

1 Analysis of discrete linear-nonlinear codes

1.1 Code definition

We consider the following family of code,

$$r = f_L(x) + f_{NL}(x) + \epsilon$$

where the N -dimensional neural response r to stimulus $x \in [1, \dots, n]^K$ is modeled in terms of a linear part, $f_L(x)$, a nonlinear part, $f_{NL}(x)$, and additive Gaussian noise, $\epsilon \sim \mathcal{N}(0, \sigma^2 I_N)$. The linear part has the form,

$$f_L(x) = Mx_z$$

where M is an $N \times K$ matrix and x_z is z-scored version of x . The magnitude of vectors in M are chosen so that $\mathbb{E}_x f_L(x)^2 = P_L$. The nonlinear part has the form,

$$f_{NL}(x)_i = \sqrt{P_{NL}} n_i [x_1 = t_i(1)] \dots [x_K = t_i(K)]$$

for $i \in [1, \dots, n^K]$ where $[...]$ is the indicator function, t_i returns a values in $[1, \dots, K]$, and n_i is a unit vector chosen so that $n_i \cdot n_j = 0 \forall i \neq j$. That is, there is a dimension in f_{NL} that is one for a single, unique stimulus and zero for all other stimuli.

A set of stimuli is defined by two parameters: K , the number of distinct features, or latent variables; n , the number of a discrete values that each of the K features takes on. There are n^K unique stimuli. A code is described by four parameters: P_L and P_{NL} are the power of the linear and nonlinear code parts, respectively – the total power is $P = P_L + P_{NL}$; the number of neurons N ; and the noise variance σ^2 . The SNR = $\sqrt{P/\sigma^2}$.

Current results:

1. Analytic approximation of the error rate.
2. Analytic worst-case lower bound on XOR performance.
3. Analytic approximation of the CCGP.

Current applications:

1. Can show the SNR required for both high flexibility and CCGP for a given stimulus set.
2. Can simulate responses to multiple simultaneous stimuli and get error rate.

Possible extensions/applications:

1. Can show how many examples are necessary to learn a generalizable category boundary (a la Ben’s paper).
2. Can make claims about modularity for some kinds of tasks.
3. Can use to model transfer in perceptual learning.
4. Can characterize tradeoff between CCGP and multi-stimulus capacity (high multi-stimulus capacity may imply low CCGP).
5. Additional code parts (e.g., f_{I1}).
6. Continuous stimuli.

1.2 Preliminaries

To estimate the error rate and other quantities in these codes, we need to understand how the code parameters scale the distances in representation space between nearby stimuli. This distance will be important for deriving error and generalization rates. We do this separately for the different code parts and show that the distances are close to additive later.

1.2.1 The linear code part

In the linear code, all stimuli lie on a K -dimensional rectangular lattice. We derive the distance between points on that lattice as a function of our stimulus and code parameters, as well as the number of neighbors a representation has at minimum and next-from-minimum distance.

Linear distance derivation: First, we find how the distance between adjacent stimuli in the linear code depends on the number of features K and the number of values that each feature takes on n along with the linear power of

the code (P_L),

$$d_L = \sqrt{\frac{12P_L}{K(n^2 - 1)}}$$

We approach this by computing the variance (i.e., the linear power P_L) of a uniformly sampled K -dimensional lattice with n points spaced at distance d_L along each dimension. Then, we invert the expression for the variance to find an expression for the distance between the points. First, we write the variance P_L as

$$\begin{aligned} n^K P_L &= \sum_{i=0}^{n-1} \left[\left(i - \frac{n-1}{2} \right)^2 d_L^2 + \sum_{j=0}^{n-1} \left[\left(j - \frac{n-1}{2} \right)^2 d_L^2 + \dots \right] \right] \\ &= \sum_i^{n-1} \sum_j^{n-1} \dots \sum_k^{n-1} \left(i - \frac{n-1}{2} \right)^2 d_L^2 + \left(j - \frac{n-1}{2} \right)^2 d_L^2 + \dots + \left(k - \frac{n-1}{2} \right)^2 d_L^2 \\ &= K n^{K-1} \sum_i^{n-1} \left(i - \frac{n-1}{2} \right)^2 d_L^2 \\ &= K n^{K-1} d_L^2 \sum_i^{n-1} i^2 - (n-1) \sum_i^{n-1} i + n \frac{n-1}{2} \end{aligned}$$

