Modeling Semantic Fluency Data as Search on a Semantic Network

Jeffrey C. Zemla (zemla@wisc.edu) Joseph L. Austerweil (austerweil@wisc.edu) Department of Psychology, University of Wisconsin-Madison, 1202 W. Johnson Street Madison, WI 53706 USA

Abstract

Psychologists have used the semantic fluency task for decades to gain insight into the processes and representations underlying memory retrieval. Recent work has suggested that a censored random walk on a semantic network resembles semantic fluency data because it produces optimal foraging. However, fluency data have rich structure beyond being consistent with optimal foraging. Under the assumption that memory can be represented as a semantic network, we test a variety of memory search processes and examine how well these processes capture the richness of fluency data. The search processes we explore vary in the extent they explore the network globally or exploit local clusters, and whether they are strategic. We found that a censored random walk with a priming component best captures the frequency and clustering effects seen in human fluency data.

Keywords: memory; search; semantic networks; fluency

Introduction

One important task for the human mind is to retrieve knowledge from memory when it is needed. To investigate how the mind solves this task, psychologists have used the semantic fluency task (Bousfield, & Sedgewick, 1944), in which participants generate as many unique items as they can from a category (e.g., "Name as many animals as you can") in a fixed amount of time (e.g., one minute). Semantic fluency data are richly structured. For example, researchers have found a *frequency* effect: Items that occur more often in the world are also produced more often in fluency lists. Under the assumption that knowledge is represented as a semantic network, we evaluate a number of possible process models of memory retrieval by determining how well they can reproduce the rich structure of human fluency data.

A debate has emerged as to whether a censored random walk over a semantic network (Abbott, Austerweil, & Griffiths, 2015) or a strategic search in a high-dimensional space (Hills, Jones, & Todd, 2012) better describes how the mind retrieves knowledge from memory. Central to the debate has been one property of fluency data: People tend to retrieve items in clusters in a manner consistent with optimal foraging (Hills et al., 2012). This is the tendency to search in memory within a cluster (sub-category) until the gains from further search in that cluster are outweighed by the benefits of switching to a new cluster. Both models can retrieve items in a manner consistent with optimal foraging. One issue in resolving this debate is that the semantic network account has only made computational-level claims and qualitative comparisons to human data. To advance the debate, we explore different possible search processes and compare their retrieval behavior to human retrieval behavior by examining two traditional effects seen in semantic fluency data: *frequency* and *clustering* effects. We test processes that vary in the extent that they search strategically and explore the network.

The semantic fluency task is often scored by a simple count of the number of items named. While healthy controls typically have no trouble generating many items, patients with memory deficits such as Alzheimer's disease or semantic dementia recall fewer items (Troyer, Moscovitch, Winocur, Leach, & Freedman, 1998). In addition, items that are typical of a category are reported more frequently than items that are atypical (Henley, 1969). For instance, cat is more likely to be named by a participant than lynx. This is particularly pronounced in patients with memory deficits, who often only recall items that are frequent in natural statistics (Sailor, Antoine, Diaz, Kuslansky, & Kluger, 2004). Another well-studied property of semantic fluency data is clustering (Troyer, Moscovitch, & Winocur, 1997). Healthy control participants tend to cluster items together in recall. For example, a participant may list cat, dog, and hamster in sequence because all three items belong to a common sub-category: pets.

In this paper, we evaluate how well different search processes on a semantic network reproduce frequency and clustering effects found in human fluency data. To do so, we first describe several search processes. Next, we measure frequency and clustering performance in humans using previously collected semantic fluency data (Zemla, Kenett, Jun, & Austerweil, 2016). Then, we implement several network search procedures on a standardized semantic network and calculate those same measures for comparison.

Memory Retrieval as Search

A model of memory retrieval is a search process over some representation. A search process begins with a cue (e.g., a category label) and uses that cue to locate relevant information (e.g., category members). A crucial component to any search process is the procedure used to navigate the representation (i.e., how the next item in search is determined). Search processes vary in whether they are local or global, or have aspects of both. A local search process will move from the current location to one nearby (in the representation) using information associated with the current location. A global search process may move from the current location to one far away using information encoded across the representation. Search is typically performed in conjunction with some executive process, to determine the relevance of the information encountered. For example, when searching through memory for "animals", it is necessary to recognize that some items in the search path are not animals, and that some animals have already been found. This executive process does not need to be overt; in the search processes described below, it is assumed that search can traverse over items without conscious awareness and without being reported. In addition, search processes vary in whether they are strategic or not. Strategic search processes involve greater use of working memory and executive functioning to determine where to direct search next. Note that these dimensions are meant to help organize search processes and are idealized (e.g., most search processes are not purely local or global).

