
What Comes to Mind? Interpretable Dimensions in Embedding Space Predict Human Ad Hoc Category Construction

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Humans rapidly construct ad hoc categories—e.g., “vegetables for painting”—by
2 recruiting task-relevant properties and retrieving items that score highly on them.
3 We test whether this behavior can be predicted directly from off-the-shelf word
4 embeddings. Across 20 composite categories, we fit per-category elastic-net binomial GLMs over fastText dimensions and evaluate on item-mention probabilities.
5 A sparse linear readout predicts human behavior with strong aggregate accuracy
6 ($r = 0.699$ across $N = 3458$ pairs; **Brier** = 0.0049) and is well calibrated (ECE
7 = 0.0198, improving slightly with an intercept-only adjustment). Beyond per-
8 category fits, we frame *leave-one-base-out* (LOBO) as a retrieval transfer test:
9 we learn a single modifier direction from non-target bases and blend it with the
10 held-out base using a mixing weight γ tuned for *average precision*. This yields
11 coherent semantic shifts and small, mixed retrieval gains (**mean $\Delta\mathbf{AP}$** = 0.0051,
12 median 0.0025; 55.0% of 20 categories improved; mean $\Delta\mathbf{P@10}$ = 0.0150). The
13 learned ad hoc axes align with human-rated properties (mRSA up to $R^2 = 0.227$,
14 median $R^2 = 0.157$), supporting interpretability. Overall, simple, interpretable
15 shifts in embedding space capture key regularities in what comes to mind under
16 situational constraints.
17

18 1 Introduction

19 Humans can improvise new categories on the fly to meet immediate goals. Faced with an unusual
20 constraint—e.g., needing “vegetables for painting”—people quickly recruit task-relevant properties
21 (pigment yield, surface texture, shape as a brush) and retrieve items that score highly on those
22 properties[1]. This kind of ad hoc categorization suggests a search over a structured representational
23 space: items that align on contextually relevant dimensions come to mind first.

24 Traditional accounts model such spaces via behavioral feature elicitations and similarity judgments,
25 which are powerful but slow to scale and necessarily incomplete[2]. The feature space that matters
26 for a novel composite category may be high-dimensional, sparsely verbalizable, and only partially
27 accessible to introspection. Meanwhile, contemporary language models (LMs) organize words in
28 dense vector spaces that capture rich relational structure learned from text. If the relational structure
29 in these spaces partially mirrors structure in the world, then it may be possible to predict which items
30 people generate for ad hoc categories directly from embeddings—without hand-crafting features case
31 by case.

32 This paper investigates that possibility. We ask whether interpretable linear paths in off-the-shelf
33 embedding spaces can predict what “comes to mind” across a set of everyday base categories (e.g.,
34 kitchen appliances, vegetables, musical instruments) composed with task modifiers (e.g., “that could
35 fit in your pocket”, “you would use for self-defense”). Concretely, we frame item generation frequen-
36 cies as binomial outcomes and fit per-category L1-regularized GLMs over embedding dimensions,

37 selecting penalty strength via cross-validation and calibrating post-hoc. Across 20 ad-hoc categories,
38 simple models on fastText embeddings capture substantial variance in human item generation proba-
39 bilities (Pearson $r = 0.699$; Brier = 0.0049), with small but reliable gains after calibration (ECE:
40 0.0198 \rightarrow 0.0189).

41 Beyond per-category probes, we also test *transfer*: can a *single* modifier axis learned on non-target
42 bases help retrieve items in a held-out base (LOBO)? We treat this as a retrieval problem and tune
43 the blend weight γ for *average precision* on validation; we then report Δ AP and P@10 as primary
44 LOBO metrics. The learned direction yields coherent geometric shifts and small, heterogeneous
45 improvements (mean Δ AP = 0.0051; median 0.0025; see §4.4). Finally, we assess interpretability
46 via multi-feature RSA, where distances along each ad hoc axis are partially explained by distances in
47 human feature ratings (up to $R^2 = 0.227$, median $R^2 = 0.157$).

