ANALYZING REPRESENTATION AND TRANSFORMATIONS

• Reminder: Spafford's lab Mar 3, B1 172, 11:30a

Definitions

- A *representational space* is a specific kind of vector space.
- Vectors are simply collections of mathematical objects (numbers, functions, or other vectors)
- A *vector space* is a set of vectors that is *closed* under addition and multiplication (i.e., a sum or product of any vectors in the set is also in the set)
- A *basis* is an *independent* set of vectors that *span* the vector space

Definitions: independent

• A set of vectors \mathbf{x}_n is independent if $a_1\mathbf{x}_1 + a_2\mathbf{x}_2 + \ldots + a_n\mathbf{x}_n = 0$

only when

$$a_1 = a_2 = \ldots = a_n = 0$$

• where *a_n* are scalars.

• So, *n* must be equal to the dimension of the vectors in **V** in order for **x**_n to be independent.

Definitions: span

- A set of vectors *spans* a vector space if any vector in that space can be written as a linear sum of those vectors.
- That is, if for all $\mathbf{x} \in \mathbf{V}$ there are some a_n , $a_1\mathbf{x}_1 + a_2\mathbf{x}_2 + \ldots + a_n\mathbf{x}_n = \mathbf{x}$
- then the set of vectors **x**^{*n*} span the vector space **V**.

Example

• The standard *Cartesian basis* in any number of dimensions is comprised of the unit vectors along the principle axes (obviously independent and spanning V).



Definitions: orthonormal

 Note also, the vectors in this basis are orthogonal, i.e., the dot product of any two will be zero; i.e.,

$$\mathbf{x} \cdot \mathbf{y} = \langle \mathbf{x}\mathbf{y} \rangle_n = \sum x[n]y[n] = 0$$

n

- If all the vectors in a basis are orthogonal, it is called an orthogonal basis
- Further, if the length of all the vectors in an orthogonal basis is equal to one then it is an orthonormal basis: the Cartesian basis is such

Overcomplete basis

- If we relax the constraint that the vectors have to be independent, we have what is called an *overcomplete* basis (or sometimes `frame').
 - overcomplete bases are redundant (hence not as succinct) for defining the vector space
 - in noisy, physical systems, this redundancy can prove invaluable for error correction and the efficient use of available resources

A longer example

- Let the vector **x** be written in a standard Cartesian basis $\mathbf{x} = [x_1, x_2] = x_1 \mathbf{i} + x_2 \mathbf{j}$
- Suppose we have a different basis, rotated by 45 degrees (φ₁, φ₂)
- 'Encoding' $a_i = \langle \mathbf{x} \boldsymbol{\phi}_i \rangle_n$
- Projection onto the new basis a = [a₁, a₂] = a₁φ₁ + a₂φ₂
 'Decoding'

$$\mathbf{x} = \sum_i a_i \boldsymbol{\phi}_i$$



Comments

- This way, we can move back and forth between orthonormal bases.
- Notice that if we substitute the encoding into the decoding, we recover the same vector we originally encoded
- Thus we can think of the coefficients *a_i* as 'representing', or 'carrying the same information', or 'encoding' the original coefficients *x_i*

Overcomplete representation

• Suppose we do not know what the decoding basis is, but we know the encoding basis

- We can guarantee that the encoding basis is overcomplete by using redundant, nonorthogonal encoders
- Let's choose the encoding basis as equally spaced at 120 degree intervals (i.e., $\tilde{\phi}_1 = [\frac{\sqrt{3}}{2}, \frac{1}{2}], \tilde{\phi}_2 = [-\frac{\sqrt{3}}{2}, \frac{1}{2}], \text{ and } \tilde{\phi}_3 = [0, -1]).$



Finding decoders

- We need to identify the set of vectors that span the space into which our vector is being encoded (i.e., the decoding basis).
- To find this basis, we first note, as before, that

$$\mathbf{x} = \sum_{i} a_i \boldsymbol{\phi}_i$$

• and (notice the tilde! unlike earlier)

$$a_i = \left\langle \mathbf{x} \tilde{\boldsymbol{\phi}}_i \right\rangle_n$$

Finding decoders

Substitute the encoding into decoding to give

$$\mathbf{x} = \sum_{i} \left\langle \mathbf{x} ilde{\phi}_{i} \right\rangle_{n} \phi_{i}$$

