

# 1 The tradeoff between flexibility and generalization for discrete codes

## 3 1.1 Approximating code flexibility

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

13 In the nonlinear code for  $n^K = N_s$  stimuli, all of the stimuli are  $\sqrt{P_N}$  from  
14 the origin in representation space and  $\sqrt{2P_N}$  from each other. In this case,  
15 the vector corresponding to the optimal hyperplane for a linear decoder that  
16 implements an arbitrary partition of such stimuli has a constant magnitude  
17  $c$  in the direction of all stimuli – and the magnitude is positive for stimuli  
18 in one category and negative for stimuli in the other category. Using this  
19 understanding, we can calculate the performance of the linear decoder where  
20  $r$  is the decoding vector,  $x$  is particular stimulus representation in the positive  
21 category, and  $\sigma^2$  is the variance of normally distributed output noise for the  
22 neurons in the code:

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

23 where  $Q$  is the cumulative distribution function of the standard normal dis-  
24 tribution.

## 25 1.2 Approximating code generalization

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

35 We set out to approximate CCGP for two sets of two stimulus representations  
36 each. Here, we consider a linear code that is distorted by a nonlinear code.  
37 Thus, we can consider distances in the purely linear code and distances in  
38 the purely nonlinear code separately.

### 39 1.2.1 Preliminaries

40 First, we find how the distance between adjacent stimuli in the linear code  
41 depends on the number of features  $K$  and the number of values that each  
42 feature takes on  $n$  along with the linear power of the code ( $P_L$ ),

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

43 **Linear distance derivation:** We approach this by computing the variance  
44 (i.e., the linear power  $P_L$ ) of a  $K$ -dimensional lattice with  $n$  points spaced  
45 with distance  $d_L$  along each dimension when the points are sampled with  
46 uniform probability. Then, we invert the expression for the variance to find  
47 an expression for the distance between the points. First, we write the variance

48  $P_L$  as

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

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

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

51 Now, we rewrite in terms of  $d_L$ ,

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

52 which is the expression given above.

53 This distance can be used as a scaling factor that allows translation between  
54 distance in stimulus space and distance in representation space – that is, if

two stimuli have distance  $l$  in stimulus space, then they have distance  $ld_L$  in representation space. However, this assumes that each feature is encoded with the same fidelity, which may not always be true.

Second, we find that the nonlinear distance is,

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

Further, we can also observe that each representation in the linear code undergoes a distortion of magnitude  $P_N$  in a random direction.

Third, we remind ourselves that the dot product of two random unit vectors  $u_1 \cdot u_2$  in a  $D$ -dimensional space follows the distribution

$$u_1 \cdot u_2 \sim \mathcal{N}(0, 1/D)$$

for large  $D$ .

### 1.2.2 Main derivation

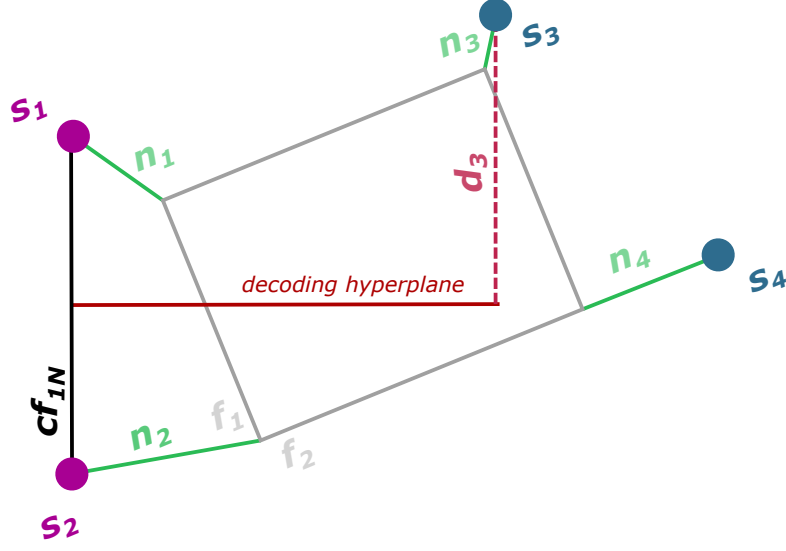
We use  $d_{LL}$  to denote the distance in the linear code between the two stimulus representations used to learn the classification and  $d_{LG}$  to denote the distance between the two stimulus representations that are generalized to. This distance is along the same unit vector  $f_1$ . We use  $d_{LA}$  to denote the distance between the two pairs of stimuli along the axis that they are to be generalized over, which we denote as the unit vector  $f_2$ . Each of the four stimuli also undergoes a distortion of magnitude  $\sqrt{P_N}$  due to the nonlinear code, we denote the direction of these four distortions as the unit vectors  $n_i$  for  $i \in [1, 2, 3, 4]$ . 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.

