

Ezra Miller:

- Goal: statistics where sample space is a stratified space  $X$   
 e.g. trees, shapes, persistence diagrams (inf. dimensional  $\rightarrow$  not covered), diffusion tensors
- ordinary statistics: object = vector, sample space = vector space  
 methods: mean, variance, PCA, LLN, CLT
- stratified stat.: all of these can raise fundamental problems in geometric probability
- Def. A manifold is a second-countable, Hausdorff topological space  $M$  s.t. every point  $x \in M$  has a neighborhood  $\cong$  open ball in  $\mathbb{R}^n$ .

• chart:  $\pi: U \rightarrow M$   
 $U$  open ball in  $\mathbb{R}^d$

$M_{op} = \pi_a(U_a) \cap \pi_b(U_b)$



$\pi_a^{-1} \cap \pi_b^{-1}(\pi_a \cap \pi_b) \rightarrow \pi_a^{-1}(M_{op}) \cap \pi_b^{-1}(M_{op})$ ,  $\pi_{op}(x) = \pi_b^{-1} \circ \pi_a(x)$

$\pi_a$  topological / smooth / analytic / algebraic  $\rightsquigarrow$  topological / smooth / analytical manifold / algebraic variety  
 {all  $\pi_a$ } = atlas

- data examples: angles (circle), rotations ( $SO(3)$ ), lines ( $\mathbb{R}P^{n-1}$ ), subspaces (Grassmann)

• stratified spaces : Def. A topologically stratified space is a Hausdorff topological space  $X$  that is a disjoint union  $X = M_0 \cup M_1 \cup \dots \cup M_l$  of manifolds called strata

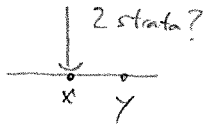
0) each stratum has finitely many connected components

1)  $M_i \cup \dots \cup M_l$  is closed in  $X \quad \forall i \leq l$

2)  $\forall x, y \in M_i: \exists$  homeomorphism  $\varphi: X \rightarrow X$  with

•  $\varphi(M_i) = M_i \quad \forall i$       •  $\varphi(x) = \varphi(y)$

→ points "look the same" on the same stratum

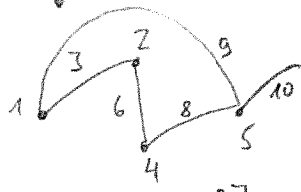
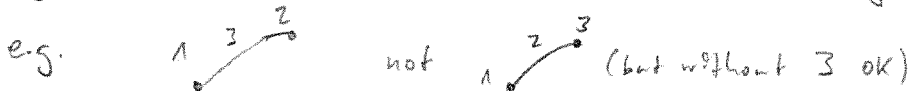


• examples:

dim 0 : finite set with discrete topology

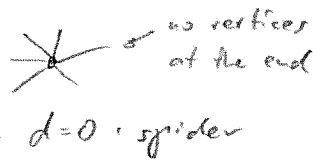
→ 1 disconnected stratum or lots of connected strata or same space but different stratifications

dim 1 : graphs no first vertices (condition 1) then edges

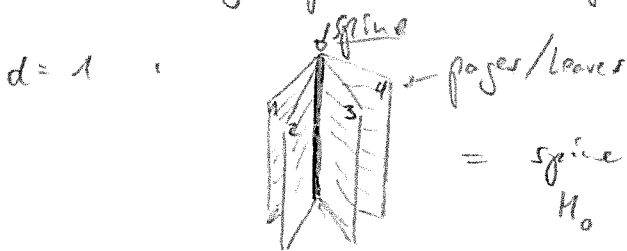


e.g. trees : connected, no cycles

• spider, or star:

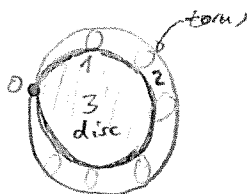


higher dim. 1 e.g. open book = spider  $\times \mathbb{R}^d$



GM example:

Goresky & MacPherson (?) → book



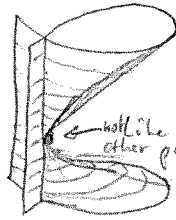
dim	0	1	2
#	1	1	2

polyhedron : intersection of finitely many half-spaces in  $\mathbb{R}^d$



dim	0	1	2
#	2	3	1

Whitney cusp:



$$\text{sing}(X) = \{ \}$$

$$\text{but } X = \text{sing}(X) \cup (X \setminus \text{sing}(X))$$

is not a stratification

shape spaces, e.g. face recognition ( $d=3$ )  
 $n$  (landmarks), labelled

$$d \left[ \begin{array}{c} | \\ | \\ | \\ | \\ | \end{array} \right] = A \in \mathbb{R}^{d \times n}, \quad A' \text{ same shape if } A' = gA,$$

