# Flies and regular subdivisions

Michael Joswig

TU Berlin & MPI-MiS Leipzig

14 Dec 2023



©Thomas Endler

joint w/ Holger Eble
Lisa Lamberti
Will Ludington

**❶ Mathematics**
    epistasis
    fitness landscapes
    cluster partitions and dendrograms

**❷ Statistics**
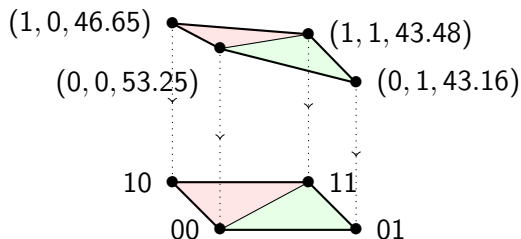    significance test

**❸ Biology**
    *E.coli* evolution
    *Drosophila* microbiome

# Regular subdivision of a point configuration

Epistasis [Bateson 1909]

- point set $V$ in $\mathbb{R}^n$
  - $n$-biallelic genetic system $V = \{0, 1\}^n$

- lift to $\mathbb{R}^{n+1}$ via height function $h : V \to \mathbb{R}$
  - phenoytpe

- take upper convex hull and project back
  - yields subdivision $\mathcal{S}(V, h)$ of conv$(V)$

- generic height function $\rightsquigarrow$ triangulation
  - lifted points coplanar $\iff$ no biological interaction

$(1, 0, 46.65)$  $(1, 1, 43.48)$
$(0, 0, 53.25)$  $(0, 1, 43.16)$

10  11
00  01

# Epistasis and shapes of fitness landscapes
Beerenwinkel, Pachter & Sturmfels 2007

Consider $n$-biallelic system $V = \{0,1\}^n$ with phenotype $h : V \to \mathbb{R}$.

- (relative) population = map $p : V \to \mathbb{R}_{\geq 0}$ with $\sum_{v \in V} p(v) = 1$
- allele frequency vector $\rho(p) := \sum_{v \in V} p(v) v$ contained in $[0,1]^n$
- $\Delta_V :=$ set of all relative populations = simplex of dimension $2^n - 1$
- for fixed $w \in [0,1]^n$:

$$\begin{array}{ll} \text{maximize} & h \cdot p \\ \text{subject to} & p \in \Delta_V \text{ and } \rho(p) = w \end{array} \qquad (\text{LP}(h,w))$$

- if $h$ and $w$ generic then $\text{LP}(h,w)$ has unique optimal solution, the fittest population $p^* = p^*(h,w) = $ vertex of $\{p \in \Delta_V \mid \rho(p) = w\}$
- optimal value of $\text{LP}(h,w)$ is $h \cdot p^* = \sum \lambda_i(h(v_i))$
- piecewise linear function $h^* : [0,1]^n \to \mathbb{R}$, $w \mapsto h \cdot p^*(h,w)$
- regions of linearity of $h^*$ = maximal cells of $\mathcal{S}(V,h)$

# The epistatic weight of a dual edge

Let $V$ be vertex set of some $n$-polytope, equipped with generic height function $h$. Thus $\mathcal{S} = \mathcal{S}(V, h)$ is a triangulation. For

$$s \;=\; \mathrm{conv}\{v_1, v_2, \ldots, v_{n+1}\} \quad \text{and} \quad t \;=\; \mathrm{conv}\{v_2, v_3, \ldots, v_{n+2}\}$$

two adjacent $n$-simplices of $\mathcal{S}$ define

$$E_h(s, t) \;:=\; \begin{pmatrix} 1 & v_{1,1} & v_{1,2} & \ldots & v_{1,n} & h(v_1) \\ 1 & v_{2,1} & v_{2,2} & \ldots & v_{2,n} & h(v_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & v_{n+2,1} & v_{n+2,2} & \ldots & v_{n+2,n} & h(v_{n+2}) \end{pmatrix} \;.$$

The epistatic weight of the dual edge $(s, t)$ is

$$e_h(s, t) \;:=\; |\det E_h(s, t)| \cdot \frac{\mathrm{nvol}(s \cap t)}{\mathrm{nvol}\, s \cdot \mathrm{nvol}\, t} \;.$$

▶ statistics

# Cluster partitions and epistatic filtrations

Consider $\mathcal{S} = \mathcal{S}(V, h)$, with dual graph $\Gamma$.

