

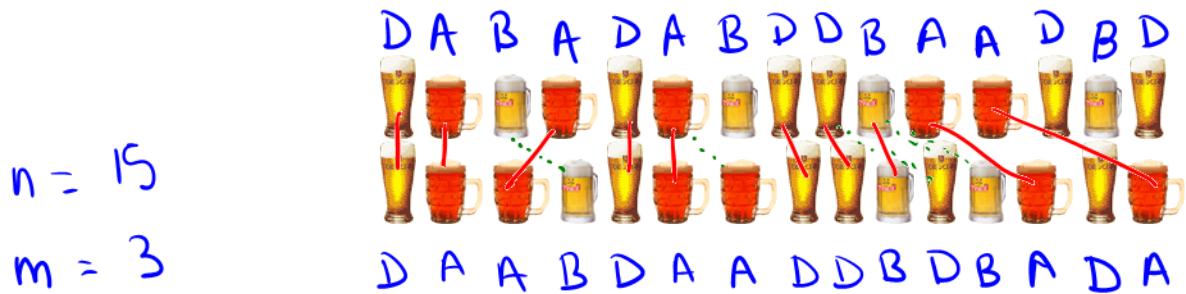
A symptoms in
Sequences Comparisons

J.C. Breton

Ü Isländer

IESC, 15/05/2018.

Longest Common Subsequence



Longest Common Subsequence

D	A	B	A	D	A	B	D	D	B	-	A	A	D	B	D
D	A		A	B		D		A		D	D		B		D

D	A		B	D	-	-	-
D	A		A	B	P

Length of the Longest Common Subsequence

$$L C_n = 10$$

Longest common subsequence

= alignment (optimal) with holes

Algorithm (Dynamic Programming)

$$|\text{LCS}(x_1 \dots x_{e_2}; y_1 \dots y_e)|$$

$$= \begin{cases} |\text{LCS}(x_1 \dots x_{e_2-1}; y_1 \dots y_{e-1})| + 1 & \text{if } x_{e_2} = y_e, \\ \max(|\text{LCS}(x_1 \dots x_{e_2}; y_1 \dots y_{e-1})|, |\text{LCS}(x_1 \dots x_{e_2-1}; y_1 \dots y_e)|) & \text{if } x_{e_2} \neq y_e. \end{cases}$$

$$|\text{LCS}(x_1 \dots x_{e_2}; 0)| = |\text{LCS}(0; y_1 \dots y_e)| = 0$$

$$\forall l_2 = 1, 2, \dots, n.$$

$$\forall l = 1, 2, \dots, n$$

Edit/Levenshtein distance

$$d(x_1 \dots x_n; y_1 \dots y_n) = 2(n - \text{LC}_n)$$

Penalties

christian
krystyaan

c	h	n	i	y	s	t	i	a	a	n
h	r	r	y	s	t	y	g	a	l	n

$$X_1 \dots X_n = H T H H T$$

↓ ↓ || |

$$Y_1 \dots Y_n = H H T H T$$

6 Common Subsequences, $LC_n = 4$.

							(n, n)
T	O	I	O	O	I		
H	I	O				O	
T	O		O	O	I		
H	I	O		I	I	O	
H	I	O	I	I	I	O	
(0, 0)	H	T	H	H	H	T	

LC_n : length of the longest strictly increasing north-east path from $(0, 0)$ to (n, n) .

$X_1, X_2, \dots, X_n, \dots$ iid with values in

$$\mathcal{C}_m = \{\alpha_1, \alpha_2, \dots, \alpha_m\}$$

$Y_1, Y_2, \dots, Y_n, \dots$ iid with values in

$$\mathcal{D}_m = \{\beta_1, \beta_2, \dots, \beta_m\}$$

$$X = (X_i)_{i \geq 1} \quad \text{and} \quad Y = (Y_i)_{i \geq 1}$$

LC_n = Length of the longest common
subsequence of

X_1, \dots, X_n and Y_1, \dots, Y_n .

$$LC_n = \max_{k=1, \dots, n} b_k : \exists \begin{array}{l} 1 \leq i_1 < i_2 < \dots < i_{b_k} \leq n \\ 1 \leq j_1 < j_2 < \dots < j_{b_k} \leq n \end{array} \text{ s.t. } X_{i_s} = Y_{j_s}, s = 1, 2, \dots, b_k.$$

$$\mathbb{E} LC_n \sim ?, \quad \text{Var } LC_n \sim ?, \quad \frac{\mathbb{E} LC_n - \mathbb{E} LC_n}{(\text{Var } LC_n)^{\frac{1}{2}}} \Rightarrow ?$$

Asymptotic behavior in Mean,

Variance, Distribution of LC_n ?