48 2 Data

49 2.1 Tasks and responses

50 We analyze 20 ad hoc categories (e.g., “kitchen appliances you would use for self-defense”), generated
51 by participants recruited on Prolific (N=121). Each row is an item–category pair with: count (number
52 of unique participants who generated the item) and n (number of participants shown that category).

53 2.2 Cleaning & normalization

54 We (i) lowercase/trim responses; (ii) merge near-duplicates with a conservative string-similarity
55 procedure; and (iii) compute per-category totals n_c . To avoid $\pm\infty$ logits for 0/1 extremes we evaluate
56 metrics on probabilities rather than logits.

57 2.3 Embeddings

58 Backend: FastText `get_sentence_vector` with L2 norm and global mean-centering.

59 3 Method

60 For category c and item i with embedding $x_i \in \mathbb{R}^D$ and mentions $y_i \sim \text{Binomial}(n_c, p_i)$,

$$\text{logit}(p_i) = \beta_0^{(c)} + x_i^\top \beta^{(c)}.$$

61 We fit, separately for each category c , a binomial GLM on L2-normalized embedding dimensions.
62 All feature columns are z-scored within category. We estimate $(\beta_0^{(c)}, \beta^{(c)})$ with an elastic-net penalty
63 (mixing parameter $\alpha=0.90$; 90% ℓ_1 , 10% ℓ_2), keeping the intercept unpenalized:

$$\min_{\beta_0^{(c)}, \beta^{(c)}} \sum_i \left[-y_i \log p_i - (n_c - y_i) \log(1 - p_i) \right] + \lambda \left(\alpha \|\beta^{(c)}\|_1 + \frac{1-\alpha}{2} \|\beta^{(c)}\|_2^2 \right).$$

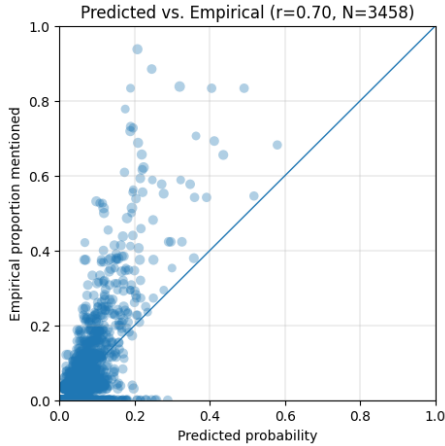
64 **Model selection and stability.** We select λ by *stratified 5-fold CV* within category (stratifying on
65 $y > 0$), with LOOCV as a fallback for very small categories. After choosing λ , we *refit* on all items
66 and report metrics from these refit predictions.

67 **Calibration.** Optionally, we apply a *one-parameter intercept-only* recalibration per category on the
68 refit predictions (minimizing n -weighted Brier). Coefficients $\beta^{(c)}$ remain fixed, preserving sparsity.

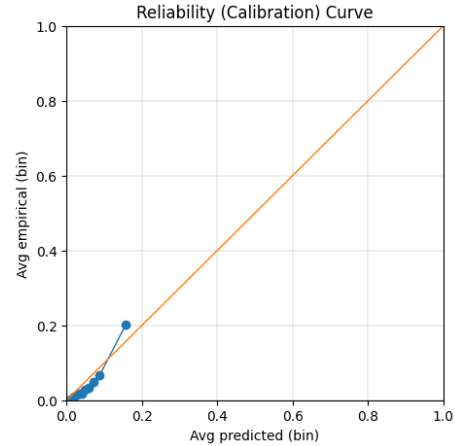
69 4 Results

70 4.1 Overall performance

71 Across all categories and items ($N=3458$ pairs), per-category L1-penalized binomial GLM predicts
72 human probabilities with strong aggregate accuracy ($r = 0.699$, **Brier**= 0.0049, **ECE**= 0.0198;
73 Fig. 1A). A one-parameter, intercept-only recalibration preserves sparsity and yields a small ECE
74 gain without changing overall error (ECE 0.0198 \Rightarrow 0.0189; Fig. 1B).

