

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 σ^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 This distance can be used as a scaling factor that allows translation between
44 distance in stimulus space and distance in representation space – that is, if
45 two stimuli have distance l in stimulus space, then they have distance ld_L
46 in representation space. However, this assumes that each feature is encoded
47 with the same fidelity, which may not always be true.

48 Second, we find that the nonlinear distance is,

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

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

51 Third, we remind ourselves that the dot product of two random unit vectors
 52 $u_1 \cdot u_2$ in a D -dimensional space follows the distribution

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

53 for large D .

54 1.2.2 Main derivation

55 We use d_{LL} to denote the distance in the linear code between the two stim-
 56 ulus representations used to learn the classification and d_{LG} to denote the
 57 distance between the two stimulus representations that are generalized to.
 58 This distance is along the same unit vector f_1 . We use d_{LA} to denote the
 59 distance between the two pairs of stimuli along the axis that they are to be
 60 generalized over, which we denote as the unit vector f_2 . Each of the four
 61 stimuli also undergoes a distortion of magnitude $\sqrt{P_N}$ due to the nonlinear
 62 code, we denote the direction of these four distortions as the unit vectors
 63 n_i for $i \in [1, 2, 3, 4]$. They are chosen such that $n_i \cdot n_j = 0$ for $i \neq j$ –
 64 however, $n_i \cdot f_j$ is not constrained to be zero. From above, we know that
 65 $n_i \cdot f_j \sim \mathcal{N}(0, 1/D)$ where D is the full dimensionality of the space (i.e., the
 66 number of neurons in the code). Additionally, for convenience, we also use
 67 n_{ij} for any number of indices to refer to the following

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

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

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

70 First, we find the center points between our two pairs of stimuli in the full
 71 code with reference to the “bottom left” stimulus, s_2 . In particular, $s_1 =$
 72 $f_1 d_{LL} + d_N n_{12}$, $s_2 = 0$, $s_3 = d_{LL} f_1 + d_{LA} f_2 + d_N n_{23}$, and $s_4 = d_{LA} f_2 + d_N n_{24}$.
 73 Thus,

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

74 Next, we find the axes of the distorted space in the full code. The vector
 75 between s_1 and s_2 is the distorted version of f_1 , and is given by

$$f_{1N} = \frac{1}{c} (d_{LL} + d_N n_{12})$$

76 where c is a random variable representing the distance between the two stim-
 77 ulus representations in the full code,

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

78 Now, we find the optimal f_2 in the distorted code, which is also the category
 79 boundary that would be learned by a decoder trained with both stimulus
 80 representation pairs. This is given by,

$$\begin{aligned} f_{2N}^{\text{opt}} &= \hat{s}_{34} - \hat{s}_{12} \\ &= \frac{1}{a} \left(d_{LA} f_2 + \frac{\sqrt{2}}{2} d_N n_{1234} \right) \end{aligned}$$

81 where a is a random variable given by

$$a = \sqrt{d_{LA}^2 + \frac{1}{2} d_N^2 \sqrt{2} d_N d_{LL} \mathcal{N}(0, 1/D)}$$

82 which is the distance between the centers of the two sets of points.

83 The decoding hyperplane used by a decoder trained on only the first set of
 84 stimulus representations is given by the unit vector orthogonal to it, which
 85 is simply f_{1N} . If $f_{1N} \cdot f_{2N} = 0$, then the decoding hyperplane includes f_{2N}^{opt}
 86 and generalization performance will be as good as possible. Thus, we can
 87 approximate CCGP by studying the dot product of these two vectors. The
 88 dot product can be written as,

$$\begin{aligned} b &= f_{2N}^{\text{opt}} \cdot f_{1N} \\ &= \frac{1}{ac} \sqrt{\frac{3}{2}} d_N d_{LA} \mathcal{N}(0, 1/D) \end{aligned}$$

89 Geometrically, $b_{\frac{c}{2}}$ is the distance along f_{2N}^{opt} that s_1 and s_2 are from the center
90 point \hat{s}_{12} . Here, we have two sides of a right triangle, the hypotenuse has
91 length $c/2$ and the other with length $b_{\frac{c}{2}}$. We use these to find the angle
92 between the learned hyperplane and the optimal hyperplane,

$$\begin{aligned}\theta &= \frac{\pi}{2} - \arccos\left(\frac{2b}{c}\right) \\ &= \arcsin\left(\frac{2b}{c}\right)\end{aligned}$$

93 In most cases, the larger θ , the worse generalization performance will be.

94 We can also directly approximate CCGP. To do this, we need to find the
95 position of s_3 and s_4 along the decoding vector defined by f_{1N} , and then
96 we evaluate whether that magnitude is greater or smaller than the threshold
97 $c/2$. So, to find this distance relative to the threshold, we need d such that

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

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

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

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

$$\begin{aligned}cd_3 &= d_{LL}^2 + \frac{1}{2}d_N^2 + \frac{d_{LL}d_N}{\sqrt{2}} (f_1 n_1 + f_1 n_2 + f_1 n_2 + f_1 n_3) + d_{LA}d_N n_{12} f_2 - \frac{c^2}{2} \\ &= d_{LL}^2 + \frac{1}{2}d_N^2 + \frac{d_{LL}d_N}{\sqrt{2}} (f_1 n_1 + f_1 n_2 + f_1 n_2 + f_1 n_3) + d_{LA}d_N n_{12} f_2 \\ &\quad - \frac{1}{2}d_{LL}^2 - \frac{1}{2}d_N^2 - \frac{d_{LL}d_N}{\sqrt{2}} (f_1 n_1 + f_1 n_2) \\ &= \frac{1}{2}d_{LL} + d_{LL}d_N f_1 n_{23} + d_{LA}d_N f_2 n_{12} \\ d &= \frac{\frac{1}{2}d_{LL} + d_{LL}d_N f_1 n_{23} + d_{LA}d_N f_2 n_{12}}{\sqrt{d_{LL}^2 + d_N^2 + 2d_{LL}d_N f_1 n_{12}}}\end{aligned}$$

$$a_{3,4} = a \pm \frac{\sqrt{2}}{2a} d_N \left(\frac{1}{2} d_{LA} \mathcal{N}(0, 1/D) + d_{LL} \mathcal{N}(0, 1/D) \right)$$

100 where the second term captures the distortion relative to the mean distance
 101 between the two stimuli. Then, we can find the distance of s_i for $i \in [3, 4]$
 102 from the learned decoding hyperplane by finding

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