Chvátal-Sankoff (1975)

$$LCS(X_1 \cdots X_{n_1+n_2}; Y_1 \cdots Y_{n_1+n_2}) \geq LCS(X_1 \cdots X_{n_1}; Y_1 \cdots Y_{n_1})$$

$$+ LCS(X_{n_1+1} \cdots X_{n_1+n_2}; Y_{n_1+1} \cdots Y_{n_1+n_2})$$

So taking expectation and iid (stationarity)

$$\mathbb{E}LC_{n_1+n_2} \geq \mathbb{E}LC_{n_1} + \mathbb{E}LC_{n_2} . \text{ So}$$

Fekete's Lemma, $\mathbb{E}\frac{LC_n}{n} \xrightarrow{n \rightarrow +\infty} \gamma^* = \sup_{k \geq 1} \frac{\mathbb{E}LC_k}{k}$

$$\gamma^* \approx 0.812, \gamma^* \approx 0.717, \gamma^* \approx 0.654, \dots$$

(iid uniform case).

Kiwi, Loebel, Matoušek (2005),

$$\lim_{m \rightarrow +\infty} \sqrt{m} \gamma_m^* = 2 .$$

(Conjecture of Sankoff and Mainville, 1983)
Link with Ulam's Problem.

Speed of convergence

$$n\tau_m^* - K_A \sqrt{n \log n} \leq \mathbb{E} LC_n \leq n\tau_m^* \quad (\text{Alexander 1994})$$

Variance?

Efron-Stein Inequality / Tensorisation

$$X_1, X_2, \dots, X_n \perp\!\!\!\perp, \quad X_i \sim \mu_i, \mu = \bigotimes_{i=1}^n \mu_i$$

$$f: \mathbb{R}^n \longrightarrow \mathbb{R},$$

$$\text{Var } f(X_1, \dots, X_n) \leq \mathbb{E} \sum_{i=1}^n \text{Var}_{\mu_i} f(X_1, \dots, X_n)$$

$$= \frac{1}{2} \mathbb{E} \sum_{i=1}^n \mathbb{E} \left(f(X_1, \dots, X_n) - f(X_1, \dots, \hat{X}_i, \dots, X_{i+1}, \dots, X_n) \right)^2$$

$$LC_n = f(X_1, \dots, X_n; Y_1, \dots, Y_n) \leq 1$$

Hence,

$$\text{Var } LC_n \leq \frac{1}{2} 2n = n \quad (\text{Steele 1986})$$

$\text{Var } LC_n \geq Kn$? Yes in some biased cases (Lemler, Matzinger, Amsalu, Gang, Ma, C.H., ...)
IID Uniform?

Largest Increasing Subsequences in Random Words

Order 1 = D, 2 = B, 3 = A

$n = 15$

$m = 3$



Instead of two sequences, just one:

Largest Increasing Subsequence

1 1 1 1 2 3 3

Alphabet $\mathcal{A}_5 = \{d_1 < d_2 < d_3 < d_4 < d_5\}$ $m = 5$

Mot: $d_1 \underline{d_3} \underline{d_2} \underline{d_4} \underline{d_2} \underline{d_3} \underline{d_5} \underline{d_4} \underline{d_2} \underline{d_4} \underline{d_3} \underline{d_5} \underline{d_4} \quad n = 12$