From behavioral evidence alone, it is difficult to infer properties of human search. For example, clustering in semantic fluency data has been taken as evidence of a strategic search process that leverages local cues (Troyer et al., 1997). To list animals we may start search at a random animal, say *elephant*. At this point, it may be more efficient to limit search to only African animals, as these items are more accessible. Once search exhausts its store of African animals, we switch back to a global search of any animal. This strategic cluster-and-switch process produces clustered fluency data as seen in participants. However, the switch between global and local cues does not need to be strategic if the underlying memory representation is organized in clusters (Abbott et al., 2015). Under this view, even a simple search process can produce clustered data simply by listing items in the order they are encountered in memorythe burden of efficiently retrieving items is built into the representation, rather than the process.

In this paper, we assume semantic memory is best represented as a semantic network and evaluate different search processes on it. Without specifying both a representation and a process, it is not possible to make claims from behavioral results (Abbott et al., 2015). We do so as a first step towards resolving the semantic network vs. space debate.

Search Processes Over a Semantic Network

In this section, we define a semantic network and outline different possible search processes on it. Some processes, such as node degree search (NDS), rely on global cues, selecting the next item independent of the current item. Others, such as cluster-based depth first search (CbDFS) or the censored random walk (CRW), exploit local clusters by constraining search to nearby nodes. Most processes use a mixture of these cues, including three variations of the CRW which implement random jumping (CRW+RJ), strategic jumping (CRW+SJ), and priming (CRW+PV) We also implement a basic spreading activation model (SA), which has been conceptually influential in the history of semantic networks (e.g., Collins & Loftus, 1975). Figure 1 illustrates how the search processes are approximately distributed over the local-global dimension, and Figure 2 depicts a hypothetical semantic network and a possible search outcome for each search process.



Figure 1. Search processes ranked in terms of how far they tend to move on each step. Note that the ordering is given for the parameter values used in the paper and the precise ordering depends on the parameter values (e.g., CRW+RJ is more global-like when the jump probability is close to 1).



Figure 2. (top) A hypothetical network. (bottom) Example observed search paths for each method.

Semantic networks A common way to represent knowledge in memory is using a semantic network. A semantic network consists of nodes and edges. A node encodes a specific item in memory, such as *dog* or *cat*. Two nodes are connected to each other via an edge whenever those two items are semantically similar. For instance, *dog* and *cat* may be connected by an edge because they are both pets, but *dog* and *elephant* are unlikely to be connected.

In this work, we examine a semantic network that is both undirected and unweighted. Undirected means that if there is an edge between *camel* and *horse*, search could go from *camel* to *horse* or *horse* to *camel* (even though people might be more likely to say *horse* after *camel* than vice versa). Unweighted means that all edges imply the same amount of relatedness between two nodes. For example, if the network has one edge connecting horse and pony and another edge connecting horse and camel, the network encodes that a horse is as related to a ponv as it is to a camel. Nodes sharing an edge are called neighbors. Although these assumptions may seem unrealistic, previous work found that a random walk over an undirected and unweighted semantic network captures optimal foraging behavior (Abbott et al., 2015). These distinctions may influence the clustering and frequency properties of semantic fluency data, but using directed and weighted edges increases the complexity of model substantially. So, we use an undirected network without weights and determine whether fluency data can be approximated by different search processes over it.

Node Degree Search (NDS) Node degree search selects nodes with probability proportional to a node's degree (the number of edges connected to a node). This corresponds to the relative frequency each node is visited by an "infinite" length random walk, which is a predictor of phonetic fluency data (Griffiths, Steyvers, & Firl, 2007). Nodes with a large degree have many neighbors, and are typically encountered more frequently than nodes with a small degree. This search process chooses items based on their approximate frequency within the network regardless of the current location. Thus, it is a global and non-strategic process.

