

A toy model of ensemble induced epistasis

Brandon H. Schlomann

August 20, 2018

High-order epistasis is quite an enigmatic phenomenon. While most commonly discussed in the context of mutations in proteins and macromolecules, it is in fact a general phenomenon and arises in other systems, most famously community ecology, where it is framed as “high-order interactions”, or, “non-linear interactions”. In both contexts, people argue hotly over its existence, origins, consequences, and how best to measure it.

Following Zach Sailer’s work on ensemble induced epistasis in a lattice protein model, I develop here a toy ensemble model, abstracted from the context of molecules, to study the emergence of high-order epistasis as a purely statistical property of ensemble systems. Mathematically, epistasis arises due to the normalization of probabilities, which introduces non-additive changes to ensemble-averaged observables. This is analogous to the central role that partition functions play in thermal models of ensemble induced allostery (Cooper + Dryden 1984).

I also explore what ensemble induced epistasis looks like in the limit of many states, analogous to recent efforts by physicists working in community ecology to draw parallels with disordered ordered systems and use tools from mean field theory. I show that for a general class of random perturbations (mutations), pairwise epistasis limits to a Gaussian distributed variable and calculate its mean and variance.

Model foundations

Consider a system of S states specified with probabilities p_i , $i = 1, 2, 3, \dots, S$. Let \mathcal{O} be a generic observable that is an arbitrary function of the states with values \mathcal{O}_i at each state i . We will refer to this observable as a phenotype. The phenotype measured in an ensemble system will be the average phenotype over the ensemble,

$$\langle \mathcal{O} \rangle = \sum_j \mathcal{O}_j p_j \tag{1}$$

Consider now a set of perturbations (mutations) to the system, each identified by a label

$n = 1, 2, 3, \dots, N$, that can generally alter both the probability and the observable of each state. We will take the perturbations to be additive so that an additive model of combined perturbations (the most commonly used non-interacting model) is exact when there is only one state.

It will be convenient to define the perturbation to probabilities as a properly normalized distribution of its own, denoted by $\alpha_{i,n}$, where i refers to the state and n labels the perturbation. Other than being non-negative and properly normalized, the $\alpha_{i,n}$ s are arbitrary.

Let the original probabilities be denoted by $p_{i,0}$. Following perturbation number n , the new probability of the i^{th} state, denoted by $p_{i,n}$, is now just the average of the original and perturbing probabilities:

$$p_{i,n} = \frac{p_{i,0} + \alpha_{i,n}}{\sum_j (p_{j,0} + \alpha_{j,n})} = \frac{p_{i,0} + \alpha_{i,n}}{2}. \quad (2)$$

For the phenotypes, let the original phenotypes be denoted $\mathcal{O}_{i,0}$. These will be perturbed by terms denoted by $\beta_{i,n}$, which don't have to be normalized. Following perturbation number n , the new phenotypes are

$$\mathcal{O}_{i,1} = \mathcal{O}_{i,0} + \beta_{i,1}. \quad (3)$$

When multiple perturbations act, say perturbations 1, 2, and 3, we will use the notation $p_{i,123}$ and $\mathcal{O}_{i,123}$.

While extremely simple, this model can be used to study when and how epistasis emerges in a general context. When the α_i 's and β_i 's are all independent, any epistasis that arises is due purely to the ensemble.

Calculations

1. Epistasis arises even when perturbations only alter the probabilities of each state and not their phenotypes

Consider a set of perturbations that don't alter the phenotype of each state, so all $\beta_i = 0$, but do alter the probabilities of each state. The measured phenotype after each single perturbation is

$$\langle \mathcal{O} \rangle_n = \sum_j \mathcal{O}_j p_{j,n} = 2^{-1} \sum_j \mathcal{O}_j (p_{j,0} + \alpha_{j,n}) = 2^{-1} \left(\langle \mathcal{O} \rangle_0 + \sum_j \mathcal{O}_j \alpha_{j,n} \right) \equiv 2^{-1} (\langle \mathcal{O} \rangle_0 + \langle \mathcal{O} \rangle_n^\alpha). \quad (4)$$

The ensemble averaged phenotype following perturbation n is an average of averages: the average of the phenotype averaged over the original probability measure and over the perturbing measure. We have introduced the notation

$$\langle \mathcal{O} \rangle_n^\alpha \equiv \sum_j \mathcal{O}_j \alpha_{j,n}. \quad (5)$$

