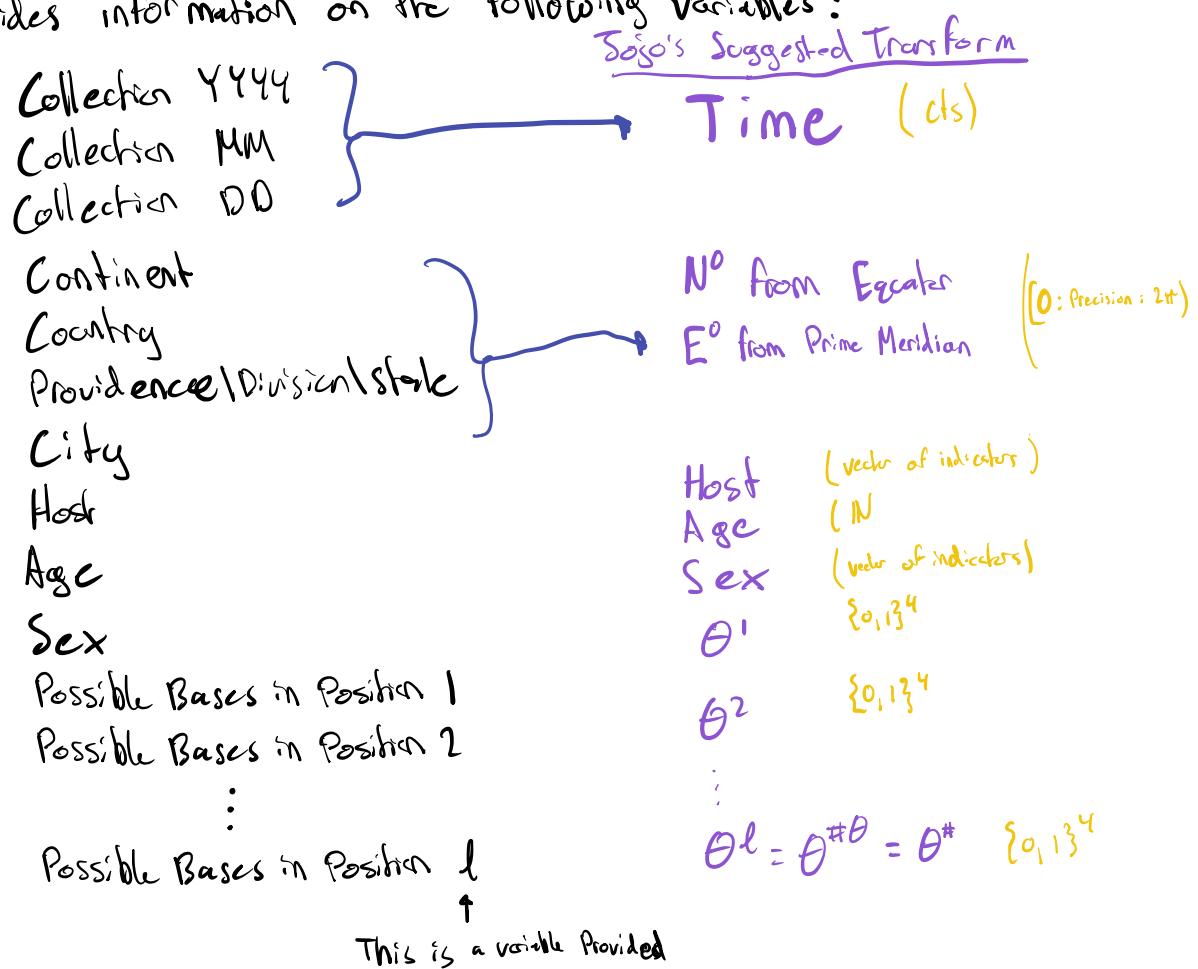


I will assume that every genomic or genetic DNA sequence is of the form, $b_1 b_2 \dots b_n$, where $H_i \in \{1, \dots, n\}$ $b_i \in B := \{A, C, G, T\}$.

$$B^+ = \bigcup_{n \in \mathbb{N}_{>0}} B^n, \quad V_B = \text{sp}\{e_A, e_C, e_G, e_T\} \text{ where } \begin{aligned} e^1 \vec{1} &= e_A, & e^2 \vec{1} &= e_C, \\ \text{over } A \text{ or } C && e^3 \vec{1} &= e_G, & e^4 \vec{1} &= e_T. \end{aligned}$$

This defines coordinates.

I am using $m \approx 2.3$ million observations from GSA ID. Each observation provides information on the following variables:



This is when i focus on the simplification where we observe $\theta = (\theta^1, \theta^2, \dots, \theta^*)$ (over time). Its interesting, the sequence part of the "data" (format) is really only a sequence when the read was "perfect" otherwise we are provided a parameter.

Onto the Models!

A fun Model!

$$F = \text{span}_{\beta \in \mathbb{B}} \left\{ e_\beta = e_{\beta_1} \otimes \dots \otimes e_{\beta_n} \mid \beta \in \mathbb{B}^{n \times n}, n \in \mathbb{N} \right\} = \left\{ c_\beta e_\beta \mid \beta \in \mathbb{B}^{n \times n}, n \in \mathbb{N}, c_\beta \in \{0,1\} \right\}, \Omega = \bigcup_{\beta \in \mathbb{B}} F$$

$$(H) = \left\{ \theta \in \{0,1\}^{\# \mathbb{B}^{n \times n}} \mid n \in \mathbb{N} \right\} = \left\{ (1), \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \dots \right\}$$

$$M_\Theta = \left(\bigoplus_{n \in \mathbb{N}} V_\Theta^{\otimes n} \right)^* = \bigcup_{n \in \mathbb{N}} L(V_\Theta, \dots, V_\Theta; \mathbb{R})$$

$P(\cdot | \theta) = P_\theta: H \rightarrow M_\Theta$ given by:

$$\theta = \begin{pmatrix} \xi_A^1 & \xi_A^2 & \xi_A^* \\ \xi_C^1 & \xi_C^2 & \xi_C^* \\ \xi_G^1 & \xi_G^2 & \dots & \xi_G^* \\ \xi_T^1 & \xi_T^2 & & \xi_T^* \end{pmatrix} \mapsto \left(\frac{\otimes \xi}{|\xi^1| \cdot \dots \cdot |\xi^*|} \right)^* \quad \text{Dual} \quad \text{i.e. } P_\theta := * \circ \theta \left(\frac{0^1}{10^1}, \dots, \frac{0^*}{10^*} \right)$$

In Particular:

$$\frac{\otimes \xi}{|\xi^1| \cdot \dots \cdot |\xi^*|} = \frac{\sum_{\beta \in \mathbb{B}^{n \times n}} \xi_{\beta_1}^1 \otimes \dots \otimes \xi_{\beta_n}^*}{|\xi^1| \cdot \dots \cdot |\xi^*|} = \frac{\sum_{\beta \in \mathbb{B}^{n \times n}} \xi_\beta}{|\xi^1| \cdot \dots \cdot |\xi^*|}$$

and the number of β such that $\xi_{\beta_1}^1 = 1$ and $\xi_{\beta_2}^2 = 1$ and ... and $\xi_{\beta_{n*}}^* = 1$ is $|\xi^1| \cdot \dots \cdot |\xi^*|$

$$\text{So } P_\theta = \frac{1}{|\xi^1| \cdot \dots \cdot |\xi^*|} \sum_{\substack{\beta \in \mathbb{B}^{n \times n} \\ \text{basis element}}} \xi^\beta \text{ satisfies } P_\theta(\Omega) = |\xi^1| \cdot \dots \cdot |\xi^*| \frac{1}{|\xi^1| \cdot \dots \cdot |\xi^*|} = 1.$$

So now we have a Random Probability Measure !!