$g \in G =$  group of transformations, e.g.  $G =$  rigid motions

or  $G$  could have rotations, translations, projective transformations, reflections

$\leadsto X = G \backslash \mathbb{R}^{d \times n} \leadsto$  (typically) algebraic variety  $\leadsto$  stratified (see below) after removing unstable orbits (keep stable and semi-stable orbits)

e.g.  $\Sigma_d^n = \mathbb{R}^{d \times n} \setminus \{\text{all columns equal}\} / \text{similarities}$

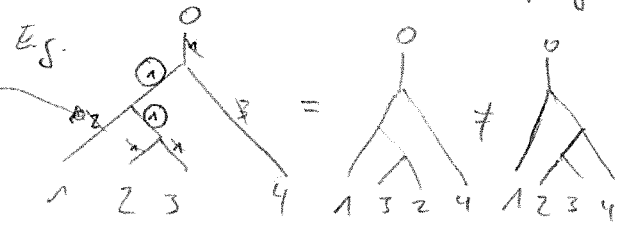
$d=2: \Sigma_2^n \cong \mathbb{C}P^{n-2}$ , else ( $d>2$ ): "non-free shapes"

Def:

A phylogenetic n-tree is a tree with <sup>n+1</sup> labelled leaves (vertices of degree 1)  $\leadsto$  lower dim. strata

and edge lengths  $> 0$  on internal edges (no pendant edges)

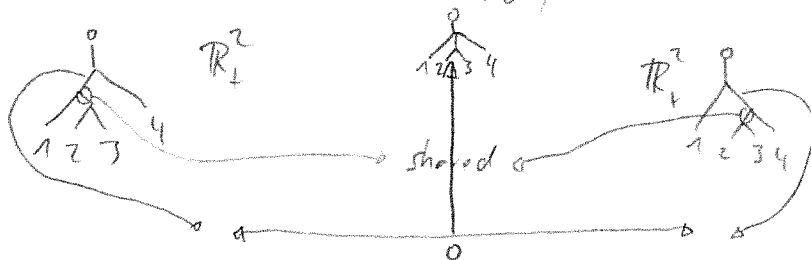
$\mathcal{T}_n = \{n\text{-trees}\}$



$\leadsto$  removing edge gives partition of leaf set into two connected components

$\leadsto$  fixed "combinatorics": specify internal edge lengths, e.g.  $\mathbb{R}_+^2$

$(\mathbb{R}_+ = \mathbb{R}_{>0}, \bar{\mathbb{R}}_+ = \mathbb{R}_{\geq 0}) \leadsto$  orbifolds  $\mathbb{R}_+^2$  are manifolds



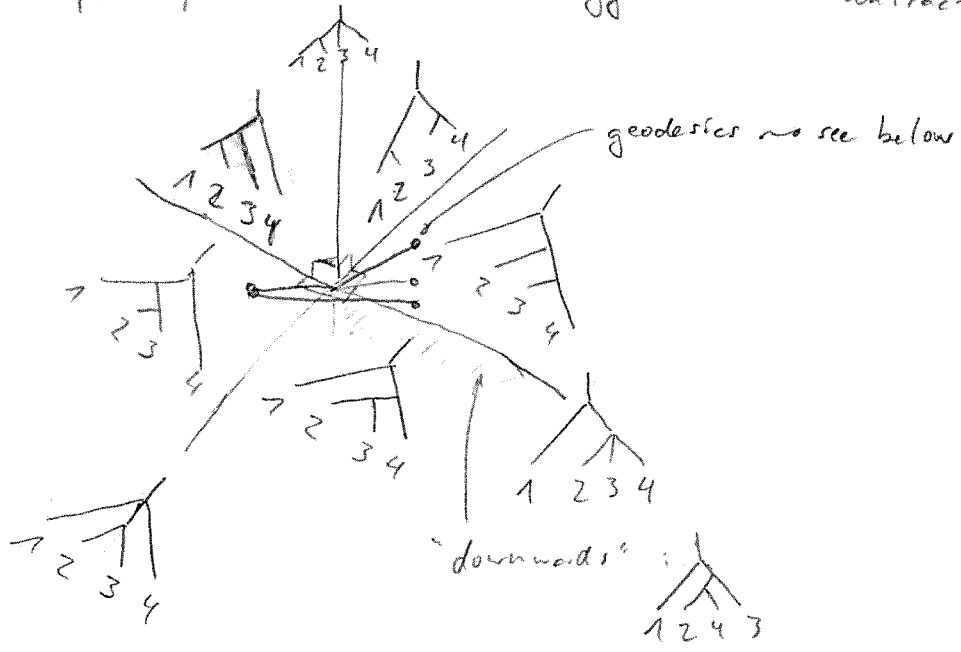
(Billera-Holmes - Vogtmann 1998)