LI_n = length of the longest increasing subsequence

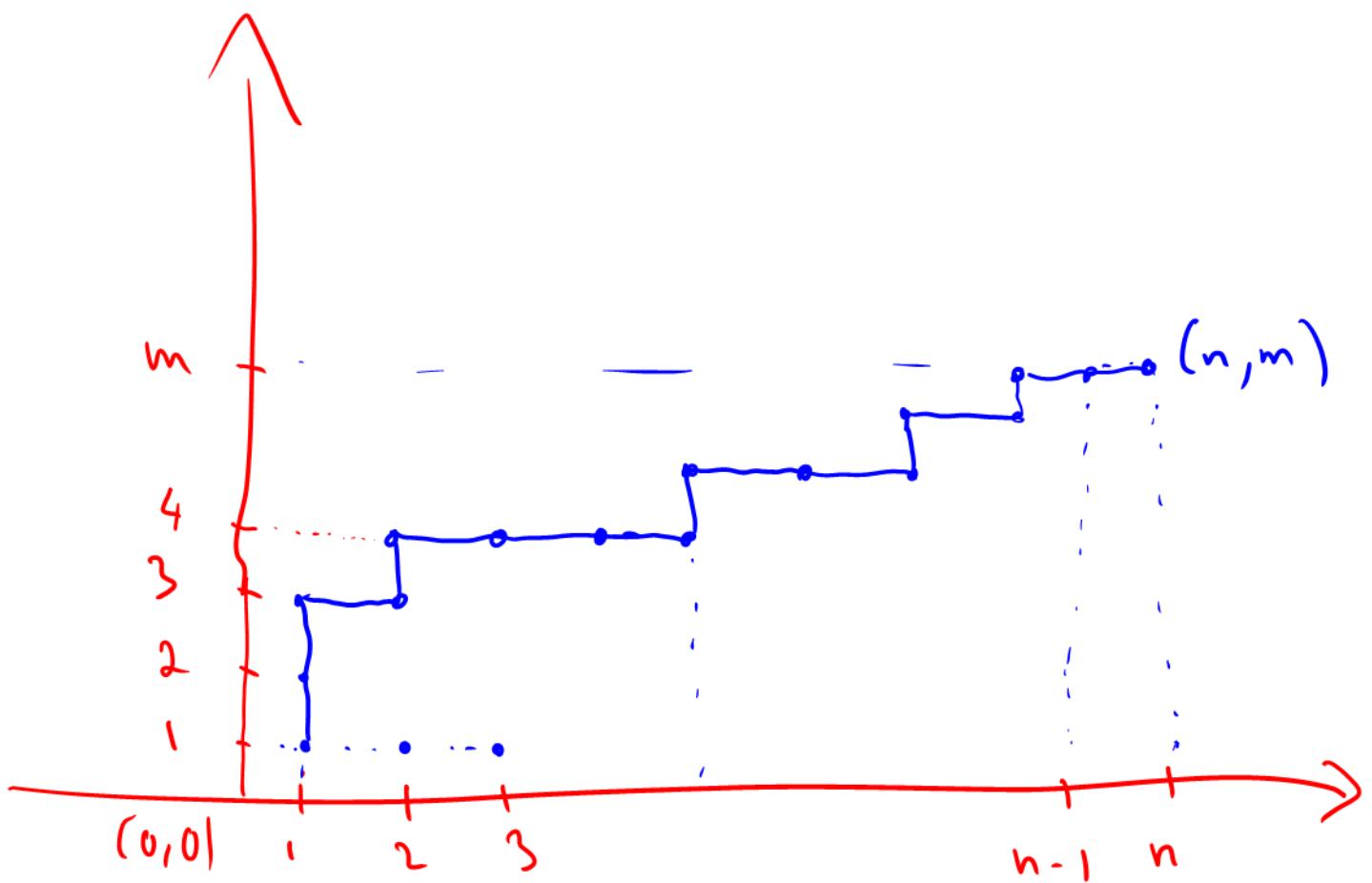
$d_1 \underline{d_2} \underline{d_2} \underline{d_3} \underline{d_4} \underline{d_4} \underline{d_5}$
 $d_1 \underline{d_2} \underline{d_2} \underline{d_3} \underline{d_4} \underline{d_4} \underline{d_4}$

$X_1, X_2, \dots, X_n, \dots$ iid $\mathcal{A}_m = \{d_1 < d_2 < \dots < d_m\}$

$LI_n = \max \{k : \exists 1 \leq i_1 < i_2 < \dots < i_k \leq n$

s.t. $X_{i_1} \leq X_{i_2} \leq \dots \leq X_{i_k}$

LI_n and Last Passage Percolation



$$LI_n = \max_{\pi \in NE} \sum_{(i,j) \in \pi} \mathbb{1}_{\{X_i = \alpha_j\}}$$

NE = North-East paths from $(1,1) \uparrow (m,n)$

The Bernoulli r.v. $\mathbb{1}_{\{X_i = \alpha_j\}}$ are dependent! (Exchangeable in the uniform case).

Kerov (1994) Let $P(X_1 = \lambda_2) = \frac{1}{m}$,

$$\lambda_2 = 1, -1, m.$$

$$\frac{\lambda_{\max} - n/m}{\sqrt{n/m}} \xrightarrow[n \rightarrow \infty]{} \lambda_{\max}^0, \text{ where } \lambda_{\max}^0$$

is the maximal eigenvalue of the
traceless $m \times m$ GUE, i.e.,

$$M - \frac{1}{m} \operatorname{Tr}(M) I, \text{ where } M \in \text{GUE}(m \times m).$$

* Tracy-Widom (2001), Johansson (2001)*

Baryshnikov (2001), Gravner, T-W (2001)

$$\lambda_{\max} \stackrel{L}{=} \max_{0=t_0 \leq t_1 \leq \dots \leq t_{m-1} \leq t_m=1} \sum_{i=1}^m (B^{(i)}(t_i) - B^{(i)}(t_{i-1})).$$

Doumerc, O'Connell, Yor, Bougerol, Jolin

Longest Common and Increasing Subsequences

Order 1 = D, 2 = B, 3 = A

$$n = 15$$

$$m = 3$$



Length of the longest common and increasing subsequence

$$LCIn = 7$$

$(X_i)_{i \geq 1}, (Y_j)_{j \geq 1},$

$$LCIn = \max_{k=1, \dots, n} l_k : \exists 1 \leq i_1 < i_2 < \dots < i_k \leq n \quad 1 \leq j_1 < j_2 < \dots < j_k \leq n$$

s.t. $X_{is} = Y_{js}, s = 1, 2, \dots, l_k$

$$X_{i_1} \leq X_{i_2} \leq \dots \leq X_{i_{l_k}}$$

$$Y_{j_1} \leq Y_{j_2} \leq \dots \leq Y_{j_{l_k}}$$

A symptotic Behavior of $LCIn$?

