

# A Neurocognitive Graph Theoretical Approach to Understanding the Relationship Between Minds and Brains

Joshua T. Vogelstein, R. Jacob Vogelstein, Carey Priebe  
Johns Hopkins University

January 6, 2010

**Current Paradigm** The dominant paradigm of quantitative neuroscience in the 20<sup>th</sup> century has been the “signal processing” framework []. Essentially, the brain is a box, that filters some stimulus, to produce some response (see Figure 1). This framework leads to the following goal:

**Goal 1.** *Learn the filter that the brain performs on stimuli to result in the actualized responses.*

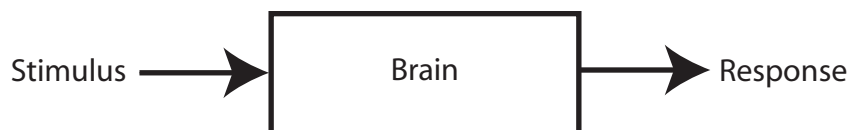


Figure 1: The signal-processing paradigm of quantitative neuroscience. The brain is a box that essentially *filters* the stimulus, outputting some response, which is often take to be multivariate time series (such as populations of spike trains or fMRI activity).

This paradigm has been appropriate, given the kind of data available to people investigating the brain. More specifically, the kind of data that has been most available has been time-series of signals related to neural activity []. Given that electrical engineers were largely the individuals obtaining and analyzing the data, it was natural to take a signal-processing approach. Often, the dynamic time-series data were used to estimate static parameters. In the previous decades, these communities have been more and more sophisticated models and algorithms to estimate these parameters []. Recently, issues such as parameter identifiability, consistency, bias, and model selection have been gaining traction as important desiderata for our models [].

**Alternate Paradigm** Here, we define a different goal, which suggests a complementary research paradigm to the current dogma:

**Goal 2.** *Construct a family of brain-models,  $\mathcal{B}$ , that is sufficient to provide causal explanations relating properties of minds with properties of brains.*

Note how the above goal is distinct in certain respects from the “filtering” goal. First, neither stimulus nor response is explicitly incorporated into this goal. While stimuli and response may be used as tools to obtain the parameters of  $\mathcal{B}$ , that is their only merit in this paradigm. Second, minds are explicitly incorporated into this goal. While spike trains or fMRI signal may indicate mental processes, they are likely not what is meant by “mind”; rather, they may be used as tools to infer mental states. Third, the inclusion of a class of models  $\mathcal{B}$  suggests a statistical *static* framework, rather than a dynamics perspective. In addition to the above stated goal, we would like the family of models  $\mathcal{P}$  to satisfy a number of desiderata:

1.  $\mathcal{B}$  should be sufficiently general to account for whatever properties of the brain are casually related to the properties of cognition under investigation.
2. The properties of any particular brain,  $b \in \mathcal{B}$ , should be either measurable or estimatable, such that experimental observation may be used to obtain them.

3.  $\mathcal{B}$  should admit algorithms that (are guaranteed to be able to) capture the relationship of interest.
4.  $\mathcal{B}$  should also admit *causal* studies, which entail modifying a particular  $b \in \mathcal{B}$ , to modify the corresponding mental property,  $m \in \mathcal{M}$ .

**NeuroCognitive Graph Theory** Here, we propose a novel approach, called NeuroCognitive Graph (NeCoG) theory, which we believe achieves the above goal and its associated desiderata. Specifically, we say that the brain may be well characterized as a labeled, attributed multigraph (which is a generalized notion of a or network<sup>1</sup>). Formally, we define a brain-graph,  $b \in \mathcal{B}$  as a 4-tuple,  $\mathcal{B} = (\mathcal{V}, \mathcal{E}, \mathcal{X}_V, \mathcal{X}_E)$ , where

- $V_i \in \mathcal{V} \subseteq \mathbb{Z}$  for  $i \in [V]$ , is the set of vertices (nodes), where  $[V] = \{1, 2, \dots, V\}$
- $E_{ijk} \in \mathcal{E} \subseteq [V] \times [V] \times \mathbb{Z}$  for  $(i, j, k) \in [V] \times [V] \times \mathbb{Z}$  is the set of (directed) edges, where  $\mathbb{Z} = \{0, 1, \dots\}$
- $\mathbf{X}_i \in \mathcal{X}_V \subseteq \mathbb{R}^{d_V}$  for  $i \in [V]$  is the set of vertex attributes, and  $d$  is the dimensionality of the vertex attributes (potentially including vertex labels)
- $\mathbf{X}_{ijk} \in \mathcal{X}_E \subseteq \mathbb{R}^{d_E}$  for  $(i, j, k) \in [V] \times [V] \times \mathbb{Z}$  is the set of edge attributes,  $d'$  is the dimensionality of the edge attributes (potentially including edge labels), and  $k$  indexes edge categories.