It would be completely random if the following condition is satisfied (Kingman
Completely Rand. Measures)

If A_1, A_2, \dots, A_n mutually disjoint, then $P_\theta(A_1), P_\theta(A_2), \dots, P_\theta(A_n)$ are independent R.V.s.

I don't believe this is the case because you can take a Partition

$$\bigcup_{i=1}^n A_i = \Omega, \text{ the joint distribution is not independent}\\ \text{because } P_\theta(A_n) = 1 - \sum_{i=1}^{n-1} P_\theta(A_i).$$

Is this condition of Kingman's well-defined? With respect to what measure should one verify this claim with?

Also this doesn't seem to be true of any probability measure ...

Anyways if you want "Completely Random Measures" in our context just don't normalize $\otimes \theta$ in the defn of P_θ . I'd rather get probability measures than do this.

Last but, one could consider Random Variables as elements of $\text{span}_K \{ e_\beta \mid \beta \in \mathbb{B}^{n \times n}, n \in \mathbb{N} \}$, where K is a field.

Another "Model"!

$$\mathcal{H} = \{ \theta \in \mathbb{R}^{B^+} \mid |\theta|_1 = 1 \}$$

Note: B^+ must be orderable for this to be valid: Partially order by length and then lexicographically.

$$M_{\mathcal{H}} = \bigcup_{n \in \mathbb{N}} L(v_0, \dots, v_n; [0, 1])$$

$$P_\theta \text{ defined by: } (f_w)_{w \in B^+} \mapsto \sum_{w \in B^+} f_w \theta^w$$

More Random Probability Measures!

What happens if we assume sequences are exchangeable?

Suppose we suppose a DNA sequence is exchangeable & define P_θ s.t. P_θ exchangeable $\forall \theta$.

Then, our Probability Distribution, $P_\theta = \sum_{n \in \mathbb{N}} P_\theta^{(n)}$, must satisfy:

$$\exists P_\theta^{(n)} = P_\theta^{(n)} \quad \forall \alpha \in G^{(n)} \quad \leftarrow \text{Symmetric group on } n \text{ letters.}$$

Since $\{e_A, e_C, e_G, e_T\} \otimes \dots \otimes \{e_A, e_C, e_G, e_T\}$

$$\begin{aligned} &= \{e_{\omega l} \mid \omega \in B^n\} \\ &= \{g(e_A^{\otimes n_1} \otimes e_C^{\otimes n_2} \otimes e_G^{\otimes n_3} \otimes e_T^{\otimes n_4}) \mid g \in G^{(n)}, (n_1, n_2, n_3, n_4) \in \mathbb{N}^4, n_1 + n_2 + n_3 + n_4 = n\} \\ &= \bigcup_{\substack{(n_1, n_2, n_3, n_4) \in \mathbb{N}^4 \\ n_1 + n_2 + n_3 + n_4 = n}} G^{(n)}(e_A^{\otimes n_1} \otimes e_C^{\otimes n_2} \otimes e_G^{\otimes n_3} \otimes e_T^{\otimes n_4}) \end{aligned}$$

is a basis for the Domain of $P_\theta^{(n)}$, if it is defined by the vals it takes on this basis.

$$\begin{aligned} \text{So } P_\theta^{(n)} &= \sum_{\alpha \in \mathbb{N}^4} \sum_{\substack{e_{\omega l} \in G^{(n)}(e_A^{\otimes n_1} \otimes e_C^{\otimes n_2} \otimes e_G^{\otimes n_3} \otimes e_T^{\otimes n_4}) \\ |\alpha|=n}} P_\theta^{(n)}(e_{\omega l}) = \sum_{\alpha \in \mathbb{N}^4} |G^{(n)}(e_A^{\otimes n_1} \otimes e_C^{\otimes n_2} \otimes e_G^{\otimes n_3} \otimes e_T^{\otimes n_4})| P_\theta^{(n)}(e_A^{\otimes n_1} \otimes e_C^{\otimes n_2} \otimes e_G^{\otimes n_3} \otimes e_T^{\otimes n_4}) \\ &= \sum_{\alpha \in \mathbb{N}^4} \alpha_1! \alpha_2! \alpha_3! \alpha_4! P_\theta^{(n)}(e_A^{\otimes n_1} \otimes e_C^{\otimes n_2} \otimes e_G^{\otimes n_3} \otimes e_T^{\otimes n_4}) \xrightarrow{\substack{\text{By Dim=Dim} \\ \text{This is not an equality.}}} \sum_{\alpha \in \mathbb{N}^4} \alpha_1! C_{\theta, \alpha}^n \frac{n_1! n_2! n_3! n_4!}{n_1! n_2! n_3! n_4!} = \sum_{\alpha \in \mathbb{N}^4} \alpha_1! C_{\theta, \alpha}^n Z^\alpha \end{aligned}$$

$$\text{In Particular, } P_\theta^{(n)}(\Omega) = \sum_{\alpha \in \mathbb{N}^4} \alpha_1! C_{\theta, \alpha}^n$$

Since $P_\theta = \sum_{n \in \mathbb{N}} P_\theta^{(n)}$ is a probability Measure: $1 = P_\theta(\Omega) = \sum_{n \in \mathbb{N}} \sum_{\alpha \in \mathbb{N}^4} \alpha_1! C_{\theta, \alpha}^n$

This follows from concluding we have the isomorphism in yellow for each n , since n was arbitrary.

Thus: Every exchangeable Probability Measure is isomorphic to some polynomial (with coeffs in $\text{Range}(P_\theta)$).

All of the requirements of this polynomial follow from those on the exchangeable Prob. Measure.

To call a probability density or distribution exchangeable is to assert that it is G -invariant. Since exchangeable probability distributions are equivalent to polynomials, it is natural to want an irreducible representation. These polynomials are not constrained to any degree, so let's consider them as elements of $\mathbb{R}[z_1, z_2, z_3, z_4]$, which has a G -invariant subring $\{c z_1^n z_2^m z_3^p z_4^q \mid c \in \mathbb{R}, n, m, p, q \in \mathbb{Z}\} = \langle z_1^n z_2^m z_3^p z_4^q \mid n, m, p, q \in \mathbb{Z} \rangle = \langle z_1 z_2 z_3 z_4 \rangle$. We can take care of this by taking the polynomial ring in 3 variables, e.g. $\mathbb{R}[z_1 - z_2, z_2 - z_3, z_3 - z_4]$.

Interestingly, since the original polynomial did not change sign under the application of a transposition, these parameterizations don't either. Thus all terms are of even degree.