and we can rewrite this with known expressions for the sum of integers and sum of squared integers up to a particular value,

$$\begin{aligned} n^K P_L &= K n^{K-1} d_L^2 \left[\frac{(n-1)n(2n-1)}{6} - \frac{n(n-1)^2}{2} + \frac{n(n-1)^2}{4} \right] \\ &= K n^K d_L^2 \left[\frac{(n-1)(2n-1)}{6} - \frac{(n-1)^2}{4} \right] \\ &= K n^K d_L^2 \left[\frac{2n^2 - 3n + 1}{6} - \frac{n^2 - 2n + 1}{4} \right] \\ &= K n^K d_L^2 \left[\frac{4n^2 - 6n + 2}{12} - \frac{3n^2 - 6n + 3}{12} \right] \\ n^K P_L &= K n^K d_L^2 \frac{n^2 - 1}{12} \\ P_L &= K d_L^2 \frac{n^2 - 1}{12} \end{aligned}$$

Now, we rewrite in terms of d_L ,

$$d_L = \sqrt{\frac{12P_L}{K(n^2 - 1)}}$$

which is the expression given above.

Following from the lattice structure, stimuli at a diagonal point on the lattice have distance $\sqrt{2}d_L$.

Linear neighbors derivation: Second, we find the average number of neighbors that a particular stimulus has at both this nearest distance N_{LA} and nearest diagonal distance N_{LD} . This is a counting problem. We observe that, in the lattice, there are two edge values for each feature and $n - 2$ non-edge values. Thus,

$$N_{LA} = \frac{1}{n^K} \sum_{c=0}^K (2K - c) \binom{K}{c} (n - 2)^{K-c} 2^c$$

and

$$N_{LD} = \frac{1}{n^K} \sum_{c=0}^K \left(4 \binom{K-c}{2} + 2(K-c)c + \binom{C}{2} \right) \binom{K}{c} (n - 2)^{K-c} 2^c$$

Nonlinear distance derivation: The nonlinear distance has been treated in detail elsewhere (for any f_{Ic} with $c \in [1, \dots, K]$). The nonlinear distance is

$$d_N = \sqrt{2P_N}$$

Nonlinear neighbors derivation: Because each nonlinear representation is along a vector that is orthogonal to all other nonlinear representations, from a particular stimulus all other representations are at minimum distance. So,

$$N_N L = n^K - 1$$

Total code distance: Naively, the total code distance would be

$$d_C = \sqrt{d_L^2 + d_{NL}^2}$$

However, because the linear and nonlinear parts are chosen to project into random subspaces, there is the chance of some alignment. Thus, the corrected total code distance is a random variable with the following form,

$$d_C = \sqrt{d_L^2 + d_{NL}^2 + 2d_{NL}d_L\eta}$$

where $\eta \sim \mathcal{N}(0, 1/N)$ due to the fact that the dot product of two unit vectors are normally distributed with variance inverse to their length (i.e., the distribution of η). The distance is similarly defined for the next nearest stimuli,

$$d_{C+1} = \sqrt{2d_L^2 + d_{NL}^2 + \sqrt{8}d_{NL}d_L\eta}$$

Total code neighbors: To combine the code neighbors, it is enough to simply take the minimum between the linear and nonlinear parts, which will always be equal to N_{LA} (or N_{LD} for next-nearest). So,

$$\begin{aligned} N_C &= N_{LA} \\ N_{C+1} &= N_{LD} \end{aligned}$$

Nonlinear distance as a function of linear distance: It is of interest to compute how, for fixed power, linear and nonlinear distance depend on each other. So,

$$\begin{aligned} P &= P_N + P_L \\ &= \frac{d_N^2}{2} + Kd_L^2 \frac{n^2 - 1}{12} \\ d_N^2 &= 2P - 2Kd_L^2 \frac{n^2 - 1}{12} \\ d_N &= \sqrt{2P - 2Kd_L^2 \frac{n^2 - 1}{12}} \end{aligned}$$