Théorème 1 (J.-C. Breton, C.H.). Let $(X_i)_{i \geq 1}$ and $(Y_i)_{i \geq 1}$ be two (independent) sequences of iid random variables uniformly distributed on $\{t_m = \{d_1 < d_2 < \dots < d_m\}\}$. Then,

$$\frac{\text{LCI}_n - n/m}{\sqrt{n/m}} \xrightarrow[n \rightarrow +\infty]{\text{max}} 0 = t_0 \leq t_1 \leq \dots \leq t_{m-1} \leq t_m = 1$$

$$\min \left(-1 \sum_{i=1}^m B_1^{(i)}(1) + \sum_{i=1}^m (B_1^{(i)}(t_i) - B_1^{(i)}(t_{i-1})), -1 \sum_{i=1}^m B_2^{(i)}(1) + \sum_{i=1}^m (B_2^{(i)}(t_i) - B_2^{(i)}(t_{i-1})) \right)$$

where

$$B_1 = (B_1^{(1)}, \dots, B_1^{(m)}) \quad \text{et} \quad B_2 = (B_2^{(1)}, \dots, B_2^{(m)})$$

are two m -dimensional standard BM (and II).

Ramdom Matrix Interpretation ??

Particular cases

• $X = Y, B_1 = B_2, \dots, B_m = B_n, LCI_n = LI_n$.

$$\frac{LI_n - \frac{n}{m}}{\sqrt{\frac{n}{m}}} \xrightarrow[n \rightarrow +\infty]{} \max_{\substack{t_0=0 \leq t_1 \leq \dots \leq t_{m-1} \leq t_m=1 \\ \text{Weyl Chamber}}} \sum_{i=1}^m (B^{(i)}(v_i) - B^{(i)}(v_{i-1})) - \frac{1}{m} \sum_{i=1}^m B^{(i)}(1).$$

Above result is very natural.

$$LI_n = \# \alpha_1 + \# \alpha_2 + \dots + \# \alpha_m.$$

$$= \max_{\substack{0 \leq k_0 \leq k_1 \leq k_2 \dots \leq k_{m-1} \leq k_m=n}} (\# \alpha_1 + \# \alpha_2 \text{ after } k_1 + \dots + \# \alpha_m \text{ after } k_{m-1})$$

$$\text{Moreover } \# \alpha_1 + \# \alpha_2 + \dots + \# \alpha_m = n!$$

$$\begin{aligned} \text{But } \# \alpha_i \text{ after } k_{i-1} &= \# \text{ of } \alpha_i \text{ before } k_i \\ \text{and before } k_i &\quad - \# \text{ of } \alpha_i \text{ before } k_{i-1}. \end{aligned}$$

$$\boxed{LI_n - \frac{n}{m} = -\frac{1}{m} \sum_{i=1}^m S_n^{(i)} + \max(S_{k_1}^{(1)} + S_{k_2}^{(2)} + \dots + S_{k_{m-1}}^{(m-1)}) + O(\sqrt{n})}$$

where $(S_{ek_i}^i)$ are random walks.

Hence Danzer's

$$\frac{I_n - \frac{n}{m}}{\sqrt{\frac{n}{m}}} \Rightarrow -\frac{1}{m} \sum_{i=1}^m B^{(i)}(1) + \max_{0 \leq t_0 \leq t_1 \leq \dots \leq t_m = 1} \sum_{i=1}^m (B^{(i)}(t_i) - B^{(i)}(t_{i-1}))$$

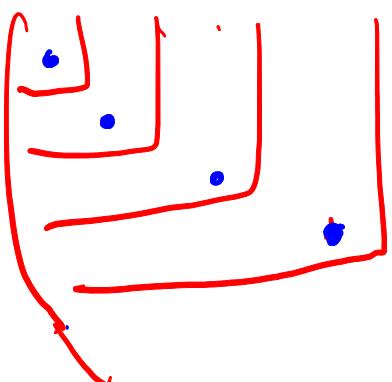
$\underbrace{\hspace{10em}}$
II d

$\lambda_{\max}(\text{GUE}_{m \times m})$

II d ?

$$-\frac{1}{m} \sum_{i=1}^m \lambda_i + \lambda_{\max} \quad \text{YES}$$

$M =$



Benaych-George
+
C.H.

- In general, the dependence structure of X with Y is reflected into a dependence structure between the two m -dimensional BM, B_1 and B_2 .