"Model" for exchangeable Sequences

$$\mathcal{H} = \{\theta \in \mathbb{R}^{(\mathbb{N})^4} \mid |\theta| = 1\} \text{ or } \{\theta \in \mathbb{C}^{(\mathbb{N})^4} \mid |\theta| = 1\}$$

$$M_{\mathcal{H}} = \left\{ \sum_{\alpha \in (\mathbb{N})^4} c_\alpha z^\alpha \mid c_\alpha \geq 0 \right\} = \mathbb{R}_+ [z_1, z_2, z_3, z_4]$$

$$P_\bullet \text{ given by } \theta \mapsto \sum_{\alpha \in (\mathbb{N})^4} \theta_\alpha z^\alpha$$

It would be appropriate to let: $\text{sort}(\theta) \sim \text{GEM}(\eta)$
 (This is more of a fact than a modelling choice).

A simple thing that should be incorporated.

$$\mathcal{H} = \{\theta \in \mathbb{R}^\infty \mid |\theta| = 1\}$$

$$M_{\mathcal{H}} = \cup L(v_0, \dots, v_B; [0, 1])$$

$$\theta \mapsto \left(\sum_{n \in \mathbb{N}} \theta_n \otimes I^{x^n} \right)^*$$

This is just a distribution over lengths of seqs. *discrete normal?*

Any (mixture) of a location/Scale dist. should do. I guess I think at it this way because $P = \sum P^{(n)}$ so Empirical Length distribution is. $(\|P^{(1)}\|, \|P^{(2)}\|, \dots, 0, \dots)$

For Messiness

$$\Phi = \left\{ (\alpha, (\phi_\eta)_{\eta=1}^B) \mid \alpha = \sum_1^{\#\theta} I_{\theta=1} \theta^i, \beta = \sum_1^{\#\theta} I_{\theta=1} \theta^i, L = \#\theta \right\}$$

$$M_\Phi = \mathbb{R}[z]$$

$$(\alpha, \phi) \mapsto z^\alpha (\sum \phi_\eta z^\eta)$$

Torsion? Lag?

$$\mathbb{H} = \bigcup_{n \in \mathbb{N}} \{\theta\}_{\theta \in \mathbb{R}^{4 \times n}}$$

$$M_{\mathbb{H}} = L(\mathbb{N}; V_B^{\otimes 4})$$

$$\theta \mapsto (\theta^{\max(1, i-4)} \otimes \dots \otimes \theta^i)_{i=1}^{\#\theta} = (\underbrace{S(\theta^{\max(1, i-4)} \otimes \dots \otimes \theta^i)}_{\text{src } \mathbb{R}[z]} \oplus \underbrace{A(\theta^{\max(1, i-4)} \otimes \dots \otimes \theta^i)}_{\text{src } \mathbb{R}[G_B]})_{i=1}^{\#\theta}$$

Anything on the closed simplex (w/ e_A, e_C, e_G, e_T as the extreme points)

$$\mathbb{H} = \left\{ \alpha = (\alpha^1, \dots, \alpha^n) = \begin{pmatrix} \alpha_x^1 & \alpha_x^2 & \dots & \alpha_x^n \\ \alpha_y^1 & \alpha_y^2 & \dots & \alpha_y^n \\ \alpha_z^1 & \alpha_z^2 & \dots & \alpha_z^n \end{pmatrix} \in [-1, 1]^{3 \times n} \mid \forall i \in \{1, \dots, n\} |\alpha_i^i| \leq 1, n \in \mathbb{N} \right\}$$

$$M_{\mathbb{H}} = \bigcup_{n \in \mathbb{N}} L(\tilde{V}_B, \dots, \tilde{V}_B; \mathbb{R}) \quad \text{where } \tilde{V}_B = \text{span}\{b_A - b_C, b_C - b_G, b_G - b_T\} \text{ (say).}$$

$$\alpha \mapsto \otimes \begin{pmatrix} \alpha_x^1 b_A + \alpha_x^2 b_C, \alpha_x^2 b_A + \alpha_x^3 b_C, \dots, \alpha_x^n b_A + \alpha_x^1 b_C \\ \alpha_y^1 b_A + \alpha_y^2 b_C, \alpha_y^2 b_A + \alpha_y^3 b_C, \dots, \alpha_y^n b_A + \alpha_y^1 b_C \\ \alpha_z^1 b_A + \alpha_z^2 b_C, \alpha_z^2 b_A + \alpha_z^3 b_C, \dots, \alpha_z^n b_A + \alpha_z^1 b_C \end{pmatrix}^*$$

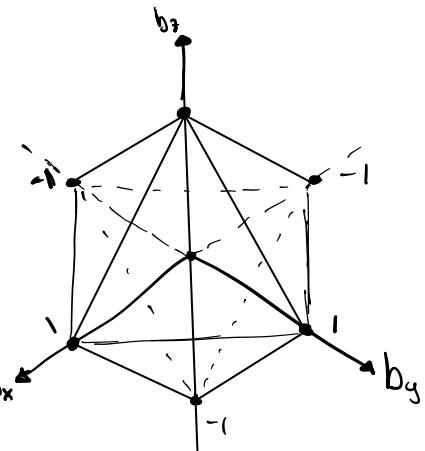
Any distribution on α_{b_1, b_2}^i induces one on G_B .

It probably makes most sense to let α_{b_1, b_2}^i be exponential or gamma or Laplace (w/ location 0).

Prob distributions on α_{b_1, b_2}^i , say w/ f, tells us the effect of sign changes, i.e. how $f(x)$ compares to $f(-x)$.

This induces a distribution on the transposition (b_1, b_2) , which has representation $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, a 2x2 real matrix or a complex number.

$\log \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} \frac{\pi i}{2} & -\frac{\pi i}{2} \\ \frac{\pi i}{2} & \frac{\pi i}{2} \end{pmatrix}$ means that real multiplicative relations have additive complex analogies. What other fun approaches are there using complex variables?



Some facts about Distributions taking $t \in G_B$ as an argument.

$$\text{Gamma}(\underbrace{\text{shape } k, \text{ scale } \theta}_{\text{i require this } \in \mathbb{R}}) \Rightarrow f(x; k, \theta) = \frac{1}{\Gamma(k)\theta^k} x^{k-1} e^{-x/\theta} \Rightarrow f(t; k, \theta) = \frac{t^{k-1} e^{-t/\theta}}{\Gamma(k)} \\ = \frac{e^{-t/\theta}}{\Gamma(k)} t^{k-1} e^t = f(t^{k-1}; 0, \theta)$$

