTAG}$



New sequences: CTA GG CT. Could this have originated from one of ①, ②, or ③?

↪ Break it up into k-mers: CTA GG, TA GG C, AG GCT.  $\rightarrow$  ②! unique mapping

AGCTT. ① & ③ are both feasible matches  $\rightarrow$  ambiguous mapping

TAGCTT. ① & ③ are both feasible matches  $\rightarrow$  ambiguous mapping

AGCTTAC. ④, ⑤, ①, intersection is ①  $\rightarrow$  resolved to unique mapping

Memory:  $O(|\Sigma|^K)$ , Time:  $O(L)$  construction,  $O(m)$  query

RNA-Sequencing (Pseudo) Alignment: Experiment quantifying RNAs in cell pop.

We want genes that the RNAs transcribed from. From references of the list of gene sequences, we can construct a colored DBGI. (Kallisto!)

↳ By aligning with DBG<sub>i</sub>, we get a vector of counts of each gene.

To resolve ambiguities, use EM to estimate  $\alpha_T$ , abundance of gene T.

Define a binary matrix  $(y)_{r,t} := \mathbb{1}\{\text{read } r \text{ could be gene } t\}$ .

---

## RNA Seq

Simplified Model: Let  $T := \{t_1, \dots, t_k\}$  be the set of transcripts.

Let  $R := \{r_1, \dots, r_n\}$  be the set of reads from a given sample/cell.

Assume: single-ended reads, error free reads, all transcripts are known, and reads originated from the reference transcriptome.

Goal: Find the transcript from which read originated.

↳ not always unambiguously determined!

Rephrased Goal: Find the feasible set of transcripts

Solution 1: align reads against transcriptomes using algorithms & heuristics.

↳ ex) STAR alignment using suffix trees & arrays

Solution 2: Pseudoalignment; only care about whether it aligns, not where it aligns. → Compatibility Matrix  $Y = (y_{r,k})$  with DBG<sub>i</sub>!

↳ One issue: some ambiguous sequences are actually not ambiguous ...

DBG does not require consecutive k-mers to appear consecutive in order!

→ however, for sufficiently large k, this is unlikely.

For matrix Y, define an equivalence class for each  $\epsilon_i = \{t \in T \mid y_{fit} = 1\}$ :

basically, if the rows look the same, they are functionally the same.

Assuming that a read is generated from a transcript  $t \in T$  w.p.  $\frac{\text{abundance of } t}{\text{length of } t}$ ,

and is a small uniformly sampled segment from t.

Unambiguous Model  $X_t : \hat{x}_t = \frac{x_t}{N}$  where  $X_t := \# \text{ of times } t \text{ is observed}$ .

↪ When reads are ambiguous, split it into each  $X_t$  relevant. But this is not optimal. → Use EM! Observed variable is matrix Y that does not commit to a single t, and latent is matrix Z that has a "committed mapping" s.t. each row has only one possible t.

↪ Intuitively, if we have a bunch of one t in multiple reads, we should weight it more when splitting into  $X_m$ .

$\{\alpha_k^{(t+1)}\}_{k=1}^K$  ↗ E step (revising rows)

Ambiguous Model  $X_t$ : Let  $Z_{ik}^{(t+1)} := E[Z_{ik}^{(t)} \mid Y_i, \Theta^{(t)}] = \Pr[Z_{ik} = 1 \mid Y_i, \Theta^{(t)}] = \frac{x_k^{(t)} y_{ik}}{\sum_k x_k^{(t)} y_{ik}}$

Then let  $\eta_k^{(t+1)} := \sum_{i=1}^N Z_{ik}^{(t+1)}$ .  $\hat{x}_k^{(t+1)} := \frac{\eta_k^{(t+1)}}{N}$ . ← M step (aggregating columns)

Hypothesis Testing

Parameterized HT: data  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_\theta$ ,  $\theta \in \Theta$ . Two candidates  $\Theta_0, \Theta_1$ .