For simplicity, we assume that  $d_{LL} = d_{LG}$ , but this can be relaxed later.

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,  $s_1 =$



**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  $d_3 = f_{1N} \cdot s_3 - \frac{c}{2}$ .

82  $f_1 d_{LL} + d_N n_{12}$ ,  $s_2 = 0$ ,  $s_3 = d_{LL} f_1 + d_{LA} f_2 + d_N n_{23}$ , and  $s_4 = d_{LA} f_2 + d_N n_{24}$ .  
 83 Thus,

$$\begin{aligned}\hat{s}_{12} &= \frac{1}{2} (d_{LL} f_1 + d_N n_{12}) \\ \hat{s}_{34} &= \frac{1}{2} (d_{LL} f_1 + 2d_{LL} f_2 + d_N n_{23} + d_N n_{24})\end{aligned}$$

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

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

86 where  $c$  is a random variable that normalizes  $f_{1N}$  to be a unit vector, and

87 corresponds to the distance between  $s_1$  and  $s_2$ ,

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

88 Now, using  $f_{1N}$ ,  $c$ , along with our understanding of  $s_3$  and  $s_4$  as a linear com-  
 89 bination of linear and nonlinear codes, we can directly approximate CCGP.  
 90 To do this, we need to find the position of  $s_3$  and  $s_4$  along the decoding vec-  
 91 tor defined by  $f_{1N}$ , and then we evaluate whether that magnitude is greater  
 92 or smaller than the threshold  $c/2$ . So, to find this distance relative to the  
 93 threshold, we need  $d$  such that

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

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

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

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

$$\begin{aligned} cd_3 &= d_{LL}^2 + \frac{1}{2}d_N^2 + \frac{d_{LL}d_N}{\sqrt{2}} (f_1 n_1 + f_1 n_2 + f_1 n_2 + f_1 n_3) + d_{LA}d_N n_{12} f_2 - \frac{c^2}{2} \\ &= d_{LL}^2 + \frac{1}{2}d_N^2 + \frac{d_{LL}d_N}{\sqrt{2}} (f_1 n_1 + f_1 n_2 + f_1 n_2 + f_1 n_3) + d_{LA}d_N n_{12} f_2 \\ &\quad - \frac{1}{2}d_{LL}^2 - \frac{1}{2}d_N^2 - \frac{d_{LL}d_N}{\sqrt{2}} (f_1 n_1 + f_1 n_2) \\ &= \frac{1}{2}d_{LL}^2 + d_{LL}d_N f_1 n_{23} + d_{LA}d_N f_2 n_{12} \\ d_3 &= \frac{\frac{1}{2}d_{LL}^2 + d_{LL}d_N f_1 n_{23} + d_{LA}d_N f_2 n_{12}}{\sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_N f_1 n_{12}}} \\ d_3 &\sim \frac{\frac{1}{2}d_{LL}^2 + d_{LL}d_N \mathcal{N}(0, 1/D) + d_{LA}d_N \mathcal{N}(0, 1/D)}{\sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_N \mathcal{N}(0, 1/D)}} \end{aligned}$$

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

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

97 The calculation for  $s_4$  is analogous and leads to the same result, thus we drop  
98 the dependence on 3 and write

$$\begin{aligned} P(\text{CCGP error}) &= \mathbb{E}_d Q(d/\sigma) \\ &= \mathbb{E} Q\left(\frac{1}{\sigma} \frac{\frac{1}{2}d_{LL}^2 + d_{LL}d_N\mathcal{N}(0, 1/D) + d_{LA}d_N\mathcal{N}(0, 1/D)}{\sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_N\mathcal{N}(0, 1/D)}}\right) \end{aligned}$$

99 For large  $D$ , we can further simplify by ignoring the random variables (which  
100 is also the case where the nonlinear and linear codes are perfectly orthogonal  
101 to each other),

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