Cluster-based Depth First Search (CbDFS) Cluster-based depth-first search is equivalent to traditional depth-first search, except that the primary unit is a node and its neighbors (cluster) rather than a single node. Search begins at a starting node and outputs all of the neighbors of that node (in random order), skipping any node that has already been output. The process then moves to the most recently output node that has new neighbors and outputs those neighbors. Search is local, always emitting the current node's nearest neighbors and traversing one edge at a time.

Censored Random Walk (CRW) A censored random walk (Abbott et al., 2015) begins at a starting node and proceeds to follow a random walk, outputting each node the first time it is traversed (subsequent traversals over the same node are not output, i.e., they are "censored"). It is a local search process because it only depends on the neighborhood of the current node and can only move to a neighbor of the current node. However, because it only outputs the nodes it observes for the first time, items adjacent in its output may be far apart on the network. Nonetheless, it is much more likely to output a sequence of nodes that are close (in network space).

Censored Random Walk with Random Jumps (CRW+RJ) A censored random walk with random jumps is equivalent to a CRW with one key exception: At each step, the walk may jump to another node in the network (possibly one unconnected to the current node) with probability θ_{RJ} . (Goñi et al., 2010). The target node is chosen proportional to the node's degree (number of edges). As such, nodes that have more edges are more probable jump targets. As with the CRW, this search process is partially local, but due to random jumps, it has a global component. The decision to switch between local and global cues is random, and not a strategic decision.

Censored Random Walk with Strategic Jumps (CRW+SJ) A censored random walk with strategic jumps is similar to one with random jumps, except jump points are not chosen at random. Rather, the jumps occur after encountering θ_{SJ} censored nodes. The number of censored nodes is a proxy for time spent without outputting a new item, and is as a metacognitive cue that the current cluster is

exhausted. As with the CRW+RJ, this model will switch between local and global search. However unlike the random jump model, this switch is strategic: the switch is performed when there is evidence that the local cluster has been exploited sufficiently.

Censored Random Walk with Priming Vector (**CRW+PV**) One artifact of collecting multiple fluency lists from the same individual is that they are not independent. This is particularly pronounced when multiple lists are collected during a single session, as in our data set (Zemla et al., 2016). This results in search being biased by transitions made in a previous search (priming effects).

The censored random walk with priming vector attempts to capture this by biasing transitions toward those transitions produced in the previous list. Search is still a random walk, but whenever it reaches a node present in the previous list, with probability θ_{PV} it transitions to the next observed node in the previous list (if such a transition exists) and with probability $1 - \theta_{PV}$ it moves to a random adjacent node. This search is primarily local, and does not have a strategic component.

Spreading Activation (SA) Classic models of semantic networks (e.g. Collins & Loftus, 1975) explain priming effects using spreading activation. Each node has an activation value attached to it. An initial activation value of 1.0 is given to the starting node, with all other nodes given an activation of 0.0. Activation spreads between all nodes through edges, decaying as it propagates through the network with proportion θ_{SA} at each step. At each step, after performing a batch update of all node activation values, the search process chooses a node with probability proportional to its activation value. This node is then assigned an activation of 1.0. Note that we bound all activation values to be between 0.0 and 1.0. As activation begins to spread throughout the network, search quickly resembles global search as every node's activation eventually reaches and stavs at 1.0. Once this happens, the search returns unobserved nodes with uniform probability.

Experiment and Simulation Details

In this section, we describe the previous data used to evaluate the search processes, as well as the simulation and parameter fitting procedures.

Human Data

We use human data from a previously reported experiment (Zemla et al., 2016). In their study, twenty participants were recruited from Amazon's Mechanical Turk. Each participant performed the semantic fluency task three times for three categories (animals, vegetables, and fruit). Participants entered items as they came to mind and hit "Enter" after each item, which notified the participant that the item was recorded and cleared it from the screen. Participants were instructed not to repeat an item within a list, but could repeat items across lists. Categories were pseudorandomized so that no participant received the same category twice in a row and each triad of lists contained each category once. For each list, participants were asked to generate as many items as they could from the category in three minutes (with a visible timer). We only analyze the results for the animal category. The data were cleaned after collection, correcting any spelling mistakes, removing pluralizations, and standardizing synonymous animals.