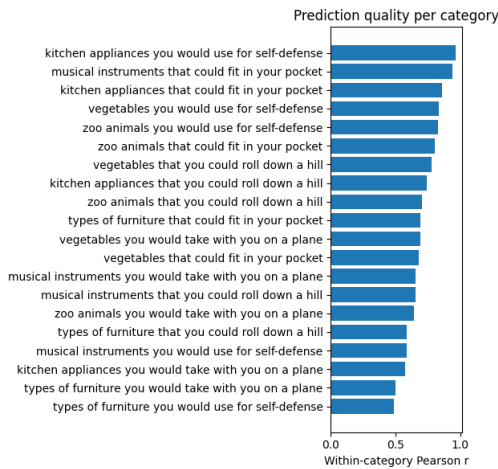


(a) Predicted vs. empirical probabilities across all items ($r = 0.699$).

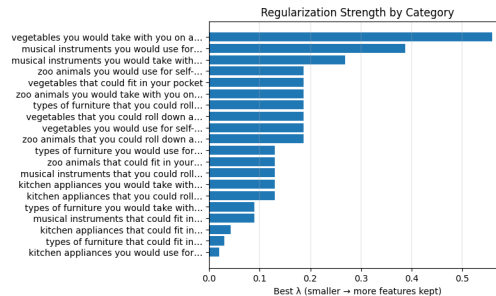


(b) Reliability before/after intercept recalibration (ECE 0.0198 \rightarrow 0.0189).

Figure 1: Overall performance and calibration.



(a) Within-category Pearson r (higher is better).



(b) Selected regularization strength λ per category.

Figure 2: Per-category performance and sparsity.

75 4.2 Within-category analysis

76 Prediction quality is high within individual categories (Fig. 2), with several exceeding $r \approx 0.8$ (e.g.,
 77 self-defense appliances, instruments in your pocket). The same sparse linear readout adapts to diverse
 78 ad hoc constraints by reweighting embedding dimensions per category.

79 4.3 Calibration, Sparsity and Interpretability

80 Reliability curves reveal mild under-confidence in the low-probability regime (≈ 0.05 – 0.20), which
 81 the intercept shift corrects while leaving the slope intact (Fig. 1B). Cross-validated λ values scale
 82 sensibly with category difficulty/sparsity (Appendix panels in Fig. 2B). For each category, only a
 83 small set of embedding axes receives non-zero weight, producing compact, readable profiles.

84 4.4 LOBO as retrieval (AP-tuned)

85 We ask whether a *single* modifier direction, learned from non-target bases, transfers to a held-out base.
 86 We blend the base model with the learned modifier via $p_{\text{enh}} = (1 - \gamma)p_{\text{base}} + \gamma p_{\text{mod}}$ and choose γ
 87 to *maximize average precision* on validation within the held-out base; we then evaluate on the target
 88 items and report ΔAP and P@10 . Geometrically, the learned axis induces coherent shifts that align
 89 with intuitive semantics (e.g., harm potential for *self-defense*; Fig. 3). Quantitatively, improvements