• Writing the dot product explicitly, we get x[m] $x[m] = \sum_{i,n} x[n] \tilde{\phi}_i[n] \phi_i[m]$ $\sum_{i,n} x[n] \sum_{i,n} \tilde{\phi}_i[n] + [m]$

$$= \sum_{n} x[n] \sum_{i} \phi_{i}[n] \phi_{i}[m]$$

 $\delta_{nm} = \sum \tilde{\phi}_i[n]\phi_i[m]_i$

• So,

Finding decoders

• In matrix notation, $I = \phi \phi$ $(a_1\phi_1 \ a_2\phi_2 \ a_3\phi_3)$ Solving for the decoders 1 $\mathrm{I}~=~ ilde{\phi}\phi$ ¢2 > 0 $ilde{\phi}^T = ilde{\phi}^T ilde{\phi} \phi,$ \$3 $oldsymbol{\phi} \;\; = \;\; \left(ilde{oldsymbol{\phi}}^T ilde{oldsymbol{\phi}}
ight)^{-1} ilde{oldsymbol{\phi}}^T$ -2 -2 Calculating gives -1 0 2 $\phi_1 = \left[\frac{\sqrt{3}}{3}, \frac{1}{3}\right], \phi_2 = \left[-\frac{\sqrt{3}}{3}, \frac{1}{3}\right], \text{ and } \phi_3 = \left[0, -\frac{2}{3}\right]$ • That is, $\phi_i = \frac{2}{3}\tilde{\phi}_i$

Comments

- The φ_i are an *overcomplete basis* for the vector space in which the a_i are coordinates. Hence are the decoding vectors defining the representation
- A set of encoders / decoders is like this is called biorthogonal because together they act as an orthogonal basis.
- Technically, the NEF encoders / decoders don't satisfy the definitions provide here, but they are close

Basis functions

- Applies equally to basis *functions*
- Just as basis vectors define a vector space, so basis functions define a function space
- E.g. sines and cosines for a Fourier decomposition. This is an orthonormal basis.
- Can have overcomplete basis functions as well.
- The tuning curves are like a set of overcomplete basis functions for the space they represent.

• Consider the noise free case

 $\gamma_{ij} = \left\langle a_i(\mathbf{x}) a_j(\mathbf{x}) \right\rangle_{\mathbf{x}}$

 This is called a Gram matrix (a correlation matrix without the means removed). Measures similarity of all neural responses.

• Tells us about the representation of the vector space by the neurons

Construct an 'activity matrix'

 $\mathbf{A}^{T} = \begin{bmatrix} a_{1}(\mathbf{x}_{min}) & a_{1}(\mathbf{x}_{min} + \Delta_{\mathbf{x}}) & \cdots & a_{1}(\mathbf{x}_{max} - \Delta_{\mathbf{x}}) & a_{1}(\mathbf{x}_{max}) \\ a_{2}(\mathbf{x}_{min}) & & a_{2}(\mathbf{x}_{max}) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{N-1}(\mathbf{x}_{min}) & a_{1}(\mathbf{x}_{min} + \Delta_{\mathbf{x}}) & \cdots & a_{1}(x_{max} - \Delta_{\mathbf{x}}) & a_{N}(\mathbf{x}_{max}) \end{bmatrix}$

• So we can write our estimate

$$\hat{\mathbf{X}}_{N_{\Delta} \times N_{v}} = \mathbf{A}_{N_{\Delta} \times N} \boldsymbol{\phi}_{N \times N_{v}}$$

• Need the optimal decoders, so take **X** as given (i.e., the actual domain)

$$\mathbf{X} = \mathbf{A}\boldsymbol{\phi}$$
$$\mathbf{A}^T\mathbf{X} = \mathbf{A}^T\mathbf{A}\boldsymbol{\phi}$$
$$(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{X} = \boldsymbol{\phi}.$$

• Taking the inverse correctly results in linear optimal decdoders, we can write

γ	=	$\mathbf{A}^T \mathbf{A}$
Υ	=	$\mathbf{A}^T \mathbf{X}$

• In other words,

 $\phi = \gamma^{-1} \Upsilon$ $= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{X}$

• Which means

 $\hat{\mathbf{X}} = \mathbf{A}\boldsymbol{\phi} \\ = \mathbf{A}(\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{X}.$

• How do we take the inverse correctly?