- $m=2$, X and $Y \perp\!\!\!\perp$, then

$$\frac{\text{LCIn} - n/2}{\sqrt{n}} \implies \max_{0 \leq t \leq 1} \min \left(B_1(t) - \frac{1}{2} B_1(1), B_2(t) - \frac{1}{2} B_2(1) \right),$$

where B_1 and B_2 are two $\perp\!\!\!\perp$ standard linear BM.

Lember, Matzinger, C.H. (2006).

Proof: "LCI_n"

(A) Combinatorial Representation of LCI_n.

$LCI_n = \max$ over random constraints of
min of random sums of randomly
stopped random variables.

(B) Derandomization. \mathbb{F} .

(C) Weak Invariance Principles
(Dvoretzky).

Back to LC_n (permutations/Words)

Ulam's Problem.

\mathfrak{A}_n : symmetric group, $\sigma \in \mathfrak{A}_n$, $LIn(\sigma)$ the length of the longest ↑ subsequence of σ . Study of the r.v. $LIn(\sigma)$, when \mathfrak{A}_n is endowed with the uniform measure.

Ulam: $E LIn(\sigma) \approx 1.7 \sqrt{n}$

$$\text{Var } LIn(\sigma) \approx$$

$$\frac{LIn(\sigma) - E LIn(\sigma)}{\sqrt{\text{Var } LIn(\sigma)}} \implies N(0, 1).$$

(Monte-Carlo Calculations in Problems of Mathematical Physics). Late 50's
- early 60's.

Vershik and Kerov (1977), Logan and Shepp (1977)

$$\frac{\mathbb{E} L I_n(\sigma)}{\sqrt{n}} \longrightarrow 2.$$

Conjecture (Simulations / Percolations) = 1980's

$$\text{Var } L I_n(\sigma) \approx n^{1/3}. \quad \left\{ \begin{array}{l} \text{Odlyzko} \\ + \text{Rains} \end{array} \right.$$

Baik - Deift - Johansson (1999) Kesten

$$\frac{L I_n(\sigma) - 2\sqrt{n}}{n^{1/6}} \Longrightarrow F_2.$$

F_2 : Tracy-Widom law, limiting law of max eigenvalue of $n \times n$ GUE.

F_2 also appears as limiting law of $L C_n$ of two random permutations?

(Aldous - Diaconis ~ 2000).

Proposition (Easy). Let σ and π be two uniform random permutations in S_n . Let $LC_n(\sigma, \pi)$ be the length of the longest common subsequence of $(\sigma_1, \dots, \sigma_n)$ and (π_1, \dots, π_n) . Then,

$$\frac{LC_n(\sigma, \pi) - 2\sqrt{n}}{n^{1/6}} \xrightarrow{\quad} F_2.$$

Theorem 2 (Hard). Let $(X_i)_{i \geq 1}$ and $(Y_i)_{i \geq 1}$ be two independent sequences of iid random variables with values in a finite alphabet \mathcal{G}_m . Let LC_n be the length of the longest common subsequence of X_1, \dots, X_n and Y_1, \dots, Y_n . If $\text{Var } LC_n \geq kn$, then

$$\frac{LC_n - \mathbb{E} LC_n}{\sqrt{\text{Var } LC_n}} \xrightarrow{\quad} N(0, 1).$$

Proof of Proposition

Fact 1. Note that $\text{LI}_n(\sigma)$, the length of the longest increasing subsequence of $\sigma \in A_n$ is equal to

$$\text{LC}_n((1, \dots, n), (\sigma_1, \dots, \sigma_n)) !$$

$$\begin{aligned} \text{id} : 1 & 2 \dots = \dots n \\ \sigma : \sigma_1, \sigma_2, \dots &= \sigma_n \end{aligned}$$

Fact 2. Let $a = (a_1, \dots, a_n)$ be a fixed permutation of $1, \dots, n$. Let σ be a uniform random permutation in A_n . Then,

$$\text{LI}_n(\sigma) \stackrel{d}{=} \text{LC}_n((a_1, \dots, a_n); (\sigma_1, \dots, \sigma_n))$$

Proof. Let $p \in A_n$ s.t $p(i) = a_i, \forall i = 1, \dots, n$.