Picking threshold value $\theta \geq 0$ yields

- $\Gamma(\theta) = \Gamma$ minus dual edges of epistatic weight $> \theta$

### Definition

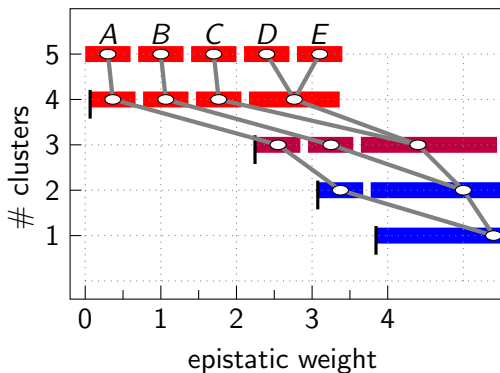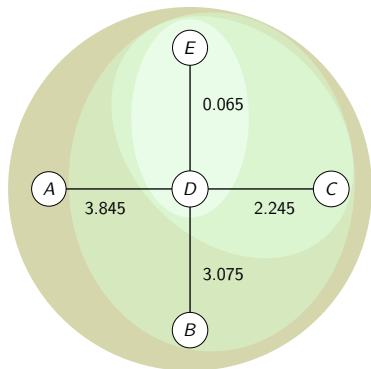A $\theta$-cluster of $\mathcal{S}$ is a connected component of $\Gamma(\theta)$.

- induces partition of $\Gamma(\theta)$ into $\theta$-clusters
- 0-cluster = single facet; $\infty$-cluster = all facets
- epistatic filtration $\Gamma(0) < \Gamma(\theta_1) < \cdots < \Gamma(\theta_\ell) = \Gamma$, linearly ordered by refinement

# Example: nonunimodular triangulation of $[0, 1]^3$

Consider triangulation $\mathcal{S}([0, 1]^3, \texttt{ttd})$ with five maximal simplices:

$A = 000\ 100\ 110\ 101$  $B = 000\ 001\ 101\ 011$  $C = 000\ 010\ 110\ 011$
$D = 000\ 110\ 101\ 011$  $E = 110\ 101\ 011\ 111$

# Height functions as random variables

Fix simplices $s$ and $t$ with joint vertices $v_1, v_2, \ldots, v_{n+2}$ and random variables $X_{v_i}$. We set

$$\lambda_i := (-1)^{n+i} \det(E_i) \cdot \frac{\text{nvol}(s \cap t)}{\text{nvol } s \cdot \text{nvol } t} \ .$$

Then the expectation of the random variable $e_X(s, t)$ satisfies

$$\left| \sum_{i=1}^{n+2} \lambda_i \, \mathbb{E}(X_{v_i}) \right| \ \leq \ \mathbb{E}\big(e_X(s, t)\big) \ \leq \ \sum_{i=1}^{n+2} \big|\lambda_i\big| \, \mathbb{E}(X_{v_i}) \ .$$

If the random variables $X_{v_i}$ are independent, then

- variance can be bounded, too.

If additionally, each random variable is normally distributed, then

- folded normal distribution

# Significance test for one epistatic weight

Let $(s, t)$ be a dual edge of $\mathcal{S}$.

- distribution mean $\mu = \mathbb{E}(e_X(s, t))$ of random variable $e_X(s, t)$ not known exactly

- wanted: one-sided test of significance with null hypothesis $\mu = 0$ vs. alternative $\mu > 0$

Assumption: random variables $X_v$ normally distributed (and independent)

- for sample mean $Z = e_{\bar{X}}(s, t)$ then

$$P(X \geq Z) \;=\; \int_Z^\infty \frac{\sqrt{2}}{\sigma_{e_{\bar{X}}(s,t)}\sqrt{\pi}} \, e^{-\frac{1}{2}\left(\frac{x}{\sigma_{e_{\bar{X}}(s,t)}}\right)^2} dx$$

## Definition

dual edge $(s, t)$ significant if $P(X \geq Z) < 0.05$

# A synthetic experiment

For $V = \{0,1\}^5$, $\eta(v) = 5$ (for $v \neq 0$), $\eta(0) = 5 - \eta_0$, $0.8 \leq \eta_0 \leq 1.2$ the regular subdivision $\mathcal{S}(V, \eta)$ is a vertex split.