- Because there is no noise, and tuning curves are very similar, the matrix is likely to be singluar
- Use singular value decomposition (SVD) to decompose and approximately invert the matrix.
- SVD decomposition of an MxN matrix, **B**, gives

 $\mathbf{B}_{M\times N} = \mathbf{U}_{M\times N} \mathbf{S}_{N\times N} \mathbf{V}_{N\times N}^T$

• S is a diagonal matrix of singular values

- When B is square and symmetrical (as for γ), this simplifies to γ = USU^T
 Also, γ⁻¹ = US⁻¹U^T
- When γ is singular (or nearly so), some elements of S are zero (or very small), so the inverse of S includes infinite (or very large) terms
- The SVD `pseudo-inverse' is defined where for S_i=0, the inverse is set to 0 (∞=0!).

- SVD can be very informative:
 - The columns of U whose corresponding singular values are non-zero form an orthonormal basis that spans the range of γ
 - The columns of U whose corresponding SVs are zero form an orthonormal basis that spans the null space (i.e., x s.t. γ x=0).

- When a vector in *Υ* lies in the range of *γ*, the SVD pseudo-inverse guarantees that the corresponding vector from φ minimizes the length of that φ vector.
- Given that γ is singular, there are an infinite number of solutions for φ. The solution that provides the shortest vector is a natural choice from the set

SVD (singular matrix)



- When a vector in *Υ* lies outside the range of *γ*,
 the SVD pseudo-inverse guarantees that the best (in the least squares sense) φ given *Υ* is found.
- I.e., this `pseudo-inverse' minimizes the error. Thus, we can use SVD to find the optimal decoding functions, which we can now write as

$\boldsymbol{\phi} = \mathbf{U}\mathbf{S}^{-1}\mathbf{U}^T\mathbf{A}^T\mathbf{X}$

• Which is gives the same ϕ as always.

• Let's do the same thing to function decoders $f(\mathbf{X}) = \mathbf{A}\phi^{f}$ $\mathbf{A}^{T}f(\mathbf{X}) = \mathbf{A}^{T}\mathbf{A}\phi^{f}$ $\phi^{f} = (\mathbf{A}^{T}\mathbf{A})^{-1}\mathbf{A}^{T}f(\mathbf{X})$ $= \gamma^{-1}\Upsilon$

 So we can find lease-squares optimal decoders using

 $\boldsymbol{\phi}^f = \mathbf{U}\mathbf{S}^{-1}\mathbf{U}^T\mathbf{A}^Tf(\mathbf{X})$

- So, the representational decoder, is found in the special case where *f*(**x**)=**x**.
- Regardless of which transformation we need decoders for, we always perform SVD on the same matrix, $\gamma = A^T A$.
- Understanding the properties of γ can provide insight into all possible decodings of the population, *a_i*.

- SVs tell us the importance of the corresponding U vector. Importance being:
- related to the error that would result if we left a particular vector out of the mapping
- related to the variance of population firing along the vectors in the γ matrix.
- the amount of (independent) information about population firing that can be extracted by looking only at data projected onto that U vector.

- In general, we can think of the magnitude of the singular value as telling us how relevant the dimension defined by the corresponding U vector is to the identity of the matrix we have decomposed.
- Since the matrix we have decomposed is like the correlation matrix of the neuron tuning curves, the large singular values are most important for accounting for the structure of those correlations.

• Notice that the vectors in **U** are orthogonal:

- they provide an (ordered) orthogonal basis for that matrix, where the original γ matrix was generated by a non-ordered non-orthogonal basis (the neuron tuning curves)
- Consider a point in 'neuron space' a = a₁e₁ + a₂e₂ + ... + a_Ne_N
 Here, e_i are axis vectors, one for each neuron with activity a_i.

- Because the neural responses are nonindependently driven by some variable, x, only a subspace of the space spanned by the e_i vectors is ever actually occupied by the population.
- The γ matrix, because it tells us the correlations between all neurons in the population, provides us with the information we need to determine what that subspace is (i.e., U)

• Let's do some math: $\hat{\mathbf{X}} = \mathbf{A}\mathbf{U}\mathbf{S}^{-1}\mathbf{U}^T\mathbf{A}^T\mathbf{X}$

• Or more simply $\hat{\mathbf{X}} = \chi \Phi$

• So, $\chi = \mathbf{AU}$

 $\chi_m(\mathbf{x}) = \sum_i a_i(\mathbf{x}) U_{im}$

• or

• and continuing $\Phi = \mathbf{S}^{-1}\mathbf{U}^T\mathbf{A}^T\mathbf{X}$ $= \mathbf{U}^T\mathbf{U}\mathbf{S}^{-1}\mathbf{U}^T\mathbf{A}^T\mathbf{X}$ $= \mathbf{U}^T\phi$

• SO

$$\Phi_m = \sum_i U_{mi} \phi_i$$

 Notice that χ and Φ are rotated versions of A and φ respectively. They are rotated into the coordinate system defined by U.