Then clearly,

$\tilde{\sigma} = \sigma p$ is a uniform random permutation in A_n .

$$\begin{aligned}
& LC_n((a_1, \dots, a_n); (\sigma_1, \dots, \sigma_n)) \\
&= LC_n((\rho(1), \dots, \rho(n)), (\tilde{\sigma}_1, \dots, \tilde{\sigma}_n)) \\
&\stackrel{d}{=} LC_n((\rho(1), \dots, \rho(n)), (\tilde{\sigma}_{\rho(1)}, \dots, \tilde{\sigma}_{\rho(n)})) \\
&= LI_n(\sigma), \text{ by Fact 1.}
\end{aligned}$$

Fact 3. Let π and σ be two uniformly random permutations in S_n . Then,

$$P(LC_n(\pi, \sigma) \leq x) = P(LI_n(\sigma) \leq x)$$

Proof. Let $\mathcal{B} := \{a = (a_1, \dots, a_n) : a_i \in \{1, \dots, n\}$
and $a_i \neq a_j \text{ for } i \neq j\}$

$$P(LC_n(\pi, \sigma) \leq x)$$

$$= \sum_{a \in \Omega} P(LC_n(\pi, \sigma) \leq x \mid (\pi_1, \dots, \pi_n) = a) P((\pi_1, \dots, \pi_n) = a)$$

$$= \frac{1}{n!} \sum_{a \in \Omega} P(LC_n(a_1, \dots, a_n); (\sigma_1, \dots, \sigma_n) \leq x)$$

Fact 2

$$\downarrow = \frac{1}{n!} \sum_{a \in \Omega} P(LI_n(\sigma) \leq x)$$

$$= P(LI_n(\sigma) \leq x).$$

Hence, the proposition follows (is equivalent)
to the theorem of Bank-Derft and Johansson.



Proof of Theorem 2

Three step method

- (a) Stein's Method (Chatterjee, Lachièze-Rey and Peccati)
- (b) Short String Genericity principle
- (c) Variance Estimates

Distances

Kolmogorov distance

$$d_K(\mu_1, \mu_2) = \sup_{x \in \mathbb{R}} |\mu_1(-\infty, x] - \mu_2(-\infty, x]|$$

$$= \sup_{h \in \mathcal{H}_1} \left| \int h d\mu_1 - \int h d\mu_2 \right|$$

$$\mathcal{H}_1 = \left\{ \frac{1}{1+x^2} : x \in \mathbb{R} \right\}$$

Monge-Kantorovich-Wasserstein distance

$$d_1(\mu_1, \mu_2) = \sup_{h \in \mathcal{H}_2} \left| \int h d\mu_1 - \int h d\mu_2 \right|$$

$$\mathcal{H}_2 = \left\{ h: \mathbb{R} \rightarrow \mathbb{R} : |h(x) - h(y)| \leq |x - y| \right\} = \text{Lip}(\beta).$$

$$d_K(\mu_1, \mu_2) \leq \sqrt{2C d_1(\mu_1, \mu_2)}, \text{ if } \mu_2 \text{ is a.c. with density bounded by } C.$$

Chatterjee (2008) let $f: \mathbb{R}^n \rightarrow \mathbb{R}$,

let $W = (W_1, \dots, W_n)$ have $\mathbb{1}$ components

let $0 < \sigma^2 = \text{Var } f(W) < +\infty$. Then,

$$d_2\left(\frac{f(W) - \mathbb{E}f(W)}{\sqrt{\text{Var } f(W)}}, \mathcal{T}\right) \leq \sqrt{\frac{\text{Var } T}{\sigma^2}} + \frac{1}{2\sigma^3} \sum_{j=1}^n \mathbb{E}|\Delta_j f(W)|^3$$

Lachièze-Rey & Peccati (2016)

$$\begin{aligned} d_K\left(\frac{f(W) - \mathbb{E}f(W)}{\sqrt{\text{Var } f(W)}}, \mathcal{T}\right) &\leq \sqrt{\frac{\text{Var } T}{\sigma^2}} + \sqrt{\frac{\text{Var } \tilde{T}}{\sigma^2}} \\ &+ \frac{1}{4\sigma^3} \sum_{j=1}^n \sqrt{\mathbb{E}|\Delta_j f(W)|^6 + \frac{\sqrt{2\pi}}{16\sigma^3} \sum_{j=1}^n \mathbb{E}|\Delta_j f(W)|^3}, \end{aligned}$$

where $\Delta_j f = \dots, \mathcal{T} = \dots$

and $\tilde{\mathcal{T}} = \dots$

Let w' be an $\perp\!\!\!\perp$ copy of w , then

$$\Delta_j f(w) = f(w) - f(w'), \text{ where}$$

$$f(w') = f(w_1, w_2, \dots, w_{j-1}, w'_j, w_{j+1}, \dots, w_n).$$

For any $A \subset \{1, \dots, n\}$, let w^A be defined via $w_i^A = \begin{cases} w'_i & \text{if } i \in A \\ w_i & \text{if } i \notin A. \end{cases}$

$$\text{Let } T_A = \sum_{j \notin A} \Delta_j f(w) \Delta_j f(w^A) \text{ and}$$

$$T = \frac{1}{2} \sum_{A \subset \{1, \dots, n\}} \frac{T_A}{\binom{n}{|A|}(n-|A|)}.$$

$$\tilde{T}_A = \sum_{j \notin A} \Delta_j f(w) |\Delta_j f(w^A)|$$

$$\tilde{T} = \frac{1}{2} \sum_{A \subset \{1, \dots, n\}} \frac{\tilde{T}_A}{\binom{n}{|A|}(n-|A|)}$$