Assume, for the moment, that we take the fundamental computational unit of the brain to be a point neuron. Then, each vertex is a neuron, and each edge is a synapse. Multi-edges potentially correspond to different edge magnitudes and/or categorically different edges (like chemical and electrical synapses). Vertex attributes could include neurotransmitter released, proteins expressed, morphological properties, receptive fields, etc. Edge attributes could include probability of release, post-synaptic potential shape, etc. This level of description, however, is not necessary. For instance,  $\mathcal{V}$  might instead correspond to neuroanatomical regions, which admit very different notions for edges and attributes. This multi-scale aspect of  $\mathcal{B}$  is an important advantage over other frameworks.

Given  $\mathcal{B}$ , what can we do with it? As stated above, our goal is to relate these models to properties of cognition. More specifically, let  $\mathcal{M}$  characterize the space of the cognitive (mental) property, and  $g \in \mathcal{G}$  be some function to learn.

- If  $\mathcal{M} = \{0, 1\}$ , then  $g$  is a two-way classifier:  $g : \mathcal{B} \rightarrow \{0, 1\}$ .
- If  $\mathcal{M} = \{0, 1, \dots, C\}$ , then  $g$  is an  $C$ -way classifier:  $g : \mathcal{B} \rightarrow \{0, 1, \dots, C\}$ .
- If  $\mathcal{M} = \mathbb{R}^a$ , then  $g$  is a (multivariate-) regressor:  $g : \mathcal{B} \rightarrow \mathbb{R}^a$ .

In general, solving the above problems—which means finding  $g$ —will depend on  $F = F_{BM}$ , the joint distribution of brains and minds. In practice, however,  $F$  is unknown, and therefore  $g$  must be estimated from the data. Assume we have a corpus of training data,  $\mathcal{D}_n = \{(B_1, M_1), \dots, (B_n, M_n)\}$ , where  $n$  is the number of training samples. Our goal then is to compute  $g_n : \mathcal{B} \times (\mathcal{B}, \mathcal{M})^n \rightarrow \mathcal{M}$ , which takes as input an observed brain-graph  $b$  and  $n$  training pairs  $\{(b_1, m_1), \dots, (b_n, m_n)\}$  and produces a prediction  $\hat{m} = g_n(b; \mathcal{D}_n)$ . The goal is to find a  $g_n$  that minimizes some loss function,  $L_F(g_n)$ . For instance, when  $|\mathcal{M}| = 2$ , a potentially reasonable loss function is  $L_F(g_n) = \mathbb{E}[P_F[g_n(B; \mathcal{D}_n) \neq M | \mathcal{D}_n]]$ .

**Finding a good  $g_n$**  Thus, given a mental property, a decision about how to represent it,  $\mathcal{M}$ , and a loss function,  $L$ , our task is to find a good  $g_n$ . Vogelstein et al. (2010) showed that a  $k_n$  nearest neighbor (knn) classifier is a universally consistent classifier (meaning, achieves the Bayes optimal performance), for  $F_{BM}$ , under a Frobenius norm distance metric. While this is a desirable property, our belief is that other  $g_n$ 's may outperform the knn classifier on finite data sets. More specifically, while knn induces no bias whatever into  $g_n$ , the variance is large. Thus, by incorporating neuroscientific knowledge about these brain-graphs into  $g_n$ , it may be possible to only marginally increase the bias, but drastically reduce the variance, yielding improved performance.

One possible strategy is to propose a class of models,  $\mathcal{P} = \{P_\theta : \theta \in \Theta \subseteq \mathbb{R}^d\}$ , whose dimensionality  $d$  depends on the data. The goal then is to find a Minimally Sufficient Model (MSM; by analogy with minimally sufficient statistics), which is the brain-graph with the least parameters that is sufficient to explain the mental property under

<sup>1</sup>I think binary graphs should probably go here, with extensions later.

investigation.  $g_n$  then operates directly on  $\hat{\theta}$ , the data dependent estimate of the model parameters,  $\theta$ . Classification/regression is then performed on  $\hat{\theta}$ , which hopefully lives in some space smaller than  $\mathcal{B}$ , effectively reducing the variance, without increasing bias too much.

**Generative Model** Consider the following generative model for attributed multigraphs:

- Let  $c$  be a *class identity*, where  $c \in \mathcal{C} = \{0, 1, \dots, C\}$
- Let  $\theta_c$  be the *parameters* for class  $c$ , where  $\theta_c \in \Theta \subseteq \mathbb{R}^d$ , for some  $d = d_V + d_E \in \mathbb{Z}$ , where the dimensionality of the parameters is implicitly a function of the data,  $\mathcal{D}_n$ .
- Let  $G$  be a *graph*, where  $G(\theta_c) \in \mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{X}_V, \mathcal{X}_E)$ , where for clarity, we restrict edges to be integer weights (i.e., only include a single category of edge attributes, but this may straightforwardly generalized)

To sample graphs from this generative model, assuming that  $C$  and  $V$  are given (the number of classes and vertices per graph, respectively), one can use the following procedure (generalizing to the unknown  $V$  case is straightforward and therefore omitted):