Also for $c \in \mathbb{R}_{>0}$ $f(ct; k, \theta) = f((ct)^{k-1}; 0, \theta)$

$$\text{Inverse Gamma}(\text{shape } k, \text{ scale } \beta) \Rightarrow f(x; k, \beta) = \frac{\beta^k}{\Gamma(k)} x^{-k-1} e^{-\frac{\beta}{x}}$$

$$\Rightarrow f(t) = \frac{\beta^k}{\Gamma(k)} t^{-k-1} e^{-\frac{\beta}{t}} = \frac{\beta^k}{\Gamma(k)} (t^{-1})^{k+1} e^{-\beta t^{-1}} = \text{gamma}(t^{-1}, k+1, \frac{1}{\beta}) \\ = \text{gamma}(t^{-k-2}, 0, \frac{1}{\beta})$$

$$\text{Chi-Squared } (K \in \mathbb{N}_+) \Rightarrow f(x, k) = \frac{1}{\Gamma(\frac{k}{2}) 2^{k/2}} x^{\frac{k}{2}-1} e^{-\frac{x}{2}}$$

$$\Rightarrow f(x, 2t) = \text{Gamma}(x; k, 2) \Rightarrow \chi_{2t}^2(t) = \text{Gamma}(t^k; 0, 2)$$

I have sort of been assuming we may treat our data in time as a stochastic process. We can do this by:

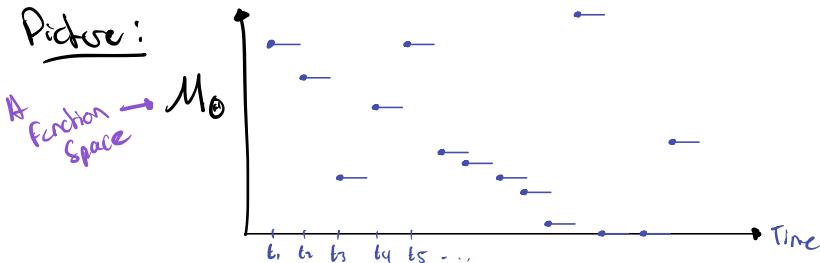
extending the random Probability Measure formulation to each point in time t_1, t_2, \dots, t_n by defining $P_{t_i} \in T_{t_i} (\bigoplus V_B^{\otimes n})^* \cong (T_{t_i} (\bigoplus V_B^{\otimes n}))^*$

i.e. the random measure determined by the data at a point in time.
 t_i is an element of the dual of the tangent space at t_i , which I think is called the cotangent space.

I would prove this but it follows from letting data take values in $\mathbb{R} \times \bigoplus V^{\otimes n}$ where \mathbb{R} is our proxy for time (taking \mathbb{C} could be useful and is appropriate because polynomials and permutations are holomorphic). It would be more reflective of the data to take "time" as $\{2011, 2020, 2021\} \times \{1, \dots, 12\} \times \{1, \dots, 31\}$. This perspective is valuable since some data points do not provide MM or DD.

Empirical distribution function up to time T : $F_T := \frac{\sum_{t \leq T} P_t}{\|\sum_{t \leq T} P_t\|}$

Picture:



Kingman's theoretical approach to genetics Problems is motivated by:

A geneticist who examines a sample of n gametes from a population of size N , and whose experimental techniques do not permit a labelling of the alleles he can distinguish, will have as data a member of ω_n .

" ω_n " is the set of all partitions of n , i.e. integer partitions of n

$$\text{In notation: } \omega_n := \{(a_1, \dots, a_n) \in \mathbb{N}^n \mid \sum_{j=1}^n j a_j = n\}$$

And a_1 is interpreted as the number of alleles represented once in the sample, a_2 alleles are represented twice in the sample, ..., a_n alleles are represented n times in the sample.

So if we observe a sample (of size n) and denote the set of observed alleles as $S = \{s_1, \dots\}$ and frequency of some s_i as f_s , then we may determine (a_1, \dots, a_n) by:

$$a_1 = \sum_{s \in S} 1_{f_s=1}, \quad a_2 = \sum_{s \in S} 1_{f_s=2}, \quad \dots, \quad a_n = \sum_{s \in S} 1_{f_s=n}$$

I am not sure what kinds of experimental techniques resulted in this kind of raw data. Anyways once one accepts this notion of a partition, there is absolutely no relation to the underlying subject or applied setting (unless this is reflective of the observation generating mechanism).

Regardless, modern techniques can provide arrays of sequences where each element of a sequence is a subset of nucleotides. Subsets of $B = \{A, C, G, T\}$ are observed because the sequencing technology cannot always determine the exact nucleotide at a location,

but in these instances it can sometimes "narrow down"

which nucleotides could be at that location. So the typical observed genomic or genetic sequence, α , of length $\# = l$, is an element of $\underbrace{2^B \times \dots \times 2^B}_{l \text{ times}} = (2^B)^{\times l}$. So each sequence may

be identified with a subset of exact sequences because

Set theory gives us: $\underbrace{2^B \times \dots \times 2^B}_{l \text{ times}} \cong \underbrace{2^{B \times \dots \times B}}_{l \text{ times}}$. For example,

$$\alpha = (\{A, C\}, \{G, T\}, \{T\}) \cong \{(A, G, T), (A, T, T), (C, G, T), (C, T, T)\} = x \alpha$$

Cartesian Product.

It is natural to wonder what distribution our data was drawn from (if there is one) or what the limiting distribution of our data is (if there is one).

Kingman takes a really neat approach to this when the sample is an element of ω_n (the set of integer partitions). The sample, a_i , is an element of ω_n because a_i denotes the number of alleles represented i times in the sample, so the sample size is $\sum i a_i$, which equals n since $a \in \omega_n$. Kingman provides a notion of a consistent family of distributions, called a "partition structure"

$P = (P_1, P_2, P_3, \dots) = (P_n \in \Pi_n)_{n \in \mathbb{N}}$ is a partition structure if for all $m < n$, P_m may be obtained by random sampling from P_n (Π_n denotes the set of probability distributions on ω_n).

Kingman then provides a construction for some Partition Structures

A very large class of partition structures can be constructed in the following way. Consider an infinite population divided into a countable number of sub-populations labelled (say) by colours 1, 2, 3, ..., and suppose that the proportion of the population which is of colour i is x_i , where

$$x_i \geq 0, \quad \sum_{i=1}^{\infty} x_i = 1. \quad (2.8)$$

For $n \geq 1$ and $\pi \in \omega_n$, let $P_n(\pi)$ be the probability that a random sample of size n has a_1 colours represented once, a_2 twice, and so on. Thus $P_n(\pi)$ is a function

$$P_n(\pi) = \phi_{\pi}(x) \quad (2.9)$$

of $x = (x_1, x_2, \dots)$, given explicitly by

$$\phi_{\pi}(x) = n! \prod_{j=1}^n (j!)^{-a_j} \sum x_1^{v_1} x_2^{v_2} \dots, \quad (2.10)$$

where the sum extends over all sequences (v_1, v_2, \dots) of integers $0 \leq v_i \leq n$ such that, for $1 \leq j \leq n$, exactly a_j of the v_i are equal to j .