- We can think of U as the rotation matrix that aligns the first axis of the coordinate system along the dimension with the greatest variance in the encoding of x, the second axis along the dimension with the second greatest variance, etc.
- So the χ vectors also end up being orthogonal and ordered by importance
- So they are basis functions for the function space computable from **A**

SVD 3D projections



- Whichever χ(x) functions have reasonably large associated singular values, are exactly the functions that we can do a good job of extracting from our encoding of the input space
- We can also extract any linear combinations of those functions quite well.
- Because these functions are ordered, the more useful the `first' function is for reconstructing some transformation, *f*(x), the better we can extract it

• Let's apply this analysis. Start with broad tuning.



• We get the following (normalized) functions. (b) is the Legendre polynomials

$$l_i(x) = \frac{(-1)^i}{2^i i!} \frac{d^i}{dx^i} \left[\left(1 - x^2 \right)^i \right]$$



- The similarity between χ(x) and l_i(x) means that this neural population supports the extraction of functions that can be well-estimated using the standard Legendre basis.
- But the χ(x) functions are ordered by their singular values. Thus, the higher-order polynomial terms are not as well encoded by our population as the lower-order ones

Comments

- Tuning curves are very broad, and the polynomial basis is also very broad
- The *l_i* basis is ordered by decreasing linearity so we expect the functions in precisely that order
- None of this would be true if the population did not evenly tile the input space.
- The basis does not just depend on the general `shape' of the neuron tuning curves, but also on the 'heterogeneity' of the population

Higher dimensions

- n-dimensional vectors have cross terms in the computable functions $f(\mathbf{x}) = c_0 + c_1 x_1 + c_2 x_2 + c_3 x_1^2 + c_4 x_2^2 + c_5 x_1 x_2 + \dots$ $= \sum_{l=0}^{N_{order}} \sum_{n=0}^{l} c_{n,l-n} x_1^n x_2^{l-n}$
- Lots of cross terms, how big are the sigular values?

Multiset and singular values



• So, higher-D means higher order is harder.

Gaussian tuning curves

Same analysis for Gaussian tuning



Gaussian basis functions



SVs: Gaussian vs Linear

• Comparing scalar singular values for the two:



Comments

- Gaussian basis looks like a Fourier basis (but pinned because of the end effects)
- The Gaussian basis is better for encoding localized functions (see singular values)
- In general, examining χ(x) determines what kinds of transformations are and are not wellsupported by that population.
- Extremely useful for constraining hypotheses

- Essentially everything is the same, except the function 'lines' become 'tubes' of uncertainty
- Importantly, however, the noise doesn't scale with the singular values, it is isometric in the vector space.
- This means that small singular values are more greatly affected by noise.

Error with noise

• Write as before

$$E = \left\langle \left[\mathbf{x} - \sum_{i} \left(a_{i}(\mathbf{x}) + \eta_{i} \right) \boldsymbol{\phi}_{i} \right]^{2} \right\rangle_{\mathbf{x},\eta}$$
$$= \left\langle \left[\mathbf{x} - \sum_{i} a_{i}(\mathbf{x}) \boldsymbol{\phi}_{i} \right]^{2} + \sigma_{\eta}^{2} \sum_{i} \boldsymbol{\phi}_{i}^{2} \right\rangle_{\mathbf{x}}$$

• We can now use our SVD expressions to give

$$E = \left\langle \left[\mathbf{x} - \sum_{m} \chi_{m} \Phi_{m} \right]^{2} + \sigma_{\eta}^{2} \sum_{i,j} \phi_{i} \delta_{ij} \phi_{j} \right\rangle_{\mathbf{x}}$$
$$= \left\langle \left[\mathbf{x} - \sum_{m} \chi_{m} \Phi_{m} \right]^{2} + \sigma_{\eta}^{2} \sum_{i,j} \phi_{i} \sum_{m} U_{im} U_{mj} \phi_{j} \right\rangle_{\mathbf{x}}$$
$$= \left\langle \left[\mathbf{x} - \sum_{m} \chi_{m} \Phi_{m} \right]^{2} + \sigma_{\eta}^{2} \sum_{m} \Phi_{m}^{2} \right\rangle_{\mathbf{x}}$$