facets of  $\mathcal{T}_n$ : orbifolds  $\bar{\mathbb{R}}_+^{n-2}$  in bijection with binary tree topologies (indexed by)

faces of  $\mathcal{T}_n$ : orbifolds  $\bar{\mathbb{R}}_+^d$  for  $d \leq n-2$

tree  $T: x_T \in \mathbb{R}_+^{\text{et}(T)}$  edges  $\hat{=}$  splits  $\leadsto$  strata,  $\mathcal{O}_T \subseteq \bar{\mathcal{O}}_T \Leftrightarrow T_{\text{of } T} = \text{contraction of } T$

$$\bar{O}_T = \bar{O}_{T_1} \cap \bar{O}_{T_2} \Leftrightarrow T \text{ is the biggest common contraction}$$



diffusion tensors: positive semidefinite matrices of size 3  
stratification by subsets of eigenvalues that are equal

ice-cream cone: homeom.  $\cong \mathbb{R}^2$   
 periodically curved Local geometry away from apex  
 $\alpha > 2\pi$ , e.g.  $5\frac{\pi}{2}$   $\rightarrow$  negatively curved: Kale  
 contained in  $\mathbb{T}_q$ , see above

• Tubular neighbourhood Theorem:  $X$  is topologically stratified  
 $\Rightarrow$  a tubular neighbourhood of each stratum  $S$  is a fiber bundle  
 over  $S$  (i.e.  $\cong S \times N_S$  it is locally a cross-product of  $S$  and  $N_S$ )  
 with fiber  $N_S$  where  $N_S$  is homeomorphic to a cone over  
 a stratified space  $L$  with  $\dim L = \dim N_S - 1$ .

$\leadsto$  tangent cone at  $p \in S$ :  $T_p X = T_p S \times N_p X$   
 even more: the tangent bundle is locally trivial  
 tangent manifold normal slice  $\cong N_S$   
 space  $\leftarrow$  a copy of  $\mathbb{R}^{\dim S}$

e.g. open book: on page  $S$ :  $N_p S = \text{one point} = \text{cone over } \emptyset$   
 on spine:  $N_p S = \text{spider} = \text{cone over } \ell \text{ points}$

• Def.  $\dim X = \max_k (\dim M_k)$

• Suppose  $p \in S$  stratum of dim.  $d-1$  where  $d = \dim X$ .

$$N_p X = ? \quad T_p X = ? \quad T_p X = \underbrace{T_p S}_{\dim: d-1} \times \underbrace{N_p S}_1$$

Thm. If  $M_i$  and  $M_j$  are strata,  $M_i \cap M_j \neq \emptyset \Rightarrow M_i \subseteq \bar{M}_j$

Cor. Strata are partially ordered:  $M_i \leq M_j \Leftrightarrow M_i \subseteq \bar{M}_j$

$p \in \bar{S} \Rightarrow \bar{S} \geq S \rightsquigarrow N_p X$  is cone over a finite set (conv. components of intersection of the half-<sup>with</sup>  $d$ -dim strata containing  $S$ )  $\rightarrow$  spider

$\rightarrow T_p X$   <sup>$d$ -dim.</sup> open book  $\rightsquigarrow$  universal for points in  $d$ -dim.  $1$  strata

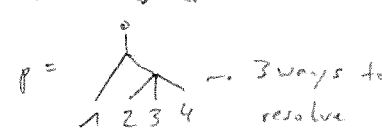
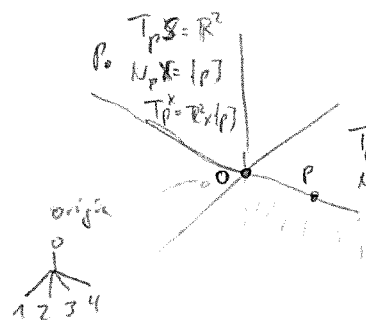
e.g.



$$\left. \begin{matrix} T_p S \\ N_p S \end{matrix} \right\} N_p S = \times$$

therefore at the Whitney cusp the point must be its own stratum

e.g.  $J_4$



$T_0 J_4 = N_0 J_4 =$  cone over Petersen graph



• Metrics, smoothness etc.:

"Riemannian stratified space" when properly defined should require top. strat. <sup>geodesic</sup> metric space st. each stratum is Riemannian

[ "smooth stratified space"  $\rightsquigarrow$  see [Pflaum] (book) ]

e.g. icecream cone / kale: apex has to be its own stratum to get Riemannian strat. space

\* geodesic has to be required (at least intrinsic), else intrinsic metric on  $\Delta$  would be allowed

exercise: embed  $J_3$  into half  $K_2$  for  $\alpha > 3\pi \rightarrow$  3 angles each larger than  $\pi$

- Def. A Whitney stratified space is a top. strat. space  $X$  embedded in a vector space where limits of secant lines joining strata  $M_i$  and  $M_j$  with  $M_i \subseteq \bar{M}_j$  are contained in limits of tangent planes to  $M_j$  as points converge to  $M_i$

e.g. real or complex alg. (semi)analytic (sub)analytic varieties

Thm. Whitney stratification  $\Rightarrow$  triangulation exists

- $X$  strat. metric space

Def. A prob. distr.  $p$  on  $X$  has Fréchet function

$$F(y) = \int_X d(x, y)^2 p(dx)$$

$p$  has Fréchet mean argmin <sub>$y \in X$</sub>   $F(y)$ .