As before, it is clear (and easy to check directly using (2.8)) that (2.9) defines a partition structure for any x satisfying (2.8). Even more general structures can be constructed by allowing x to be random. Thus let \mathcal{A} be the set of all sequences x satisfying (2.8), which may conveniently be topologized as a subset of the compact metrisable cartesian product $[0, 1]^{\omega}$ whose elements are sequences x satisfying $0 \leq x_i \leq 1$. Then ϕ_{π} is a continuous function from \mathcal{A} to $[0, 1]$, and if μ is any probability measure on (the Borel subsets of) \mathcal{A} ,

$$P_n(\pi) = \int_{\mathcal{A}} \phi_{\pi} d\mu \quad (2.11)$$

defines a member P_n of Π_n . Because (2.9) satisfies (2.7), we have

$$\begin{aligned} (\sigma_{mn} P_n)(\rho) &= \sum_{\pi \in \omega_n} \int_{\mathcal{A}} \phi_{\pi} d\mu S(\pi, \rho) \\ &= \int_{\mathcal{A}} \left(\sum_{\pi \in \omega_n} \phi_{\pi} S(\pi, \rho) \right) d\mu \\ &= \int_{\mathcal{A}} \phi_{\rho} d\mu = P_m(\rho), \end{aligned}$$

so that (2.11) defines a partition structure for any choice of μ .

Results of Watterson (1976) show that the Ewens structure is itself expressible in the form (2.11). To do this, μ is taken as the Poisson-Dirichlet distribution $\mathcal{PD}(\theta)$

As statisticians this is a compelling structure. However, this approach is far too abstract when the genomic or genetic data is not an integer partition. This is the case.

I propose a genomic method first.

① Let: \mathbb{K} be a field.

$$V_{YY44} = \text{Span}_{\mathbb{K}}\{e_{2019}, e_{2020}, e_{2021}\}$$

$$V_{MM} = \text{Span}_{\mathbb{K}}\{e_{301}, \dots, e_{307}\}$$

$$V_{DD} = \text{Span}_{\mathbb{K}}\{e_{01}, \dots, e_{33}\}$$

$$V_{Continent} = \text{Span}_{\mathbb{K}}\{e_{Antarctica}, \dots, e_{South America}\}$$

$$V_{Country} = \text{Span}_{\mathbb{K}}\{e_{Country} \mid \text{Country } \in \text{Countries in the world}\}$$

$$V_{City} = \text{Span}_{\mathbb{K}}\{e_{City} \mid \text{City } \in \text{Cities in the world}\}$$

$$V_{Host} = \text{Span}_{\mathbb{K}}\{e_h \mid h \in \text{Hosts}\}$$

$$V_{Age} = \text{Span}_{\mathbb{K}}\{e_{Age} \mid \text{Age } \in \{0, \dots, 100\}\}$$

$$V_{SEX} = \text{Span}_{\mathbb{K}}\{e_{Male}, e_{Female}\}$$

$$V_B = \text{Span}_{\mathbb{K}}\{e_A, e_C, e_G, e_T\}$$

$$\underline{\text{Require: }} \|ex\|_1 = 1 \quad \forall x \in \{2019, 2020, 2021, \dots, A, C, G, T\}$$

② Let: Each observation, X , be an element of

$$V_{Info} := V_{YY44} \times V_{MM} \times V_{DD} \times V_{Continent} \times V_{Country} \times V_{City} \times V_{Host} \times V_{Age} \times V_{SEX} \times \bigcup_{l \in \mathbb{N}_+} V_B^{\times l}$$

$$= \bigcup_{l \in \mathbb{N}_+} V_{YY44} \times V_{MM} \times V_{DD} \times V_{Continent} \times V_{Country} \times V_{City} \times V_{Host} \times V_{Age} \times V_{SEX} \times V_B^{\times l}$$

$$= \bigcup_{l \in \mathbb{N}_+} V_{YY44} \times V_{MM} \times V_{DD} \times V_{Continent} \times V_{Country} \times V_{City} \times V_{Host} \times V_{Age} \times V_{SEX} \times \underbrace{V_B \times \dots \times V_B}_{l \text{ times}}$$

$$= (V_{YY44} \times V_{MM} \times V_{DD} \times V_{Continent} \times V_{Country} \times V_{City} \times V_{Host} \times V_{Age} \times V_{SEX} \times \{1\}) \bigcup \left(\bigcup_{l \in \mathbb{N}_+} V_{YY44} \times V_{MM} \times V_{DD} \times V_{Continent} \times V_{Country} \times V_{City} \times V_{Host} \times V_{Age} \times V_{SEX} \times V_B \times \dots \times V_B \right)$$

$$\underline{\text{Require: }} \|X_u\|_1 = 1 \quad \text{if } \text{Dim}(V_u) \text{ unbounded, for } u \in \text{Vars} = \{YY44, MM, DD, \text{Continent}, \text{Country}, \text{City}, \text{Host}, \text{Age}, \text{SEX}, B_1, B_2, \dots\}$$

$$\underline{\text{Note: }} V_{Info} \subseteq \bigcup_{l \in \mathbb{N}_+} V_{YY44} \times MM \times DD \times \text{Continent} \times \text{Country} \times \text{City} \times \text{Host} \times \text{Age} \times \text{SEX} \times B^{\times l}$$

③ Denote the data by \mathbb{X} and the sample size by m .

Observe: $\mathbb{X} \in V_{Info}^{x m}$, so $\mathbb{X} = (\mathbb{X}^1, \mathbb{X}^2, \dots, \mathbb{X}^m)$ where $\mathbb{X}^i \in V_{Info}$ for $i \in \{1, \dots, m\}$

Note: $\#\mathbb{X} = m$ & It is possible that $|\{\mathbb{X}^i \mid i \in \{1, \dots, m\}\}| \neq m$

③ The empirical density of an observation is $(\otimes x)^*$.

If $|X| = 1$, then $(\otimes x)^*$ is the empirical distribution of the observation.

④ The average Density of the observations is $\frac{1}{\#X} \sum_{i=1}^{\#X} (\otimes x_i)^*$

If $|X^i| = 1 \forall i = 1, \dots, \#X$, then $\frac{1}{\#X} \sum_{i=1}^{\#X} (\otimes x_i)^*$ is the empirical d.f., Denoted \bar{F} .

④.1) If you would like, define a fn. $w: V_{\text{vars}} \rightarrow \mathbb{R}$ and then the density w.r.t. w is $\sum_{i=1}^{\#X} w(x_i) (\otimes x_i)^*$.

⑤ A process of \bar{F} is defined by a non-empty subset $S \subset \text{Vars}$, and the isomorphism

$\bar{F} \cong V_S \otimes V_{\text{vars} \setminus S}^*$. This satisfies the requirements of a vector field given by the basis of $V_{j_1} \otimes \dots \otimes V_{j_{\text{last}}}$ where the points are the indices of the basis. The value of this vector field at a point is an element of $V_{\text{vars} \setminus S}^* = L(V_{\text{vars} \setminus S}; \mathbb{K})$, in particular

\bar{F}_α where $\alpha \in \bigtimes_{j \in S} \text{Vals}(j)$ may be defined by $\bar{F}_\alpha = \frac{\partial}{\partial x} \bar{F}^*$ which really just amounts to the sum of duals of basis elements whose indices are equal to α when projected onto the variables of S . Weight these however you please.

In epidemiological settings it is probably desirable to accommodate genomic and non-genomic data at the same time because we do not get genomic data from every infected individual. More broadly, we may wish to incorporate population data with genomic data.