In our case, $W = (X_1, \dots, X_n; Y_1, \dots, Y_n)$
and

$$f(W) = LC_n(X_1 \dots X_n; Y_1 \dots Y_n).$$

But for $1 \leq j \leq n$,

$$|\Delta_j f(W)| = |LC_n(X_1 \dots X_n; Y_1 \dots Y_n) - LC_n(X_1 \dots X_{j-1}, X'_j X_{j+1} \dots X_n; Y_1 \dots Y_n)|$$

$$\leq 1,$$

and similarly for $n+1 \leq j \leq 2n$. So

$$\frac{1}{\sigma^3} \sum_{j=1}^n |\mathbb{E} |\Delta_j f(W)|^3 = \frac{1}{\sigma^3} \sum_{j=1}^{2n} |\mathbb{E} |\Delta_j LC_n|^3$$

$$\leq \frac{2n}{\sigma^3} = \frac{n}{\sigma^3} \leq \frac{1}{k^{3/2} \sqrt{n}}$$

If $\sigma^2 \geq k_n$. Similarly for $\mathbb{E} |\Delta_j LC_n|^6$

How to estimate

$$\sqrt{\frac{\text{Var } T}{\sigma^2}}, \sqrt{\frac{\text{Var } T}{\sigma^2}} ?$$

$$\text{Var } T = \frac{1}{4} \text{Var} \left(\sum_{A \notin \{1, \dots, 2n\}} \sum_{j \notin A} \frac{\Delta_j f(w) \Delta_j f(w^A)}{\binom{2n}{|A|} (2n - |A|)} \right)$$

$$= \frac{1}{4} \sum_{(A, B, j, l_1, l_2) \in S_1} \frac{\text{Cov}(\Delta_j f(w) \Delta_j f(w^A), \Delta_{l_1} f(w) \Delta_{l_2} f(w^B))}{\binom{2n}{|A|} (2n - |A|) \binom{2n}{|B|} (2n - |B|)}$$

$$S_1 = \left\{ (A, B, j, l_1, l_2) : A \notin \{1, \dots, 2n\}, B \notin \{1, \dots, 2n\} \right.$$

$$\left. j \notin A, l_2 \notin B \right\}.$$

General lemma, shows that

$$\text{Var } T \leq \frac{n^2}{2} \quad (\text{not good enough})$$

Since $\sigma^2 \geq Kn$ need $O(n^2)$. Hence $\frac{\sqrt{\text{Var } T}}{\sigma^2}$ needs to be estimated differently

Replace S_1 by several subsets, estimate each parts (≈ 20 pages) and use

Short Strings Genericity

(H. Matzinger and C.H.)

Assume $n = vd$, say v is fixed, for now,
(then small compared to n) and let the integers

$$n_0 = 0 \leq n_1 \leq n_2 \leq \dots \leq n_{d-1} \leq n_d = n$$

be such that

$$\text{LC}_n = \sum_{i=1}^d |\text{LCS}(x_{v(i-1)+1} \dots x_{vi}; y_{n_{i-1}+1} \dots y_{n_i})|$$

where $\text{LCS}(\dots; \dots)$ is the longest common subsequence
of the words $x_{v(i-1)+1} \dots x_{vi}$ and $y_{n_{i-1}+1} \dots y_{n_i}$.

i.e.; this is an optimal alignment aligning
 $[v(i-1)+1, vi]$ with $[n_{i-1}+1, n_i]$, for all
 $i = 1, 2, \dots, d$.