- to each vertex we assign normally distributed random variable with mean $\mu = 0$ and standard deviation $0.1 \leq \sigma \leq 2.0$
- 100 realizations per vertex
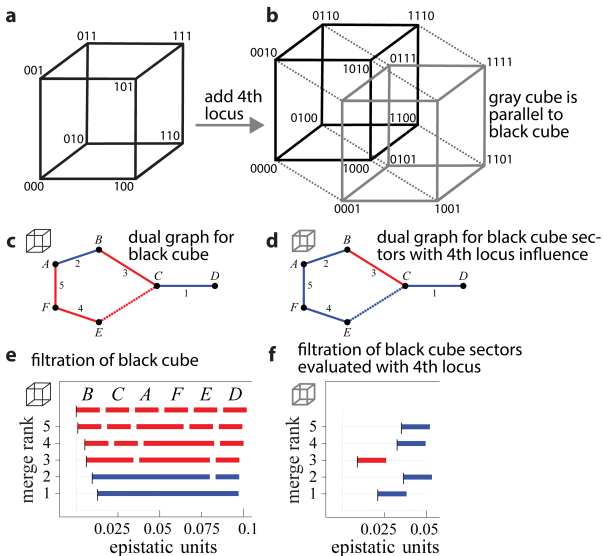- for fixed $(\eta_0, \sigma)$ repeat experiment 100 times; try $p \in \{0.05, 0.1\}$

# *E.coli* evolution. Data set: Khan et al. 2011



- significant 4D interaction: $00001 + 00000|01001|00101|00011 + 00010$
- <u>r</u>ibosome-<u>b</u>inding <u>s</u>ite mutation = master regulator

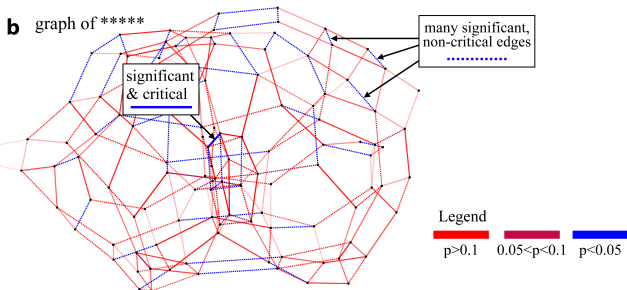# Marginal and conditional epistasis

## parallel epistatic filtration



**a**

011, 111, 001, 101, 010, 110, 000, 100

add 4th locus →

**b**

0110, 1110, 0010, 1010, 0111, 1111, 0100, 1100, 0000, 1000, 0101, 1101, 0001, 1001

gray cube is parallel to black cube

**c** dual graph for black cube

B, A, 2, 3, C, D, 5, F, 4, E, 1

**d** dual graph for black cube sectors with 4th locus influence

B, A, 2, 3, C, D, 5, F, 4, E, 1

**e** filtration of black cube

B C A F E D

merge rank 5 4 3 2 1

0.025  0.05  0.075  0.1

epistatic units

**f** filtration of black cube sectors evaluated with 4th locus

merge rank 5 4 3 2 1

0.025  0.05

epistatic units

# *Drosophila* microbiome. Data set: Ludington lab



**a** Experimental design

**b** graph of *****

many significant, non-critical edges

significant & critical

Legend

$p > 0.1$    $0.05 < p < 0.1$    $p < 0.05$

**c** filtration of *1****

dual edges

*L. plantarum* removed

epistatic units

**d** parallel filtration of *0****

parallel epistatic units

**e** filtration of *1***

*L. brevis* removed

epistatic units

**f** parallel filtration of *0***

parallel epistatic units

- Lactobacilli = master regulators

# Conclusion

- new method to process epistatic data in biology
  - ties in with previous approaches
  - provides a test for statistical signifance
  - agrees with established biological results
- new way to visualize high-dimensional data
  - works for arbitrary regular subdivisions
  - e.g., tropical hypersurfaces (which are dual to regular subdivisions)

📄 Holger Eble, Michael Joswig, Lisa Lamberti, and William B. Ludington, Cluster partitions and fitness landscapes of the *Drosophila* fly microbiome, J. Math. Biol. **79** (2019), no. 3, 861–899.

📄 _____, Master regulators of biological systems in higher dimensions, Proc. Natl. Acad. Sci. USA **120** (2023), no. 51.

# Epistatic Filtrations Calculator

This is an online client for computing higher-order epistatic interactions as detailed in the articles [1] and [2]. It was implemented as polymake extension and can be found online on GitHub. If you found this useful for your scientific work, please cite our paper [1].

The input is a sequence of genotype-phenotype maps, where several phenotypes for the same genotype are considered as independent measurements, thus giving rise to a distribution of phenotypes. Genotypes are 0/1-vectors (i.e., here we are treating the biallelic case only), and phenotypes are real numbers. The entire dataset is supposed to be contained in a single file of type csv (ASCII text, comma separated values). Such files can be exported from standard spreadsheet software.

## Upload csv file

The input csv file must be in the precise format shown in the exemplary screenshot on the right hand side:

- The genotypes are placed in the first data row. Their coordinates are separated by vertical bars, e.g. 0|0|1|0.
- Right below the genotypes, the measured data is placed accordingly. The columns are allowed to be of varying size.

Browse... No file selected.     **UPLOAD**

Please upload a file.