We can achieve both with the above formalism. By construction we have permitted a genomic observation to have 0 nucleotide bases, or equivalently, the sequence observed is in $B^{x^0} = \{\emptyset\}$. The tensor product maps this into the field \mathbb{K} . Letting it map to $1_{\mathbb{K}}$ is really the only sensible thing to do since we'd hope that $\otimes B^{x^l} = \otimes B^{x^0} \otimes B^{x^l}$. In which case, a observation can be interpreted as population obs \otimes genomic obs where population obs lives in a covariate vector space over \mathbb{K} and the genomic obs lives in $\otimes B^{x^l}$ for some $l \in \mathbb{N} = \{0, 1, \dots\}$.

In retrospect we probably want to express some certainty or uncertainty about what the true measure is, if there is one. If we have reason to believe our observations are drawn from a "true" distribution that changes over time or with other variables (observed or not), then expressing uncertainty over an indexed collection of measures is also of interest. We can achieve this by expressing priors over the θ s in the various "models" above. I recognize these are very explicit notions of ambiguity and that some more flexibility or ambiguity given by a fewer number of parameters might offer useful ways to "turn knobs" when we are restricted to a number of bands.

Below are some ideas:

$$\text{Let } \theta_b^n \sim \begin{cases} 1 & \text{w.p. } x_b \\ 0 & \text{w.p. } 1-x_b \end{cases} \quad \forall n \in \{1, \dots, \# \theta\} \quad \forall \text{ observed } \theta$$

What if we take $\theta_b^n \sim f(x) = \frac{\theta}{(1-x)^{\theta-1}} dx, \theta >$

Let $\theta_b^n \sim \text{Erlang}(k, \lambda)$ CDF: $F(x) = 1 - \sum_{j=0}^{k-1} \frac{1}{j!} e^{-\lambda x} (\lambda x)^j$

Then θ

Kingman is also desirous a "Partition Structure", which he defines as a consistent family of partition structures (P_1, P_2, \dots) such that when $m < n$, P_m may be obtained from P_n by randomly sampling m things from P_n without replacement. (I need to mention Kingman's view of a partition vs. ours). We can provide this by Hypergeometric Distribution:

$$P(m_1, \dots, m_n) = \frac{\prod_{i=1}^n \binom{n_i}{m_i}}{\binom{n}{m}} = \frac{\prod_{i=1}^n \frac{n_i!}{m_i!(n_i-m_i)!}}{\frac{n!}{m!}} = \frac{m!}{n!} \prod_{i=1}^n \frac{n_i!}{m_i!(n_i-m_i)!}$$

So if we start with P_n , Kingman desires $\frac{\partial^{\alpha}}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}} P_n \propto P_m$ where $\frac{n-m}{n} = \alpha$

we can provide him this:

$$\begin{aligned} \text{Let } P_m &= \left(\frac{m!}{n!} \prod_{i=1}^n \frac{1}{(n_i-m_i)!} \right) \frac{\partial^{\alpha}}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}} P_n \\ &= \frac{(n - |\alpha|)!}{n!} \cdot \frac{1}{\alpha!} \frac{\partial^{\alpha}}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}} P_n \\ &= \frac{(n - |\alpha|)!}{n! \alpha!} \frac{\partial^{\alpha}}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}} P_n \end{aligned}$$

$$\bigcup_{l \in \mathbb{N}} \mathbb{N}^{x l} = \bigcup_{l \in \mathbb{N}} \bigcup_{\substack{\phi \in \mathbb{N}^{x l} \\ i \in \phi \Rightarrow i \in \phi_j}} G_i \phi$$

Some Geometric Models

Data is $X = (X^1, \dots, X^n)$

Supp up to & including $j-1$ is $H_{j-1} = \{x^i \mid i \leq j-1\}$

Bernoulli Sequence (Not Exchangeable), $\theta > 0$

$$\left. \begin{array}{l} X^i \notin H_{j-1} \text{ w/ prob. } \frac{\theta}{\theta+j-1} \\ X^i \in H_{j-1} \text{ w/ prob. } \frac{j-1}{\theta+j-1} \end{array} \right\} P_X = \prod_{i=1}^{\#X} \frac{\theta \mathbb{1}(x^i \notin H_{j-1}) + (j-1) \mathbb{1}(x^i \in H_{j-1})}{\theta + j - 1}$$

note: $\mathbb{1}(x^i \notin H_{j-1}) = \prod_{i \neq j} \mathbb{1}(x^i \neq x^j)$

Weird Bernoulli Sequence (exchangeable), $\theta > 0$

$$\left. \begin{array}{l} X^i \notin H_{j-1} \text{ w/ density } \frac{\theta}{\theta+j-1} \\ X^i \in H_{j-1} \text{ w/ density } \frac{1}{\theta+j-1} \end{array} \right\} f_X = \prod_{i=1}^{\#X} \frac{\theta \mathbb{1}(x^i \notin H_{j-1}) + \mathbb{1}(x^i \in H_{j-1})}{\theta + j - 1} = \frac{\theta^{|\{H_{j-1}\}|}}{\theta \cdots (\theta + \#X - 1)}$$

So the R.V. is Really $|H_{j-1}| = |\{x^i \mid i=1, \dots, \#X\}|$ = number of unique observations.

Now, $\sum_{u=1}^{\#X} \frac{\theta^u}{\theta \cdots (\theta + \#X - 1)} = 1$ if $\#X \geq 2$

However, $\sum_{g \in G_{\#X}} \frac{\theta^{l(g)}}{\theta \cdots (\theta + \#X - 1)} = 1$ where $cls = \lambda \vdash n$ is the cycle decomposition of g

$\lambda \vdash n$ iff, $\lambda \in \mathbb{N}^n$ & $\sum_i i \lambda_i = n$

so $|\lambda|$ is number of parts when $l(g)$ is number of cycles.

λ_i may be interpreted as "number of unique vals occurring i times in $\#X$ "
 or
 "number of classes of size i "

$\sum_{g \in G_{\#X}} \frac{\theta^{l(g)}}{\theta \cdots (\theta + \#X - 1)}$ seems like a good place to start to find the d.f.

$$1 = \sum_{g \in G_{\#X}} \frac{\theta^{l(g)}}{\theta \cdots (\theta + \#X - 1)} = \sum_{u=1}^{\#X} \sum_{\substack{g \in G_{\#X} \\ l(g)=u}} \theta^{l(g)} = \sum_{u=1}^{\#X} \sum_{\substack{\lambda \vdash \#X \\ |\lambda|=u}} \sum_{\substack{g \in G_{\#X} \\ cls=\lambda}} \theta^{l(g)} = \sum_{u=1}^{\#X} \sum_{\substack{\lambda \vdash \#X \\ |\lambda|=u}} \frac{|\{g \in G_{\#X} \mid cls=\lambda\}|}{\theta \cdots (\theta + \#X - 1)} \theta^{|\lambda|}$$

$$= \sum_{u=1}^{\#X} \sum_{\substack{\lambda \vdash \#X \\ |\lambda|=u}} \frac{\#X!}{\prod_{i=1}^u i \lambda_i!} \theta^{|\lambda|}$$

$$\text{which gives } P(|H_{j-1}| = u) = \frac{\sum_{\substack{\lambda \vdash \#X \\ |\lambda|=u}} \theta^u}{\theta \cdots (\theta + \#X - 1)} = \frac{\binom{\#X!}{\prod_{i=1}^u i \lambda_i!} \theta^{|\lambda|}}{\theta \cdots (\theta + \#X - 1)}$$

This also gives us a density over partitions:

$$P(\text{cls}) = \frac{\#X!}{\prod_{i=1}^u i \lambda_i!} \frac{\theta^{|\lambda|}}{\theta \cdots (\theta + \#X - 1)}$$

The last 2 formulas can be found in Johnson, Kotz, Balakrishnan CH41,

I just wanted to derive them from $\sum_{g \in G_{\#X}} \frac{\theta^{l(g)}}{\theta \cdots (\theta + \#X - 1)} = 1$.