Let's now consider the effect of double perturbations. Specifically, let's consider the effect of perturbation m following perturbation n on the original system:

$$\begin{aligned} \langle \mathcal{O} \rangle_{nm} &= \sum_j \mathcal{O}_j p_{j,nm} \\ &= 2^{-1} \sum_j \mathcal{O}_j (p_{j,n} + \alpha_{j,m}) \\ &= 2^{-1} (\langle \mathcal{O} \rangle_n + \langle \mathcal{O} \rangle_m^\alpha) \\ &= 2^{-1} (2^{-1} [\langle \mathcal{O} \rangle_0 + \langle \mathcal{O} \rangle_n^\alpha] + \langle \mathcal{O} \rangle_m^\alpha) \\ &= 4^{-1} \langle \mathcal{O} \rangle_0 + 4^{-1} \langle \mathcal{O} \rangle_n^\alpha + 2^{-1} \langle \mathcal{O} \rangle_m^\alpha. \end{aligned} \quad (6)$$

A straightforward calculation shows that if the phenotypes are identical for all the states, the perturbations have no effect, as expected, since there is no variation to average over.

To see if pairwise epistasis appears in this system, we compare this result to the prediction of an additive model. When there is only one state and the phenotypes are perturbed directly via β terms, the additive model is

$$\begin{aligned} \mathcal{O}_{nm}^{\text{add}} &= \mathcal{O}_0 + \beta_n + \beta_m \\ &= \mathcal{O}_0 + (\mathcal{O}_0 + \beta_n) + (\mathcal{O}_0 + \beta_m) - 2\mathcal{O}_0 \\ &= \mathcal{O}_n + \mathcal{O}_m - \mathcal{O}_0. \end{aligned} \quad (7)$$

By analogy, the additive model in a multi-state system is

$$\begin{aligned}
\langle \mathcal{O} \rangle_{nm}^{\text{add}} &= \langle \mathcal{O} \rangle_n + \langle \mathcal{O} \rangle_m - \langle \mathcal{O} \rangle_0 \\
&= 2^{-1}(\langle \mathcal{O} \rangle_0 + \langle \mathcal{O} \rangle_n^\alpha) + 2^{-1}(\langle \mathcal{O} \rangle_0 + \langle \mathcal{O} \rangle_m^\alpha) - \langle \mathcal{O} \rangle_0 \\
&= 2^{-1}(\langle \mathcal{O} \rangle_n^\alpha + \langle \mathcal{O} \rangle_m^\alpha).
\end{aligned} \tag{8}$$

The difference between measured and predicted observables is

$$\langle \mathcal{O} \rangle_{nm} - \langle \mathcal{O} \rangle_{nm}^{\text{add}} = 4^{-1} (\langle \mathcal{O} \rangle_0 - \langle \mathcal{O} \rangle_n^\alpha). \tag{9}$$

Epistasis arises anytime the probabilities of each state change, and its magnitude depends on the difference between the average phenotype over the original and perturbing probability distributions. This derivation was completely agnostic to the nature of the perturbations and the phenotype. Even if the perturbations are all taken to be independent (say, independent random variables), epistasis still arises. The additive model will not predict correctly, despite the the fact that the perturbations to each state’s phenotype are purely additive.

TO DO: more mutations + high-order epistasis, effect of perturbations to the phenotypes.

2. The many state limit and the Central Limit Theorem

Much of the recent surge in physics-inspired community ecology is motivated by the idea that in the limit of large ecosystems (many species), universal behavior emerges, akin to universality in non-living systems. Mathematically, this is studied using a disordered systems approach: model parameters are drawn from probability distributions to represent “typical” systems, calculations are performed, and then the results are averaged over random draws of the parameters, leading to general conclusions. In the limit of many species (analogous to thermodynamic limit), community observables (often involving sums over species) take on universal properties via the Central Limit Theorem and become simpler to evaluate, since they involve only Gaussian variables. I’m beginning to explore this approach in the context of this toy ensemble model—what does ensemble induced epistasis look like in the limit of many states? In the following, I show that pairwise epistasis limits to a Gaussian variable and compute its mean and variance.

Let’s start with a set of probabilities and phenotypes, $\{p_{i,0}, \mathcal{O}_i\}$, $i = 1, 2, 3, \dots, S$, taken as given. Let’s draw all the $\alpha_{i,n}$ parameters from the same well-behaved (no fat tail) probability distribution. To enforce normalization, we’ll construct the α ’s via

$$\alpha_{i,n} \equiv \frac{\tilde{\alpha}_{i,n}}{\sum_j \tilde{\alpha}_{j,n}} \tag{10}$$

where the $\tilde{\alpha}$ s come from a probability distribution with a finite mean $\tilde{\mu}_\alpha/S$ and variance $\tilde{\sigma}_\alpha^2/S$. The $1/S$ scaling is required to ensure a well defined limit, analogous to the thermodynamic limit in physics.