1.3 Approximating code error rate

First, we ask how often this family of codes will make errors. Here, we define an error as the most likely stimulus under a maximum likelihood decoder \hat{x} not being the original stimulus x . Following the logic that the nearest stimuli

are most likely to be errored toward, we develop the following expression for the error rate,

$$P(\text{error}) = N_C Q\left(-\frac{d_C}{2\sigma}\right) + N_{C+1} Q\left(-\frac{d_{C+1}}{2\sigma}\right)$$

Thus, we can see that the error rate depends most strongly on the distances. From our distance definitions, we know that increasing nonlinear power is the most efficient way to increase distance. As a consequence, to drive the error rate down, it is best to put all code power toward the nonlinear part.

For multiple stimuli, we hypothesize that we can write the code error rate as,

$$P(\text{error}) = SN_C Q\left(-\frac{d_C}{2\sigma}\right) + SN_{C+1} Q\left(-\frac{d_{C+1}}{2\sigma}\right) + N_S Q\left(-\frac{d_S}{2\sigma}\right)$$

where N_S is the number of equiprobable stimulus sets under the linear part of the code and $d_S = 2\sqrt{P_N}$ (maybe).

The number of swaps

$$N_S = \frac{1}{2} \binom{S}{2} \sum_{i=0}^{K-1} \binom{K}{i} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{K-i} \sum_{j=1}^{K-i} \binom{K-i}{i}$$

1.4 Approximating code flexibility

We consider a code for a discrete set of stimuli to be flexible if arbitrary binary partitions of that set can be read out with a linear decoder. When a code has a mixture of linear and nonlinear stimulus representations, some partitions are orthogonal to the linear structure in the representation and can be implemented only if the nonlinear components of the representation are strong enough – one such partition is the parity or XOR partition. Thus, to approximate code flexibility, we will focus on this case. It allows us to ignore any contribution to the representation from the linear code and focus only on the nonlinear code. Further, we will take a lower bound on this case.

In the nonlinear code for $n^K = N_s$ stimuli, all of the stimuli are $\sqrt{P_N}$ from the origin in representation space and $\sqrt{2P_N}$ from each other. In this case,

the vector corresponding to the optimal hyperplane for a linear decoder that implements an arbitrary partition of such stimuli has a constant magnitude c in the direction of all stimuli – and the magnitude is positive for stimuli in one category and negative for stimuli in the other category. Using this understanding, we can calculate the performance of the linear decoder where r is the decoding vector, x is particular stimulus representation in the positive category, and σ^2 is the variance of normally distributed output noise for the neurons in the code:

$$\begin{aligned}
E_f &\geq P(rx > 0) \\
&= P(\mathcal{N}(\sqrt{P_N}, N_s \sigma^2) > 0) \\
&= Q\left(-\frac{\sqrt{P_N}}{\sqrt{N_s \sigma^2}}\right) \\
&= Q\left(-\frac{\sqrt{P_N}}{n^{K/2} \sigma}\right)
\end{aligned}$$

where Q is the cumulative distribution function of the standard normal distribution.

1.5 Approximating code generalization

We consider a code to have good generalization performance when a linear decoder aligned with some combination of code features that is learned on one part of the stimulus space provides good performance on another part of the stimulus space. In a simple case with two stimulus features that each take on two values, this means that a linear decoder that discriminates the value of the second feature (0 or 1) learned for a fixed value of the first feature (say, 0) will generalize with minimal loss of performance to other values of the first feature (in this case, 1). This notion of generalization performance is referred to as cross-condition generalization performance (CCGP).

We set out to approximate CCGP with a pair of stimuli used for training and a third stimulus that must be generalized to. Here, we consider a linear code that is distorted by a nonlinear code. Thus, we can consider distances in the purely linear code and distances in the purely nonlinear code separately.