Simulations

Following previous work (Abbott et al., 2015), we used the University of South Florida (USF) free association data to construct a semantic network (Nelson, McEvoy, & Schreiber, 2004). This network was constructed by pooling the free association data of 149 participants. Given a set of cue words, participants were asked to generate the first word that came to mind. A semantic network was constructed by drawing edges between each cue-response pair. For our simulations, we used only the largest connected component of the animal subset of the USF network. This network contains 160 nodes, 786 edges, and has an average node degree of 4.91.

Simulated fluency data was generated for each participant using every search process. The simulated data were yoked to real participant data in two ways: First, the simulated fluency lists were matched in length to real participant lists. For instance, if a participant generated lists containing 25, 30, and 35 items, a corresponding set of simulated fluency lists would also contain 25, 30, and 35 items. Second, the yoked list always started with the first item of the participant's real list. In some cases, participants generated items that were not in the USF network. For these cases (15% of items), the simulated lists were instead seeded with a close semantic neighbor (as judged by the first author). For example if a participant list started with *beagle* (not in the USF network), the yoked list would start with *dog*.

This seeding process ensures that the lists explore different parts of the USF network when applicable. Moreover, it mimics the strong primacy effects seen in the experimental data: thirteen of twenty participants started at least two lists with the same animal, six of whom started all three lists with the same animal. Removing this constraint is likely to overestimate the extent to which participants are able to generate novel animals from list to list. One hundred yoked data sets (sets of three lists, matched for length) were generated for each participant. Clustering and frequency metrics were calculated as the average across all 100 data sets for each participant.

Parameter fitting

Four of the seven search processes contained one free parameter. The best-fit parameter was found using a grid search which minimized the maximum z-score compared to the human data across all clustering and frequency measures (described below). CRW+RJ, CRW+PV, and SA models searched parameters 0.0 through 1.0 in intervals of 0.05. For CRW+RJ, the best fitting parameter was 0.0 (no jumping)

and so we chose the second best fit for comparison. The best fit parameters were $\theta_{RJ} = .6$, $\theta_{SJ} = 1$, $\theta_{PV} = .75$, and $\theta_{SA} = .25$. CRW, NDS, and CbDFS have no free parameters.

Quantifying Cluster and Frequency Effects

Although there are many possible statistics based on clusters and frequency, we opt for simple, transparent measures to evaluate and compare the above search processes.

Clustering The clustering of fluency data is evaluated with three measures: cluster size, number of cluster switches, and number of cluster types. Cluster size is the average number of items output from a given cluster before switching to a new cluster. The number of cluster switches is the average number of times a participant switches clusters within a list.

Clusters are determined by assigning each animal to different categories as coded by Troyer et al. (1997) and extended by Hills et al. (2012). We have further extended this coding scheme by including any animals in the data that were not in the coding scheme (14% of animals). Because each animal may belong to multiple categories, determining cluster switch points can be done in multiple ways. We used a fluid switch measure, which counts a cluster switch as any two adjacent items that do not share *any* categories.

Our third cluster-based statistic is the number of unique cluster types. This is calculated by counting the total number of categories within a list (counting all categories to which an item belongs). Intuitively, the number of unique clusters appears to measure the same thing as cluster switches—but note that a cluster switch does not imply switching to a *novel* cluster. That is, a participant may switch back and forth between the same two clusters. Nonetheless, participant cluster switches and number of cluster types are highly correlated (r = .74, p < .001).

Frequency We evaluate frequency effects in three ways: the number of unique animals named (unigrams), the number of unique ordered pairs of animals named (bigrams), and the distribution of unigrams. The number of unigrams and bigrams was counted across all three lists. Finally, we calculated the distribution of unigrams in the data: How many items appeared only once, twice, or three times?

Results

Because simulations were performed on an idealized (USF) semantic network, our discussion of the results focus on the relative patterns of fit. Figure 3 shows error bars (standard error of the mean) to help gauge quantitative fit.

Cluster Size Figure 3a depicts cluster sizes for participants and search processes. On average, participants generated clusters with 2.1 items (SD .44). Processes that behaved like global search (node-degree search and spreading activation) strongly underestimated cluster sizes. Cluster-based depth first search also underestimated cluster size, despite using a local search procedure. This is surprising because clusters are close in a semantic network.