Now let's consider what happens to the average phenotype over the perturbing distribution, $\langle \mathcal{O} \rangle_n^\alpha$, as the number of states becomes large. In terms of the $\tilde{\alpha}$ s,

$$\begin{aligned} \langle \mathcal{O} \rangle_n^\alpha &= \sum_j \mathcal{O}_j \alpha_{j,n} \\ &= \frac{1}{\sum_k \tilde{\alpha}_{k,n}} \sum_j \mathcal{O}_j \tilde{\alpha}_{j,n} \end{aligned} \quad (11)$$

As $S \rightarrow \infty$, the normalization term can be replaced by $\tilde{\mu}_\alpha$,

$$\langle \mathcal{O} \rangle_n^\alpha \rightarrow \frac{1}{\tilde{\mu}_\alpha} \sum_j \mathcal{O}_j \tilde{\alpha}_{j,n}. \quad (12)$$

Since the sum is over S terms, by the Central Limit Theorem, $\langle \mathcal{O} \rangle_n^\alpha$ limits to a Gaussian variable regardless of the nature of the α 's. This means that the magnitude of pairwise epistasis,

$$\langle \mathcal{O} \rangle_{nm} - \langle \mathcal{O} \rangle_{nm}^{\text{add}} = 4^{-1} (\langle \mathcal{O} \rangle_0 - \langle \mathcal{O} \rangle_n^\alpha), \quad (13)$$

limits to a Gaussian distributed variable, regardless of the nature of the α s.

Defining an average over realizations of the α parameters as $\mathbb{E}[\dots]_\alpha$, the mean of $\langle \mathcal{O} \rangle_n^\alpha$ is

$$\begin{aligned} \mathbb{E}[\langle \mathcal{O} \rangle_n^\alpha]_\alpha &= \frac{1}{\tilde{\mu}_\alpha} \sum_j \mathcal{O}_j \mathbb{E}[\tilde{\alpha}_{j,n}]_\alpha \\ &= \frac{1}{\tilde{\mu}_\alpha} \frac{\tilde{\mu}_\alpha}{S} \sum_j \mathcal{O}_j \\ &\equiv \bar{\mathcal{O}}, \end{aligned} \quad (14)$$

or the empirical mean phenotype of the S states. In other words, the average phenotype over a uniform distribution. Therefore, the mean pairwise epistasis is

$$\lim_{S \rightarrow \infty} \mathbb{E} [\langle \mathcal{O} \rangle_{nm} - \langle \mathcal{O} \rangle_{nm}^{\text{add}}]_\alpha = 4^{-1} (\langle \mathcal{O} \rangle_0 - \bar{\mathcal{O}}). \quad (15)$$

The limiting mean is truly universal—it doesn’t depend on the perturbations at all, only on what we took as given, namely the original probabilities and observables, $\{p_{i,0}, \mathcal{O}_i\}$. The fact that the mean of the $\tilde{\alpha}$ s doesn’t even appear here reflects the fact that the average perturbation gets normalized out, so is unphysical.

The variance of the limiting distribution depends on the variance of the underlying α s:

$$\begin{aligned}
\lim_{S \rightarrow \infty} \text{Var} [\langle \mathcal{O} \rangle_{nm} - \langle \mathcal{O} \rangle_{nm}^{\text{add}}]_{\alpha} &= \text{Var} [4^{-1} (\langle \mathcal{O} \rangle_0 - \langle \mathcal{O} \rangle_n^{\alpha})]_{\alpha} \\
&= \frac{1}{16} \text{Var} [\langle \mathcal{O} \rangle_n^{\alpha}]_{\alpha} \\
&= \frac{\sigma_{\mathcal{O}}^2 \sigma_{\alpha}^2}{16 \mu_{\alpha}^2}.
\end{aligned} \tag{16}$$

Here $\sigma_{\mathcal{O}}^2$ is the empirical variance of the phenotypes. It’s encouraging that this appears: if all the phenotypes are the same ($\sigma_{\mathcal{O}}^2 = 0$), there should be no effect of perturbing the probabilities, and we see the variance of epistasis going to zero, along with the mean. The ratio $\sigma_{\alpha}^2 / \mu_{\alpha}^2$ should be thought of as the strength of the perturbation, and so naturally sets the scale of pairwise epistasis.