We use d_{LL} to denote the distance in the linear code between the two stimulus representations used to learn the classification. This distance is along the unit

vector f_1 . Further, we use d_{LG} to denote the distance along that unit vector f_1 of the third stimulus s_3 which is generalized to. We use d_{LA} to denote the distance between the pair of stimuli that is used for training and the third stimulus along the axis that they are to be generalized over, which we denote as the unit vector f_2 . Each of the stimuli also undergoes a distortion of magnitude $\sqrt{P_N}$ due to the nonlinear code, we denote the direction of these distortions as the unit vectors n_i . They are chosen such that $n_i \cdot n_j = 0$ for $i \neq j$ – however, $n_i \cdot f_j$ is not constrained to be zero. From above, we know that $n_i \cdot f_j \sim \mathcal{N}(0, 1/D)$ where D is the full dimensionality of the space (i.e., the number of neurons in the code). Additionally, for convenience, we also use n_{ij} for any number of indices to refer to the following

$$n_{ij} = \frac{n_i + n_j}{\sqrt{2}}$$

and similarly for more indices, so that the end vector is a unit vector.

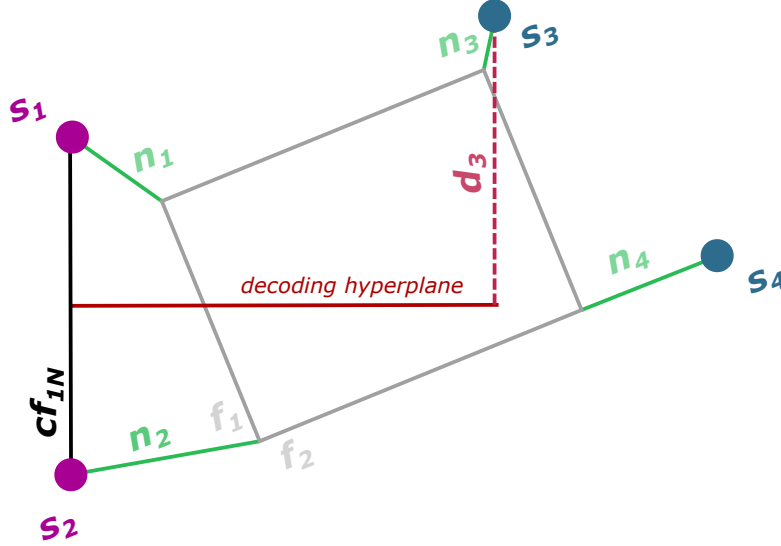


Figure 1: Illustration of the idea behind the CCGP approximation. The purple points are the representations used for training, the blue are those to be generalized to. The representations used for training define a hyperplane (red) that depends on both the linear (f_1) and nonlinear (n_1 and n_2) parts of the code. To approximate CCGP, we can simply take the dot product between the vector defined by this hyperplane and the positions of the other representations (blue). Thus, we find distance $d_3 = f_{1N} \cdot s_3 - \frac{c}{2}$.

First, we find the center points between our two pairs of stimuli in the full code with reference to the “bottom left” stimulus, s_2 . In particular, $s_1 =$

$f_1 d_{LL} + d_N n_{12}$, $s_2 = 0$, and $s_3 = d_{LG} f_1 + d_{LA} f_2 + d_N n_{23}$. Thus,

$$\hat{s}_{12} = \frac{1}{2} (d_{LL} f_1 + d_N n_{12})$$

Next, we find the vector pointing between the two representations used for learning (that is, s_1 and s_2), which is given by

$$\begin{aligned} f_{1N} &= \frac{1}{c} (s_1 - s_2) \\ f_{1N} &= \frac{1}{c} (d_{LL} f_1 + d_N n_{12}) \end{aligned}$$

where c is a random variable that normalizes f_{1N} to be a unit vector, and corresponds to the distance between s_1 and s_2 ,

$$\begin{aligned} c &= \sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_N f_1 \cdot n_{12}} \\ &= \sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_N \mathcal{N}(0, 1/D)} \end{aligned}$$

Now, using f_{1N} , c , along with our understanding of s_3 as a linear combination of linear and nonlinear codes, we can directly approximate CCGP. To do this, we need to find the position of s_3 along the decoding vector defined by f_{1N} , and then we evaluate whether that magnitude is greater or smaller than the threshold $c/2$. So, to find this distance relative to the threshold, we need d_3 such that

$$\begin{aligned} d_3 &= f_{1N} \cdot s_3 - \frac{c}{2} \\ &= \frac{1}{c} (d_{LL} f_1 + d_N n_{12}) (d_{LG} f_1 + d_{LA} f_2 + d_N n_{23}) - \frac{c}{2} \end{aligned}$$