Figure 3. Each clustering and frequency measure is shown for the human data (left) and each search process. (a) Average cluster size per list, (b) Average number of cluster switches per list, (c) Average number of cluster types per list, (d) Number of unique unigrams across three lists, (e) Distribution of unigrams across three lists, (f) Number of bigrams across three lists



The censored random walk produced clusters close in size to the actual human data. Variations of the CRW that included a priming vector (CRW+PV) or strategic jumps (CRW+SJ) showed very little difference in cluster size compared to the CRW. However the censored random walk with random jumps (CRW+RJ) produced smaller clusters, as compared to CRW and compared to human data. Abbott et al. (2015) had previously explored censored random walks with and without random jumps and found no discernible difference with respect to optimal foraging. Our results suggest that when cluster size is taken into account, the random jump model fits worse.

Cluster Switches Figure 3b shows the average number of cluster switches for participants and search processes. Participants switched clusters an average of 16.85 times per list (SD 5.71). The pattern of cluster switches mimicked the inverse of cluster size—models that underestimated the cluster size overestimated the number of cluster switches. This is not surprising, as participant cluster size and cluster switches are negatively correlated in the data (r = -.45, p = .045). Again, processes that relied strictly on global (NDS, SA) or local (CbDFS) cues overestimated the number of cluster switches. In contrast, the CRW, CRW+SJ, and CRW+PV all closely resembled human performance. However the CRW+RJ suffered from the inclusion of random jumps, overestimating the number of cluster switches.

Number of Cluster Types Figure 3c shows the number of cluster types for participants and search processes. Participants produced an average of 16.7 cluster types per list (SD 2.7). In contrast to the cluster size and switch data,

counting the average number of cluster types per list produced a different pattern of results. Search processes that behaved more like global search (NDS, SA, CRW+RJ) performed closest to the actual data, while other search processes tended to underestimate the number of cluster types encountered.

Of course, processes that rely on global cues will tend to have more breadth than processes that rely on local cues. It is interesting that our participants were able to generate a breadth of cluster types resembling global search, but switch clusters less often (as expected by local search). This suggests that participants try to exploit local clusters, but that when they do switch, they tend to avoid old clusters.

Number of Unique Items (Unigrams) Figure 3d plots the average number of unique items listed by participants and the search processes. Across three lists, participants generated an average of 54.4 unique items (SD 19.8), though an average of 99.4 token items (SD 29.9). CRW+PV and CbDFS both produced lists containing a similar number of unique items as the human participants. Both of these models do so by limiting exploration in different ways. CbDFS will tend to generate similar fluency lists on successive trials as there is no mechanism to make long-distance transitions within the network. CRW+PV will tend to produce similar lists because it will make the same transitions as it has in previous lists with high probability.

Distribution of Unigrams Figure 3e depicts the distribution of unigrams for participants and the search processes. On average, participants listed 23.2 items once, 17.3 twice, and 13.9 three times (SDs 16.3, 7.1, 6.1). The large number of

items listed twice and three times is indicative of priming effects from earlier lists. CRW+PV was the only search process that produced a similar distribution of unigrams (in particular, the large number of items produced three times). The other search processes strongly overestimate the number of items that appear only once, and strongly underestimate the number of items that appear three times.

Number of Bigrams Figure 3f plots the number of bigrams produced by participants and the search processes. Participants listed 87.1 unique bigrams (SD 28.4). Nearly all models produced a similar number of bigrams as people, except for CRW+PV, which produced too few bigrams. Because CRW+PV often follows the same transitions as in previous lists, fewer unique bigrams are generated.

General Discussion

We explored whether several search processes on a semantic network adequately captured the frequency and clustering effects seen in semantic fluency data. We found that local search processes captured the average cluster sizes and number of clusters, but failed to capture the number of cluster types produced by people. A priming component to the search process is needed to capture the unigram frequency statistics, but this priming component interferes with producing the appropriate number of bigrams. Although none of the search processes captured human behavior on every measure, the censored random walk with priming performed well across many.