$H_0$  is when  $\theta \in \Theta_0$  (null),  $H_1$  is when  $\theta \in \Theta_1$  (alternative).

ex) coin flips  $X_1, \dots, X_n \sim \text{Ber}(p)$ .  $H_0: p=0.5$ .  $H_1: p>0.5$ . (head prob.)

a natural statistic is  $T(X) = \sum_{i=1}^n X_i$ . → how do we interpret this value?

Def) Power Function:  $\Pi(\theta|\delta) := \Pr[T(X) \in R | \delta]$

Neyman-Pearson Framework:  $(\Pi(\delta) = 1 - \beta(\delta))$

reality \ decision	$H_0$ is true	$H_1$ is true
decide $H_0$	good	Type II ( $\beta$ ) Error
decide $H_1$	Type I ( $\alpha$ ) Error	good

We will consider  $H_0$  as status quo and only reject when the evidence

is strong enough to change my world view!

Define  $\alpha = \Pr[T(X) \in R | \theta \sim \Theta_0]$ ,  $\beta = \Pr[T(X) \notin R | \theta \sim \Theta_1]$ .

We want to keep  $\alpha < c$ , say  $c = 0.05$ , while minimizing  $\beta$ .

Lemma) Let  $\delta^*$  be a test that  $H_0$  is not rejected if:

$a f_0(x) > b f_1(x)$  & rejects if  $<$ . Then for every other test  $\delta$ ,

$a\alpha(\delta^*) + b\beta(\delta^*) \leq a\alpha(\delta) + b\beta(\delta)$ . (for simple tests)

Proof: Let  $S_1$  be the rejection region of some  $\delta$ .  $S_0 := \bar{S}_1$ .

$$a\alpha(\delta) + b\beta(\delta) = a \sum_{x \in S_1} f_0(x) + b \sum_{x \in S_0} f_1(x) = a \sum_{x \in S_1} f_0(x) + b \left[ 1 - \sum_{x \in S_1} f_1(x) \right] \\ = b + \sum_{x \in S_1} [af_0(x) - bf_1(x)]. \text{ Then } \delta^* \text{ minimizes such value.}$$

$\rightarrow$  the ratio  $\frac{f_1(x)}{f_0(x)} > C$  determines rejection (likelihood HT).

ex)  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ .  $H_0: \mu = \mu_0$ ,  $H_1: \mu = \mu_1 > \mu_0$ .

$$\frac{f_1(x)}{f_0(x)} = \exp\left(\bar{\sigma}^2(\mu_0 - \mu_1) \sum [X_i - \frac{n}{2\sigma^2}(\mu_1^2 - \mu_0^2)]\right). \text{ We want}$$

$$P_{H_0}(T(X) > C) = \alpha \rightarrow P_{H_0}(\bar{X}_n > C) = \alpha. \rightarrow C = \mu_0 + Z_{1-\alpha} \frac{\sigma}{\sqrt{n}}.$$

$\hookrightarrow$  But what if we want to test  $H_0: \mu \leq \mu_0$ ,  $H_1: \mu \geq \mu_0$ ? (composite)

Def) Uniformly Most Powerful (UMP) Test:  $\delta^*$  at level  $\alpha_0$  where  $\forall \delta$  s.t.

$$\alpha(\delta) \leq \alpha_0, \forall \theta \in \Theta_1, \pi(\theta | \delta) \leq \pi(\theta | \delta^*).$$

ex)  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(p)$ . Let  $Y = \sum_{i=1}^n X_i$ .  $f_p(x) = p^y (1-p)^{n-y}$ .  $H_0: p = p_0$ ,

$H_1: p = p_1$ .  $\frac{f_1(x)}{f_0(x)} = \left[ \frac{p_1(1-p_0)}{p_0(1-p_1)} \right]^Y \left[ \frac{(1-p_1)}{(1-p_0)} \right]^{n-y}$ . Observe that  $\frac{f_1(x)}{f_0(x)}$  1) only depends on  $X$  through  $Y = T(X)$ , 2) is monotone on  $Y \rightarrow$  These

conditions ensure an UMP for  $\begin{cases} H_0: p \leq p_0 \\ H_1: p > p_0 \end{cases}$  and  $\begin{cases} H_0: p \geq p_0 \\ H_1: p < p_0 \end{cases}$ . For example,

if  $p_0 = 0.5$ ,  $n = 20$ , we can set  $Y \geq 18$  for rejection to get some  $\alpha_0$ .