First, we focus on the first term and drop c for now,

$$\begin{aligned} t_1 &= (d_{LL} f_1 + d_N n_{12}) (d_{LG} f_1 + d_{LA} f_2 + d_N n_{23}) \\ &= d_{LL} d_{LG} + d_{LL} d_N f_1 n_{23} + d_{LG} d_N f_1 n_{12} + d_{LA} d_N f_2 n_{12} + d_N^2 n_{23} n_{12} \\ &= d_{LL} d_{LG} + \frac{1}{2} d_N^2 + \frac{d_N}{\sqrt{2}} (d_{LG} f_1 n_1 + d_{LG} f_1 n_2 + d_{LL} f_1 n_2 + d_{LL} f_1 n_3) + d_{LA} d_N n_{12} f_2 \end{aligned}$$

Next, we bring back the full expression and multiply everything by c ,

$$\begin{aligned}
cd_3 &= d_{LL}d_{LG} + \frac{1}{2}d_N^2 + \frac{d_N}{\sqrt{2}}(d_{LG}f_1n_1 + d_{LG}f_1n_2 + d_{LL}f_1n_2 + d_{LL}f_1n_3) + d_{LA}d_Nn_{12}f_2 - \frac{c^2}{2} \\
&= d_{LL}d_{LG} + \frac{1}{2}d_N^2 + \frac{d_N}{\sqrt{2}}(d_{LG}f_1n_1 + d_{LG}f_1n_2 + d_{LL}f_1n_2 + d_{LL}f_1n_3) + d_{LA}d_Nn_{12}f_2 \\
&\quad - \frac{1}{2}d_{LL}^2 - \frac{1}{2}d_N^2 - \frac{d_{LL}d_N}{\sqrt{2}}(f_1n_1 + f_1n_2) \\
&= d_{LG}d_{LL} - \frac{1}{2}d_{LL}^2 + (d_{LG} - d_{LL})d_Nf_1n_{12} + d_{LL}d_Nf_1n_{23} + d_{LA}d_Nf_2n_{12} \\
d_3 &= \frac{d_{LG}d_{LL} - \frac{1}{2}d_{LL}^2 + (d_{LG} - d_{LL})d_Nf_1n_{12} + d_{LL}d_Nf_1n_{23} + d_{LA}d_Nf_2n_{12}}{\sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_Nf_1n_{12}}} \\
d_3 &\sim \frac{d_{LG}d_{LL} - \frac{1}{2}d_{LL}^2 + (d_{LG} - d_{LL})d_N\mathcal{N}(0, 1/D) + d_{LL}d_N\mathcal{N}(0, 1/D) + d_{LA}d_N\mathcal{N}(0, 1/D)}{\sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_N\mathcal{N}(0, 1/D)}}
\end{aligned}$$

and this d_3 relates to the CCGP error rate in the following way,

$$P(\text{CCGP error}_1) = \mathbb{E}_{d_3} Q(-d_3/\sigma)$$

Thus, making d_3 as large as possible will minimize CCGP errors.

For D much larger than any of d_{LG}^2 , d_{LA}^2 , d_{LL}^2 , we can further simplify by ignoring the random variables (which is also the case where the nonlinear and linear codes are perfectly orthogonal to each other),

$$P(\text{CCGP error}) \approx Q\left(-\frac{1}{\sigma} \frac{d_{LG}d_{LL} - \frac{1}{2}d_{LL}^2}{\sqrt{d_{LL}^2 + d_N^2}}\right)$$

Next, we can rewrite this in terms the number of features K , the number of values n taken on by each feature, the steps on the lattice separating the trained pair from each other i (i.e., $d_{LL} = id_L$) and the generalized stimulus from the reference stimulus j (i.e., $d_{LG} = jd_L$), as well as the linear P_L and nonlinear P_N power. So,

$$P(\text{CCGP error}) \approx Q\left(-\frac{12P_L}{K(n^2 - 1)\sigma} \frac{ij - \frac{1}{2}i^2}{\sqrt{i^2 \frac{12P_L}{K(n^2 - 1)} + 2P_N}}\right)$$