A Different Approach

It is much more intuitive to speak about the unique vals and their respective frequencies, e.g., "The sequence $\underline{\quad}$ occurs $\underline{\quad}$ times in the sample". We can work from this perspective with a K -tuple, α , where

$$\alpha = (\alpha_{s_1}, \dots, \alpha_{s_K}, \dots, \alpha_{s_n})$$

Frequency of s_j in the sample.
Slot for Sequence s_j

In Particular, if $X = (X^1, \dots, X^n)$ is our data, a n -tuple of sequences, then $\alpha_{s_j} = \sum_{i=1}^n 1_{X^i=s_j}$. This approach should also be compelling because α is the tuple of coefficients on basis elements which are indexed by sequences, in notation:

$$\sum_{i=1}^n \otimes X^i = \sum_{j=1}^K \alpha_{s_j} e_{s_j} \quad \text{and} \quad \left(\sum_{i=1}^n \otimes X^i \right)^* = \sum_{j=1}^K \alpha_{s_j} \delta^{s_j}$$

If SEQS is the set of observed sequences, i.e. the basis of the support of $\sum_{i=1}^n \otimes X^i$, the above expressions are: $\sum_{s \in \text{SEQS}} \alpha_s e_s$ and $\sum_{s \in \text{SEQS}} \alpha_s \delta^s$

Note: $\text{SEQS} = H_{\#X} = H_n$

Motivation complete.

Returning to the weird bernoulli seq with $X^i \in H_{s_i-1}$ w/ density $\frac{\theta}{\theta + s_i - 1}$ • $X^i \in H_{s_i-1}$ w/ density $\frac{1}{\theta + s_i - 1}$ •

$$\text{The density function: } f_X = \frac{\theta^{|H_{\#X}|}}{\theta \cdots (\theta + |\#X| - 1)} = \frac{\theta^{\#\alpha}}{\theta \cdots (\theta + |\alpha| - 1)}$$

$$\text{The d.f. of } K \text{ given } n \text{ (same as before): } \frac{\sum_{\substack{X^i \in H \\ |\#X|=u}} \theta^u}{\theta \cdots (\theta + n - 1)} = \frac{P_{n,u} \theta^u}{\theta \cdots (\theta + n - 1)} \text{ where } P_{n,u} \text{ is the number of partitions of } n \text{ with } u \text{ parts.}$$

$$\begin{aligned} \text{The d.f. of } \alpha: P(\alpha) &= \frac{|\alpha|!}{\prod_{i=1}^{|H|} \left(\sum_{j=1}^{s_i} 1_{X^i=j} \right)!} \cdot \frac{\theta^{\#\alpha}}{\theta \cdots (\theta + |\alpha| - 1)} \\ &= \frac{|\alpha|!}{\prod_{i=1}^{|H|} \left(\sum_{j=1}^{s_i} 1_{X^i=j} \right)!} \cdot \frac{\theta^{\#\alpha}}{\theta \cdots (\theta + |\alpha| - 1)} = \frac{|\alpha|!}{|\text{Aut}(\alpha)|} \cdot \frac{\theta^{\#\alpha}}{\prod_{j=1}^{|H|} \alpha_j} \cdot \frac{\theta^{\#\alpha}}{\theta \cdots (\theta + |\alpha| - 1)} \\ &= \frac{1}{|\text{Aut}(\alpha)|} \frac{|\alpha|!}{\theta \cdots (\theta + |\alpha| - 1)} \prod_{j=1}^{|H|} \frac{\theta}{\alpha_j} = \frac{1}{|\text{Aut}(\alpha)|} \cdot \frac{|\alpha|!}{\theta \cdots (\theta + |\alpha| - 1)} \prod \left(\frac{\theta}{\alpha_j} \right) = \frac{1}{|\text{Aut}(\alpha)|} \cdot \frac{\theta + |\alpha|}{|\alpha|^{\theta}} B(\theta, |\alpha|+1) \prod \left(\frac{\theta}{\alpha_j} \right) \\ &\qquad\qquad\qquad \text{just notation.} \qquad\qquad\qquad \text{Eulerian integral of the second kind.} \qquad\qquad\qquad \text{Beta Fn.} \end{aligned}$$

The last line is just 3 ways to write the same thing. I will provide more on these as well as some other formulas for the integer partition version. Just so we're all on the same page:

	Data	Integer Partitions	Type Frequencies
Variable used	X	λ	α
Sample Size (n)	$\#X$	$\sum_{i=1}^{\#X} i \lambda_i$	$ \alpha $
Number of unique elements (x_i)	$ \{x_i : i \in \{1, \dots, \#X\}\} $	$ \lambda $	$\#\alpha$
Where Variable "Lives", Given n	$(\bigoplus_{i=1}^n V_B^{\otimes i})^{\otimes n}$	$\mathbb{N}^{n \times n}$	$\bigcup_{k=1}^n \mathbb{N}_s^{n \times k}$
Where Variable "Lives", Given k	$(\bigoplus_{i=1}^k V_B^{\otimes i})^{\otimes n}$	$\bigcup_{n \geq k} \mathbb{N}^{n \times n}$	$\mathbb{N}_s^{n \times k}$
Where Variable "Lives", Given $n \neq k$	$(\bigoplus_{i=1}^k V_B^{\otimes i})^{\otimes n}$	$\mathbb{N}^{n \times n}$	$\mathbb{N}_s^{n \times k}$
Where Variable "Lives", unknown $n \neq k$	$\bigcup_{n=1}^{\infty} (\bigoplus_{i=1}^k V_B^{\otimes i})^{\otimes n}$	$\bigcup_{n=1}^{\infty} \mathbb{N}^{n \times n}$	$\bigcup_{k=1}^{\infty} \mathbb{N}_s^{n \times k}$

... Lots of \mathbb{N} .

If we have perfect sequences, ^{ie every read is perfect,} we are in business!

This is not the case.

If we are willing to take $V_{2^B \setminus \{\emptyset\}}$ as the underlying vector space and let our random measures live in $\bigcup_{k \in \mathbb{N}} V_{2^B \setminus \{\emptyset\}}^{\otimes k}$, we are in business. This is not really reflective of actual real-life sequences and doesn't provide a distribution over true sequences \therefore . So this does not seem like the move. (we could provide a connection between $2^B \setminus \{\emptyset\}$ and B to solve this).