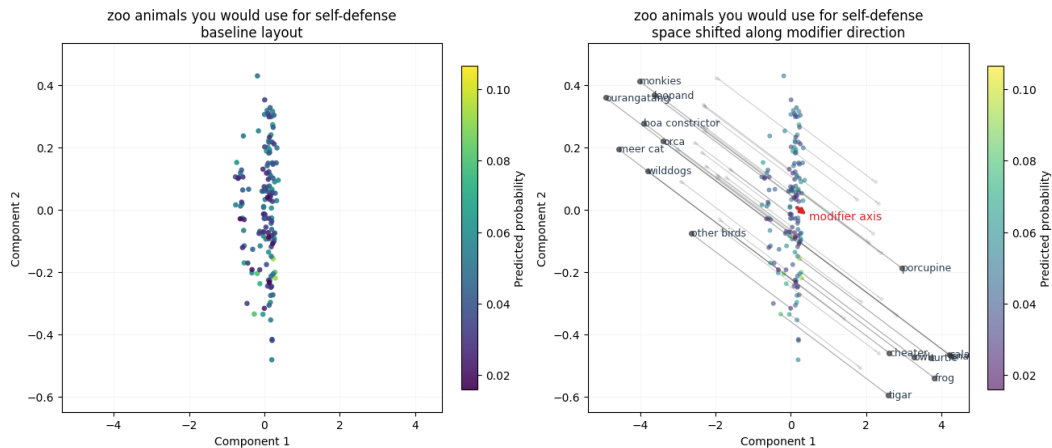


Figure 3: LOBO modifier shift for *zoo animals* \rightarrow *use for self-defense* (single modifier direction learned from other bases; γ tuned for AP).

90 are small and heterogeneous across 20 categories: **mean Δ AP** = 0.0051, median 0.0025 (**55.0%**
 91 improved). P@10 shows a modest average lift (mean Δ P@10 = 0.0150, median 0.0000; **25.0%**
 92 improved). Table A2 lists the largest gains.

93 4.5 Representational structure

94 To test whether the learned 1-D ad hoc axis aligns with human-rated properties, we performed a multi-
 95 feature RSA (mRSA): pairwise distances along the GLM axis were regressed on pairwise distances
 96 in human feature ratings for the corresponding base category, with permutation tests for significance.
 97 Several categories show robust explanatory power (Fig. 4). The associated feature-weight patterns
 98 are intuitive: for “kitchen appliances that could fit in your pocket,” differences in *requires electricity*,
 99 *common*, and *expensive* account for much of the axis structure; for “zoo animals you would use for
 100 self-defense,” *dangerous*, *striking*, and *large* are dominant. Across categories, mRSA explained up to
 101 $R^2 = 0.227$ (median $R^2 = 0.157$).

102 5 Limitations

103 Single-embedding view ignores compositional phrasing and pragmatic constraints; sparsity aids
 104 interpretability but may miss weak multi-axis signals; calibration is mainly in the low- p regime due
 105 to task prevalence; cross-base transfer effects are small on average.

106 6 Conclusion

107 We asked whether a simple, interpretable readout from off-the-shelf text embeddings can approximate
 108 which items people retrieve when composing ad hoc categories. Across 3458 item–category pairs,
 109 per-category elastic-net binomial GLMs achieved strong aggregate accuracy ($r = 0.699$, **Brier**
 110 $= 0.0049$, **ECE** = 0.0198), and a one-parameter intercept shift improved calibration without
 111 changing coefficients. The learned axes are compact and readable, and relate to human properties
 112 (mRSA up to $R^2 = 0.227$, median $R^2 = 0.157$). In a transfer test (LOBO) framed as retrieval and
 113 tuned for AP, we observe coherent semantic shifts and small, mixed gains (**mean Δ AP** = 0.0051;
 114 median 0.0025; mean Δ P@10 = 0.0150). Our results demonstrate that sparse linear readouts
 115 on standard word embeddings can serve as a simple interpretable baseline for ad hoc category
 116 composition, with the modest transfer gains indicating the need for richer composition models.

Table A1: Overall and per-category metrics (mean±sd across categories).

Setting	Pearson r	Brier ↓	ECE ↓
L1-penalized binomial GLM (raw)	0.699	0.0049	0.0198
L1-penalized binomial GLM (+recal)	0.698	0.0049	0.0189

Table A2: LOBO (AP-tuned): top improvements by Δ AP.