However, we do not have UMP for  $\begin{cases} H_0: p = p_0 \\ H_1: p \neq p_0 \end{cases}$ .

$\Rightarrow$  How do we explicitly calculate the  $\alpha_0$  for each  $Y$ ?

ex)  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$  where  $\sigma^2$  is known,  $\mu$  is unknown. Recall that  $f(X) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X-\mu)^2}{2\sigma^2}\right) \rightarrow f_p(\bar{X}) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2\right)$   
 $\rightarrow \frac{f_p(x)}{f_p(x_0)} = \exp\left(\frac{n(\mu_0 - \mu_0)}{\sigma^2} (\bar{X}_n - \frac{1}{n}(\mu_0 + \mu_0))\right).$  reject when  $\bar{X}_n > C$  is a UMP for  $H_0: \mu \leq \mu_0, H_1: \mu > \mu_0.$  But what is  $C?$  It depends on  $\alpha_0.$   
 $\Rightarrow$  Fix  $\alpha_0.$  We want  $\Pr[\bar{X}_n > C | H_0] = \alpha_0,$  and equivalently,  
 $\Pr[\bar{X}_n \leq C | H_0] = 1 - \alpha_0 \rightarrow \Pr\left[\frac{\bar{X}_n - \mu_0}{\sigma/\sqrt{n}} \leq \frac{C - \mu_0}{\sigma/\sqrt{n}}\right] = 1 - \alpha_0.$  Then,  $\frac{C - \mu_0}{\sigma/\sqrt{n}}$  is a RV  $\sim N(0, 1).$  Thus  $\frac{C - \mu_0}{\sigma/\sqrt{n}} = \phi^{-1}(1 - \alpha_0) \rightarrow C = \mu_0 + \frac{\sigma}{\sqrt{n}} \phi^{-1}(1 - \alpha_0),$   
 $(\ast \phi^{-1}(k)$  is the inverse Gaussian Cdf.)

Def) P-Value: For a sample  $X$ , the lowest  $\alpha$  for which we will reject the null for test  $S(X).$  Equivalently, it is also the probability of receiving a statistic "more extreme" than  $T(X), \Pr[T \geq T(X) | H_0].$

ex) For coin flipping, say  $H_0: p = \frac{1}{2}, H_1: p > \frac{1}{2}.$  We get  $\frac{18}{20}$  heads.

Then p-value is  $\Pr[T = 18] + \Pr[T = 19] + \Pr[T = 20]$  for  $H_0.$

ex) For mean-finding, if  $\bar{X}_n = M,$  p-value =  $\Pr[\bar{X}_n \geq M | \mu = \mu_0]$   
 $= 1 - \Pr\left[\frac{\bar{X}_n - \mu_0}{\sigma/\sqrt{n}} \leq \frac{M - \mu_0}{\sigma/\sqrt{n}}\right] = 1 - \phi^{-1}\left(\frac{M - \mu_0}{\sigma/\sqrt{n}}\right).$

( $\ast$  In real experiments, we only report p-values instead of  $\alpha$ s.)

$\rightarrow$  Can we find a way to test  $H_0: \mu = \mu_0, H_1: \mu \neq \mu_0?$  (two-sided)

We want to find two thresholds  $c_1, c_2$  s.t.  $\Pr[X_n \leq c_1 | H_0] + \Pr[X_n \geq c_2 | H_0] = \alpha$ . We could vary  $c_1, c_2$ , but an intuitive choice would be to set each of the probabilities to  $\frac{1}{2}\alpha$ . Also,  $\sigma$  is unknown now.

→ We can estimate  $S^2 \approx \sigma^2$  via  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ . Then, we have

$T = \frac{\bar{X}_n - \mu_0}{S/\sqrt{n}} \stackrel{H_0}{\sim} t_{n-1}$ , the "Student t" distribution.

(\*  $S^2$  is defined with  $\frac{1}{n-1}$  for unbiased estimating!)

Inferences on samples: rows are genes, columns are samples.

Samples are split into control / treatment groups. Consider gene 1.