More broadly, the censored random walk model (CRW), and its variations that include strategic jumping (CRW+SJ) or recent memory (CRW+PV), captured much of the clustering behavior seen in human data. Search processes that heavily favored only global or local search tended to produce too many cluster switches and underestimate cluster size. Thus, people balance between local and global search. While previous research was not able to discriminate between the CRW model with or without jumps (Abbott et al., 2015), our results suggest that a random jumping model does not capture human performance well.

The search processes were less successful at modeling frequency effects in the data, though CbDFS and CRW+PV produced the best fits. Overall, this suggests that for our experimental procedure—collecting multiple fluency lists from a single participant in a single setting—the censored random walk with priming vector produces the closest fits to human data with respect to clustering and frequency effects. Its performance may be improved by including a strategic component, where it jumps to unvisited nodes or clusters after censoring multiple items in a row. Future research should investigate this and whether the different processes also replicate optimal foraging.

The current work is limited in several ways. It relies on the validity of the USF semantic network, and the assumption of unweighted and undirected edges (De Deyne, Navarro, & Storms, 2013). This network is constructed from an aggregate of participants, and does not reflect the variability across participants. While these may not be unreasonable assumptions, we have established a baseline to compare additional search processes in the future. One possibility is to use crossvalidation: with enough lists from each participant, some lists could be used to estimate an individual's semantic network using U-INVITE (Zemla et al., 2016), while the remaining lists are used to evaluate the search processes.

Acknowledgements

Support for this research was provided by NIH R21AG0534676 and the Office of the VCGRE at UW-Madison with funding from the WARF.

References

- Abbott, J. T., Austerweil, J. L., & Griffiths, T. L. (2015). Random walks on semantic networks can resemble optimal foraging. *Psychological Review*, *122*(3), 558-569. doi:10.1037/a0038693
- Bousfield, W. A., & Sedgewick, C. H. W. (1944). An analysis of sequences of restricted associative responses. *The Journal of General Psychology*, 30(2), 149-165.
- Collins, A. M., & Loftus, E. F. (1975). A spreading-activation theory of semantic processing. *Psychological Review*, 82(6), 407-428.
- De Deyne, S., Navarro, D. J., & Storms, G. (2013). Better explanations of lexical and semantic cognition using networks derived from continued rather than single-word associations. *Behavior Research Methods*, 45(2), 480-498.
- Goñi, J., Martincorena, I., Corominas-Murtra, B., Arrondo, G., Ardanza-Trevijano, S., & Villoslada, P. (2010). Switcherrandom-walks: A cognitive-inspired mechanism for network exploration. *International Journal of Bifurcation and Chaos*, 20(3), 913-922.
- Griffiths, T. L., Steyvers, M., & Firl, A. (2007). Google and the mind: Predicting fluency with PageRank. *Psychological Science*, 18(12), 1069-1076.
- Henley, N. M. (1969). A psychological study of the semantics of animal terms. *Journal of Verbal Learning and Verbal Behavior*, 8(2), 176-184.
- Hills, T. T., Jones, M. N., & Todd, P. M. (2012). Optimal foraging in semantic memory. *Psychological Review*, *119*(2), 431-440. doi:10.1037/a0027373
- Nelson, D. L., McEvoy, C. L., & Schreiber, T. A. (2004). The University of South Florida free association, rhyme, and word fragment norms. *Behavior Research Methods, Instruments, & Computers*, 36(3), 402-407.
- Sailor, K., Antoine, M., Diaz, M., Kuslansky, G., & Kluger, A. (2004). The effects of Alzheimer's disease on item output in verbal fluency tasks. *Neuropsychology*, 18(2), 306-314.
- Troyer, A. K., Moscovitch, M., & Winocur, G. (1997). Clustering and switching as two components of verbal fluency: Evidence from younger and older healthy adults. *Neuropsychology*, 11(1), 138-146.
- Troyer, A. K., Moscovitch, M., Winocur, G., Leach, L., & Freedman, M. (1998). Clustering and switching on verbal fluency tests in Alzheimer's and Parkinson's disease. *Journal of the International Neuropsychological Society*, 4(02), 137-143.
- Zemla, J. C., Kenett, Y. N., Jun, K-S., & Austerweil, J. L. (2016). U-INVITE: Estimating individual semantic networks from fluency data. In *Proceedings of the 38th Annual Meeting of the Cognitive Science Society* (pp. 1907-1912). Austin, TX: Cognitive Science Society.