Category	AP (base)	AP (enh)	Δ AP	AUC (base)	AUC (enh)	Δ AUC	γ^*
vegetables that you could roll down a hill	0.742	0.761	0.019	0.803	0.795	-0.008	0.50
kitchen appliances you would take with you on a plane	0.791	0.808	0.017	0.813	0.814	0.001	0.35
zoo animals you would use for self-defense	0.643	0.659	0.015	0.695	0.695	-0.000	0.20
types of furniture that could fit in your pocket	0.900	0.912	0.013	0.896	0.895	-0.001	0.45
kitchen appliances that could fit in your pocket	0.794	0.804	0.010	0.824	0.836	0.012	0.10
zoo animals that could fit in your pocket	0.780	0.786	0.006	0.696	0.698	0.001	0.10

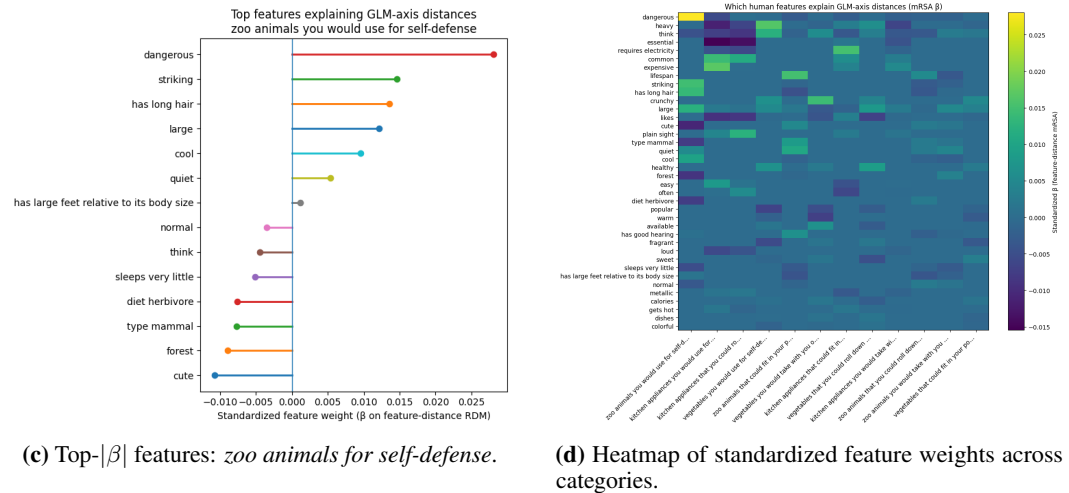
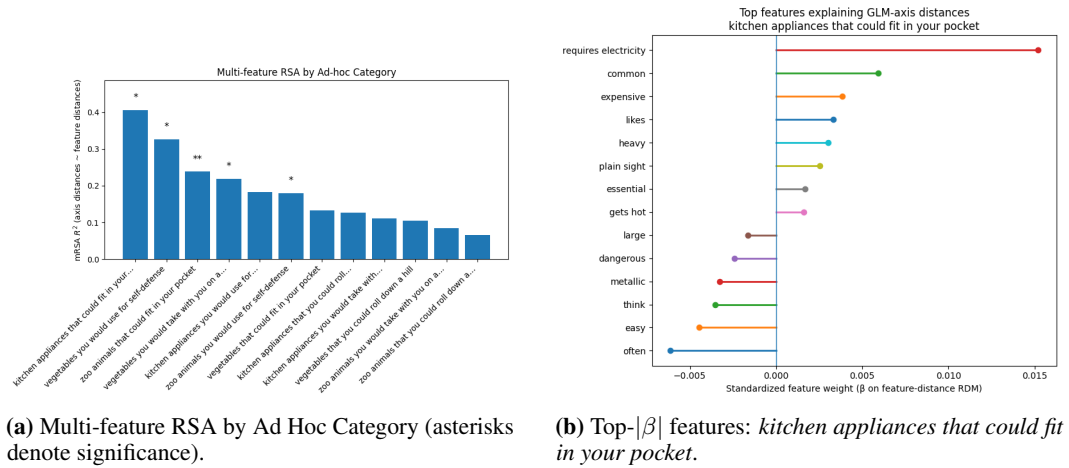


Figure 4: Representational structure and feature explanations.

118 **References**

119 [1] Lawrence W. Barsalou. Ad hoc categories. 11(3):211–227.

120 [2] Tracey Mills and Jonathan Phillips. Locating what comes to mind in empirically derived

121 representational spaces. 240:105549.