- sample  $c \sim f_c$
- sample  $\theta_c \sim f_\theta(\cdot|c)$
- for  $i \in [V]$ , sample  $\mathbf{X}_i \sim f_{X_V}(\cdot|\theta_c^V)$
- for  $(i, j) \in [V] \times [V]$ 
  - sample  $\mathbf{X}_{ij} \sim f_{X_E}(\cdot|\theta_c^E)$
  - sample  $E_{ij} \sim f_E(\cdot|\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_{ij})$

Note that we have partitioned  $\theta_c$  into  $\theta_c^V$  and  $\theta_c^E$ . The probability of obtaining any graph, when using this procedure, is therefore given by:

$$P(G|C, V) = \left( \prod_{i,j \in [V]} f_E(E_{ij}|\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_{ij}) f_{X_E}(\mathbf{X}_{ij}|\theta_c^E) \right) \left( \prod_{i \in [V]} f_{X_V}(\mathbf{X}_i|\theta_c^V) \right) f_\theta(\theta_c|c) P(c) \quad (1)$$

Evaluating Eq. 1 requires defining  $f = \{f_c, f_\theta, f_{X_V}, f_{X_E}, f_E\}$ . The most natural choice for  $f_c$  is a multinomial, that is,  $f_c = MN(\vec{\pi})$ , where  $\vec{\pi} = \{\pi_1, \dots, \pi_C\}$ ,  $\pi_c > 0$ , and  $\sum_c \pi_c = 1$ . Our task is then to choose  $f$  that yields: (i) models that display the properties of the data, and (ii) are statistically tractable, meaning that  $\theta$  may be estimated consistently from the data. The above generative framework generalizes and unifies previously proposed stochastic graph models. Consider a few examples:

**Random Dot Product Graphs** Let  $\Theta = \mathbb{R}^d$ , and  $f_{X_V}$  be the identity function,  $f_{X_E}$  is a constant, and  $f_E = f(\langle \mathbf{X}_i, \mathbf{X}_j \rangle)$ , where  $\langle \cdot, \cdot \rangle$  indicates a dot product

To make this a fully *generative mode*, we must propose a distribution from which the  $\mathbf{X}_i$ 's are sampled. Assuming  $\mathbf{X} \in \mathbb{R}^d$ , then we could let each  $\mathbf{X}_i$  be sampled from a multivariate normal, that is, let  $\mathbf{X}_i \sim \mathcal{N}(\mu_i, \Sigma_i)$ . To be more general, we could concatenate all  $\mathbf{X}_i$ 's,  $\vec{\mathbf{X}} = [\mathbf{X}_1, \dots, \mathbf{X}_V]$ , and proposed that  $\vec{\mathbf{X}} \sim \mathcal{N}(\vec{\mu}, \vec{\Sigma})$ . Imposing constraints on  $\Sigma$  (such as block-diagonal) reduces this more general model to the previous one. If we are in the classification setting, with  $z$  classes, the goal then is to estimate  $\vec{\mu}_z$  and  $\vec{\Sigma}_z$ , for all  $z = 1, \dots, Z$ . The class of  $b$  is then the  $m_z$  with the highest posterior. Thus, a natural advantage of a fully generative model is that it admits a principled loss function.

**Random Dot Product Graphs** One particularly compelling model is the Random Dot Product Graph (RDPG), where the value of an edge between two vertices is a function of the dot product between the vertex attributes. For binary adjacency matrices, for example, we have  $P[E_{ij}|\mathbf{X}_i, \mathbf{X}_j] = f(\mathbf{X}_i, \mathbf{X}_j)$ . More generally, we have  $\mathbb{E}[E_{ijk}|\mathbf{X}_i, \mathbf{X}_j, k] = f(\mathbf{X}_i, \mathbf{X}_j, k)$ . This RDPG should satisfy several criteria, as  $n \rightarrow \infty$ :

1. the parameters should be identifiable (up to an arbitrary scale constant) and asymptotically unbiased, that is  $\widehat{\mathbf{X}}_i \rightarrow \mathbf{X}_i$
2. the dimensionality of the vertex attributes,  $d$  should go to infinity
3. the number of distinct  $\mathbf{X}$ 's should go to infinity, that is  $|\mathcal{X}_v| \rightarrow \infty$

To ensure the first criterium,  $f(\cdot)$  must be carefully chosen. For example, letting the number of edges between any pair of vertices correspond to an integer weight on the edge, we could let  $P[E_{ij}|\mathbf{X}_i, \mathbf{X}_j] = \text{Poisson}(e^{\mathbf{X}_i^\top \mathbf{X}_j})$  []. Note that this model can implicitly include labels on vertices and edges, explicitly includes vertex attributes, and not does include edge attributes.

**Limitations/extensions** attributed vertices, hyperedges

**Simulated applications**

**Concluding thoughts**