• Minimize with respect to rotated decoders:

$$\frac{dE}{d\Phi_n} = \left\langle 2 \left[\mathbf{x} - \sum_m \chi_m \Phi_m \right] (-\chi_n) + 2\sigma_\eta^2 \Phi_n \right\rangle_{\mathbf{x}}$$
$$0 = -2 \left\langle \chi_n \mathbf{x} \right\rangle_{\mathbf{x}} + 2 \left\langle \sum_m \chi_m \chi_n \Phi_n \right\rangle_{\mathbf{x}} + 2\sigma_\eta^2 \Phi_n$$

 $\langle \chi_n \mathbf{x} \rangle_{\mathbf{x}} = S_n \Phi_n + \sigma_\eta^2 \Phi_n$ $\Phi_n = \frac{\langle \chi_n \mathbf{x} \rangle_{\mathbf{x}}}{S_n + \sigma_\eta^2}$

• So, the residual error will be

$$E_{r} = \left\langle \left[\mathbf{x} - \hat{\mathbf{x}} \right]^{2} \right\rangle_{\mathbf{x},\eta}$$

$$= \left\langle \left[\left[\mathbf{x} - \sum_{m} \chi_{m} \Phi_{m} \right]^{2} \right\rangle_{\mathbf{x},\eta}$$

$$= \left\langle \mathbf{x}^{2} \right\rangle_{\mathbf{x},\eta} - 2 \left\langle \mathbf{x} \sum_{m} \chi_{m} \Phi_{m} \right\rangle_{\mathbf{x},\eta} + \left\langle \sum_{m} \chi_{m} \Phi_{m} \right\rangle_{\mathbf{x},\eta}^{2}$$

• Substituting the optimal rotated decoders:

$$E_{r} = \langle \mathbf{x}^{2} \rangle_{\mathbf{x},\eta} - 2 \left\langle \mathbf{x} \sum_{m} \chi_{m} \frac{\langle \chi_{m} \mathbf{x} \rangle_{\mathbf{x}}}{S_{m} + \sigma_{\eta}^{2}} \right\rangle_{\mathbf{x},\eta} + \dots$$
$$\left\langle \sum_{m} \chi_{m} \frac{\langle \chi_{m} \mathbf{x} \rangle_{\mathbf{x}}}{S_{m} + \sigma_{\eta}^{2}} \right\rangle_{\mathbf{x},\eta}^{2}$$
$$= \langle \mathbf{x}^{2} \rangle_{\mathbf{x}} - 2 \sum_{m} \frac{\langle \chi_{m} \mathbf{x} \rangle_{\mathbf{x}}^{2}}{S_{m} + \sigma_{\eta}^{2}} + \sum_{m} (S_{m} + \sigma_{\eta}^{2}) \frac{\langle \chi_{m} \mathbf{x} \rangle_{\mathbf{x}}^{2}}{(S_{m} + \sigma_{\eta}^{2})}$$
$$= \langle \mathbf{x}^{2} \rangle_{\mathbf{x}} - \sum_{m} \frac{\langle \chi_{m} \mathbf{x} \rangle_{\mathbf{x}}^{2}}{S_{m} + \sigma_{\eta}^{2}}$$

 O_n^-

m

2

Comments

• Shows how much the *m*th basis function reduces the error in our estimate under noise

- As the singular value approaches the value of the variance of the noise, the corresponding element does not usefully contribute to the representation.
- When the noise becomes near the magnitude of that normalization term (i.e., SNR = 1 or less), the projection becomes `mis-normalized' and thus contributes incorrectly to the representation

Comments

- So, those basis functions whose corresponding singular value is equal to or smaller than the noise, should not be used in a good representation
- So, we can simply `lop off' some of the singular values when doing the inverse to get the same result as including a certain amount of noise in our calculation of Γ (more analysis needed)

Heterogeneity



RMS error for distributions



Heterogeneity

- The precise nature of the tuning curves matters.
- In the book we show that heterogeneity is a good balance between tiling and ease of construction
- Heterogeneity provides good reduction of error under noise.
- The representational capacity (number of perfectly represented points) is also high for heterogeneous populations.