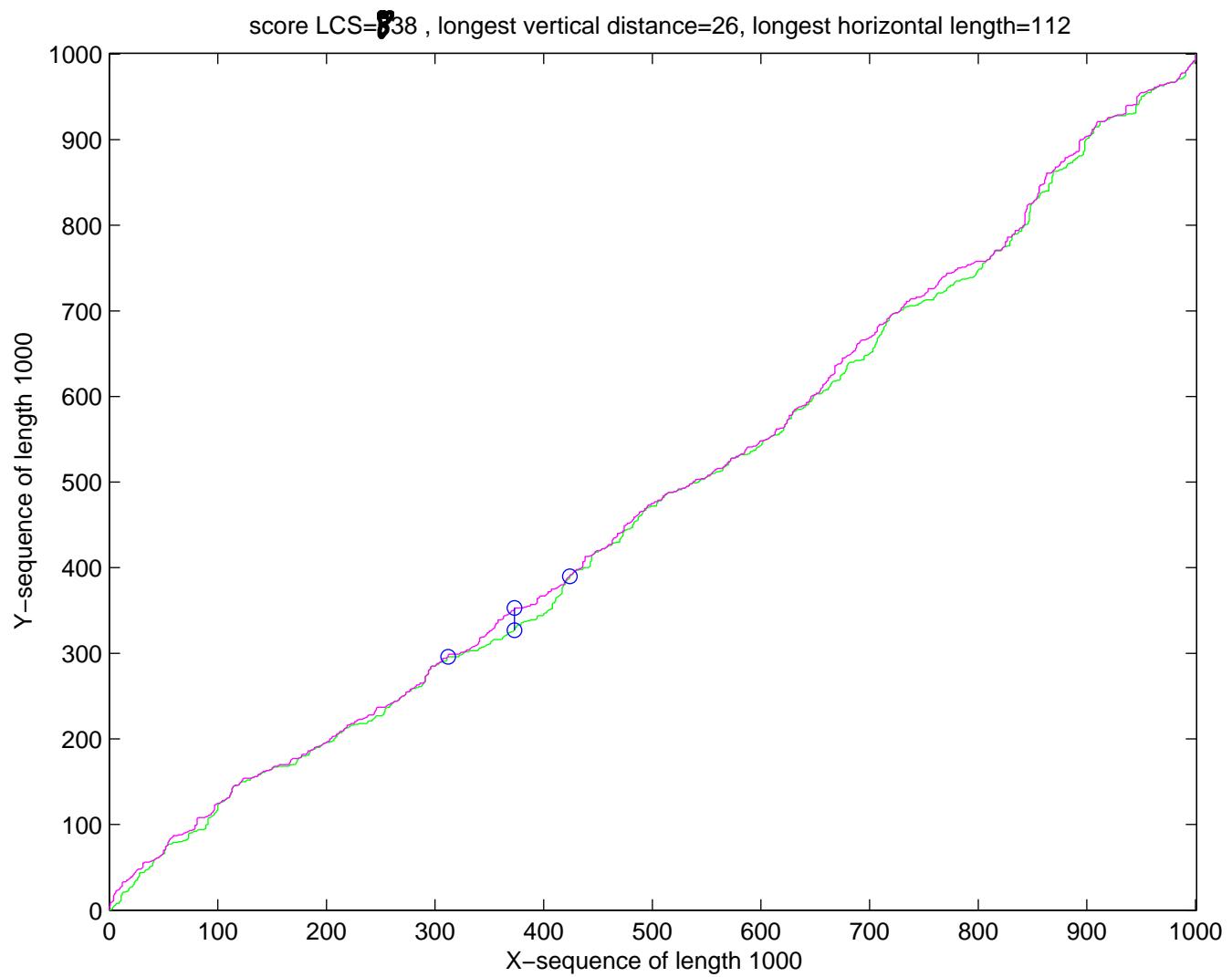
Such integers, n_0, n_1, \dots, n_d always exists.

Generality: Typically, the vast majority of the random lengths $r_i - r_{i-1}$ are close to \sqrt{n} .

Quantitatively, take $N = n^\alpha$, $0 < \alpha < 1$, let E_n be the random sets of alignments with lengths $r_i - r_{i-1}$ close to \sqrt{n} , then

$$P(E_n) \geq 1 - e^{-n^{1-\alpha} (1 + \log(1 + n^\alpha))}.$$

When representing alignments in $[N]^2$, we stay close to the diagonal at a macroscopic level.



To return to the estimation of $\text{Var} T$, we condition again, this time with respect to E_n . When E_n holds only a "small number" of terms in the sums defining the variance. When E_n does not hold, then many terms in the sum but exponentially small probability that E_n does not occur.

$$\text{Var} T \leq C \left(n e^{2^{-n^{1-\alpha}} (1 + \log(1+n^\alpha))} + n^{1+\alpha} + n^{\frac{3-2\alpha}{2}} \sqrt{\log n^\alpha} \right)$$

Similarly for $\text{Var} \tilde{T}$.

Take $\alpha = 4/5$, to get

$$d_K\left(\frac{LC_n - \mathbb{E}LC_n}{\sqrt{\text{Var } LC_n}}, \mathbb{Z}\right)$$

$$\leq C \sqrt{\frac{n^{9/5}}{\sigma^2}} + \frac{n}{\sigma^3}$$

But, $\sigma^2 \geq Kn$, so

$$\leq \frac{C}{K} \sqrt{\frac{1}{n^{1/5}}} \log n + \frac{1}{K^{3/2} \sqrt{n}}$$

$$\leq C \sqrt{\frac{1}{n^{1/5}}} \log n \xrightarrow[n \rightarrow +\infty]{} 0.$$

- Sublinear lower estimate on σ^2 is enough, $\sigma^2 > Kn^{9/10 + \varepsilon}$ is enough

- All the cases where the variance is lower estimated, we have $\sigma^2 \geq Kn$!

Thank You!

100%

1