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu_1, \sigma^2)$  be cells before treatment, and  $Y_1, \dots,$

$Y_n \stackrel{iid}{\sim} N(\mu_2, \sigma^2)$  be cells after treatment. This is a paired sample.

Then, we can take  $D_i := X_i - Y_i$ , and under  $H_0: \mu_1 = \mu_2$ ,  $D_1, \dots, D_n \stackrel{iid}{\sim}$

$N(0, \sigma_D^2)$ . →  $\bar{D}_n, S_D^2$  gives  $T = \frac{\bar{D}_n}{S_D/\sqrt{n}} \stackrel{H_0}{\sim} t_{n-1}$ .

2 sample case: In this case, no samples are paired.  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu_1, \sigma_1^2)$

and  $Y_1, \dots, Y_m \stackrel{iid}{\sim} N(\mu_2, \sigma_2^2)$ . We define a new statistic  $\bar{S}_p$ , the pooled

variance  $S_p^2 := \frac{(n-1)S_x^2 + (m-1)S_y^2}{n+m-2}$ . Then,  $T = \frac{\bar{X}_n - \bar{Y}_m}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \stackrel{H_0}{\sim} t_{n+m-2}$ .

→ Can we test multiple rows (genes)? There is a glaring problem.

Type 1 error is  $\Pr[\text{reject } H_0 | H_0 \text{ is true}] \leq \alpha$ , say 5%. If we have  $n=10$ ,  $\Pr[\text{never makes T1 error}] \leq 1 - (1 - 0.05)^{10}$ , which is bad.

Def) Family-Wide Error Rate (FWER):  $\Pr[\text{at least one T1 error}]$ .

→ Given multiple null hypotheses  $H_1, \dots, H_m$ , we compute p-values  $p_1, \dots, p_m$ .

We want the FWER over  $p_1, \dots, p_m \leq \alpha$ .

Bonferroni Method: test  $p_i \leq \frac{\alpha}{m}$ .

Proof:  $\text{FWER} = \Pr\left[\bigcup_{i \in I_0} p_i \leq \frac{\alpha}{m}\right]$  where  $I_0 := \text{set of true nulls}$ . By union

bound,  $\leq \sum_{i \in I_0} \Pr[p_i \leq \frac{\alpha}{m}] = M_0 \cdot \underbrace{\frac{\alpha}{m}}$  where  $M_0 := \# \text{ of true nulls}, |I_0|$ .

$M_0 \leq m$ , so  $\text{FWER} \leq \alpha$ .  
( $p_i \stackrel{H_0}{\sim} \text{Uni}(0,1)$ )

↪ This is nice, but we lose in statistical power (T2 error blows up).

Bonferroni-Holm Method:  $H_1, \dots, H_m \rightarrow p_1, \dots, p_m$ . Sort in ascending p-values,  $H_{(1)}, \dots, H_{(m)} \rightarrow p_{(1)} \leq \dots \leq p_{(m)}$ . If  $p_{(j)} \leq \frac{\alpha}{m-j+1}$  for  $j = 1, \dots, i$ , reject  $p_{(i)}$ . Then,  $\text{FWER} \leq \alpha$  while having more statistical power!

↪ However, even this does not give enough power.

decision \ reality	accept $H_0$	reject $H_0$	total
$H_0$ true	U	V	$M_0$
$H_0$ false	T	S	$M - M_0$
total	$M - R$	R	M

Def) False Discovery Rate:

Observe that

$FWER = \Pr[V \geq 1]$ . Let the False Discovery Proportion  $FDP := \frac{V}{\max(R, 1)}$ .  
 $FDR := E[FDP]$ . We want to control this value.

Benjamini-Hochberg Procedure: Sort  $p_{(1)} \leq \dots \leq p_{(m)}$ . Let  $i_0$  be the largest  $i$  s.t.  $p_{(i)} \leq \frac{i}{m} \alpha$ . Reject  $H_{(1)}, \dots, H_{(i_0)}$ .

Theorem) If  $p_1, \dots, p_m$  are independent (can be relaxed a bit), then the B-H procedure guarantees  $FDR = \frac{m_0}{m} \alpha \leq \alpha$ .