The Problem is: So long as we normalize each observation α will live in \mathbb{Q} . Even if we were willing to settle for working with λ ,

Kingman is of no help here.

The "It is Possible" Solution: Given the observation Parameter $\theta = \begin{pmatrix} \theta_1^1 & \theta_1^2 \\ \vdots & \vdots \\ \theta_t^1 & \cdots & \theta_t^m \end{pmatrix}$ from the beginning of this doc, construct the random probability measure as we did, i.e. $\left(\frac{\theta_1}{10^1}, \dots, \frac{\theta_t}{10^t}\right)^*$ and then multiply it by 12^m where m is the maximum length of a sequence in the sample. We pick 12 because it is $\text{LCM}\{1, 2, 3, 4\} = \text{LCM}(\text{Possible Vals of } 10^1)$. Each observation gets equal weight and now every coefficient is in \mathbb{N} so $\alpha \in \mathbb{N}_s^{n \times k}$ and $\lambda \in \mathbb{N}^{n \times n}$.

This is computationally exhausting because m is big.

Question: Can our nice formula $\frac{|\alpha|!}{|\text{Aut}(\alpha)|} \cdot \frac{|\alpha|!}{\theta \dots (\theta + |\alpha| - 1)} \prod_{i=1}^{\theta} i^{|\alpha_i|}$ hold over \mathbb{Q} ?

Note: $|\alpha|$ will always be in \mathbb{N} because $|\alpha|$ is the number of obs.

Intuition: The formula doesn't require anything of the elements of α , so why shouldn't it hold?

Logic: It was derived with the assumption that $\alpha \in \mathbb{N}^{n \times k}$. In particular

From the formula:

$$\frac{|\alpha|!}{\prod_{i=1}^k \left(\sum_{j=1}^{|\alpha_i|} 1_{\alpha_j=i} \right)!} \cdot \frac{\theta^{|\alpha|}}{\theta \dots (\theta + |\alpha| - 1)}$$

This term may be derived via a counting argument

Note: There might be a theorem or something or someone might know the answer. Do you know / have any ideas? I don't know, but gave it a shot.

Suppose: $\alpha \in \mathbb{Q}^{n \times k}$ and that $|\alpha| \in \mathbb{N}_0$.

Let: $r = \gcd(\alpha_1, \dots, \alpha_n)$ notation: $\alpha_i = \sum_{j=1}^n 1_{\alpha_j=i}$, $[a : x : b] = \{a, a+x, a+2x, \dots, b\}$

By assumption $\frac{|\alpha|}{r} \in \mathbb{N}_0$.

The term in question is: $\frac{|\alpha|!}{\prod_{i=1}^k i^{|\alpha_i|}}$. This fails over \mathbb{Q} by many counter examples.

The term of interest is: $\frac{|\alpha|!}{\prod_{i \in \mathbb{Q}_0} i^{|\alpha_i|}}$, i.e. can we define it over \mathbb{Q} ?

Pick your product!

$$\text{Guess: } \frac{|\alpha|!}{\prod_{i \in \mathbb{Q}_0} i^{|\alpha_i|}} = \frac{|\alpha|!}{|\text{Aut}(\alpha)| \prod_{i \in \mathbb{Q}_0} i^{|\alpha_i|}} = \frac{|\alpha|!}{|\text{Aut}(\alpha)| \prod_{i \in \mathbb{Q}} i^{|\alpha_i|}}$$

This looks good. How does the result compare with that of $\frac{|\alpha|}{r}$ and is this correspondence coherent and does it share the behavior of the definitions over \mathbb{N} ?

$$\begin{aligned} \frac{|\alpha|!}{\prod_{i \in \mathbb{Q}_0} i^{|\alpha_i|}} &= \frac{|\alpha|!}{|\text{Aut}(\alpha)| \prod_{i \in \mathbb{Q}_0} i^{|\alpha_i|}} = \frac{|\alpha|!}{|\text{Aut}(\alpha)| \prod_{i \in \mathbb{Q}} i^{|\alpha_i|}} = \frac{\left(\frac{|\alpha|}{r}\right) \left(\frac{|\alpha|-r}{r}\right) \left(\frac{|\alpha|-2r}{r}\right) \dots \left(\frac{|\alpha|-(m-1)r}{r}\right)}{|\text{Aut}(\alpha)| \frac{\prod_{i \in \mathbb{Q}} i^{|\alpha_i|}}{r^{|\alpha|}}} = \frac{\frac{|\alpha| \cdot (|\alpha|-r) \cdot \dots \cdot (|\alpha|-(m-1)r)}{r^{|\alpha|}}}{|\text{Aut}(\alpha)| \frac{\prod_{i \in \mathbb{Q}} i^{|\alpha_i|}}{r^{|\alpha|}}} \\ &= \frac{|\alpha| \cdot (|\alpha|-r) \cdot \dots \cdot (|\alpha|-(m-1)r)}{|\text{Aut}(\alpha)| \frac{\prod_{i \in \mathbb{Q}} i^{|\alpha_i|}}{r^{|\alpha|}}} \end{aligned}$$

$$\text{Check: } \beta \in \mathbb{N}_0^{n \times k} \quad \frac{|\beta|!}{\prod_{i \in \mathbb{Q}_0} i^{|\beta_i|}} = \frac{|\beta|!}{|\text{Aut}(\beta)| \prod_{i \in \mathbb{Q}_0} i^{|\beta_i|}} = \frac{|\beta|!}{|\text{Aut}(\beta)| \prod_{i \in \mathbb{Q}} i^{|\beta_i|}} = \frac{(m\beta)!}{|\text{Aut}(\beta)| \frac{\prod_{i \in \mathbb{Q}} i^{|\beta_i|}}{r^{|\beta|}}} = \text{For another day}$$

Following up on the "it is possible way":

Let $\alpha \in Q_s^{X^n}$ be observed.

So: $B := 12^m \alpha \in \mathbb{N}_s^{X^n}$

$$\text{Then: } P(B) = \frac{|B|!}{\prod_{i=1}^{|B|} B_i!} = \frac{|B|!}{|\text{Aut}(B)| \prod_{i=1}^{|B|} B_i!} = \frac{|B|!}{|\text{Aut}(B)| \prod_i B_i!} = \frac{|12^m \alpha|!}{|\text{Aut}(\alpha)| 12^{m^n} \prod_i \alpha_i!}$$
$$= \frac{|12^m \alpha|!}{|\text{Aut}(\alpha)| 12^{m^n} \prod_i \alpha_i!} = \frac{\prod_{i=1}^{|12^m \alpha|} i!}{|\text{Aut}(\alpha)| 12^{m^n} \prod_i \alpha_i!} = \frac{\prod_{i=1}^{m^n} i!}{|\text{Aut}(\alpha)| 12^{m^n} \prod_i \alpha_i!} = \prod_i$$

Cluster Coefficient & Model:

$n = 61,068,458$ genetic observations $\approx |\alpha|$

$K = 2,808,098 = \#\alpha$

$\max(\alpha) = 2183654$

Maybe Let:

$$\Delta = Q_{>0}^{X^n} \text{ gives us } \sum_{\lambda \in \Delta} c_\lambda X^\lambda$$