- usual stat. : data wiggles  $\Rightarrow$  mean wiggles  
strat. stat. : mean can stick

e.g.



$$a = b = c = 1 \Rightarrow \text{mean} = 0$$

$$F(r) = a^2 + b^2 + c^2$$

$$F(r_c) = F(r) - 2\epsilon a + 2\epsilon b + 2\epsilon c + 3\epsilon^2$$

$$\sim F(r_c) - F(r) = 2\epsilon(b+c-a) + 3\epsilon^2$$

$> 0$  if  $a < b+c$  for small  $\epsilon$

higher dim. : e.g. open book

- What is a Law of Large Numbers (LLN)?

$X_1, \dots, X_n$  i.i.d. rand var. with values in  $X$  distributed according to  $p$

LLN :  $\bar{X}_n \rightarrow \mu(p)$  as  $n \rightarrow \infty$ .

$\mu(p)$  is sticky  $\Rightarrow$  ?

Fix  $X = \text{open book} = K \times \text{spider}$  where  $K = \mathbb{R}^d$  is spine.  $X$  is CAT(0)

Integrable prob. distr.  $p$  has 3 possibilities for mean  $\mu(p)$ :  $\Rightarrow$  unique Fréchet mean

1)  $\mu(p) \notin K$

2)  $\mu(p) \in K$  and  $\mu(p') \in K$  for all  $p'$  "near"  $p$

3)  $\mu(p) \in K$  but  $\mu(p') \notin K$  for some  $p'$  "arbitrarily near"  $p$

Def. Corresponding to these cases,  $\mu(p)$  is 1) non-sticky, 2) sticky, 3) partly sticky (there is also a perturbation which makes it sticky)

Thm [SAMSI WG] iid  $X_1, X_2, \dots \in X \sim \rho \Rightarrow \bar{X}_n \rightarrow \mu(\rho)$  as  $n \rightarrow \infty$

Moreover, if  $\rho$  is mostly or pretty sticky (1) + (3)  
 is supp  $\rho$  spreads over  $\mathbb{R}^d$   $\Rightarrow \exists$  page  $L$  and rand.  $N \in \mathbb{N}$  s.t.  $\bar{X}_n \in L \forall n \geq N$  a.s.  
 •  $\rho$  sticky  $\Rightarrow \exists$  rand.  $N \in \mathbb{N}$  s.t.  $\bar{X}_n \in K \forall n \geq N$  a.s.

• What is a central limit theorem (CLT)?

$$X_1, X_2, \dots \stackrel{iid}{\sim} \rho, \quad \bar{X}_n \xrightarrow{D} \delta_{\mu(\rho)}$$

idea: rescaled  $\sqrt{n} \bar{X}_n \xrightarrow{D} \tilde{\rho}$  limiting distr. (classical:  $\tilde{\rho}$  gaussian)

Thm [SAMSI WG] A CLT holds for square-int.  $\rho$  holds with limiting distribution

1) gaussian on  $\mathbb{R}^{d+1} = T_{\mu(\rho)} X = T_{\mu(\rho)} L$

2) gaussian on  $\mathbb{R}^d = K = T_{\mu(\rho)} K$

3) half gaussian on  $L$  plus its reflected projection to  $K$  ( $\frac{1}{2}$  gaussian on  $T_{\mu(\rho)} K$ )

• Def. Fix a metric space  $X$  and a top set  $\mathcal{P}$  of int. prob. distr. on  $X$  as well as a subset  $K \subseteq X$ .



The mean of  $\rho(p)$  of  $\rho \in \mathcal{P}$  sticks to  $K$  if every neighbourhood  $U$  of  $p$  in  $\mathcal{P}$  contains a nonempty open set  $U' \subseteq U$  with  $\mu(\rho') \in K \forall \rho' \in U'$ .

• Example (Huebner, Mattingly, Nolan, Miller):

$X =$  isolated plane hyperbolic singularity (Kale with  $\alpha > 2\pi$ )

$\Rightarrow$  CLT has limiting distribution:

$\rightarrow$  induces further stratification (boundary of sector)

