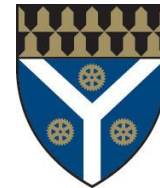




# EasyChinese: Simplifying the Chinese Writing System Using Clustering Algorithms and Ambiguity Metrics

Kai Xu, Department of Computer Science, Yale University  
Advisor: Robert Frank, Department of Linguistics, Yale University



## 1. THE PROBLEM

Chinese is a difficult language to learn to read and write, requiring the memorization of thousands of characters. Our goal is to invent a new way of writing Chinese that is both **easy to read** and **semantically unambiguous**.

We can achieve this by mapping every Chinese character to a combination of **pronunciation** and **semantic category**:

一 (one) → yī (pronunciation) + number (category) → #yī  
五 (five) → wǔ (pronunciation) + number (category) → #wǔ  
衣 (shirt) → yī (pronunciation) + clothing (category) → P<sub>yī</sub>  
裙 (dress) → qún (pronunciation) + clothing (category) → P<sub>qún</sub>

## 2. HOW TO WRITE PINYIN BETTER

Option 1: Pinyin without spaces: ambiguous 25% of the time.

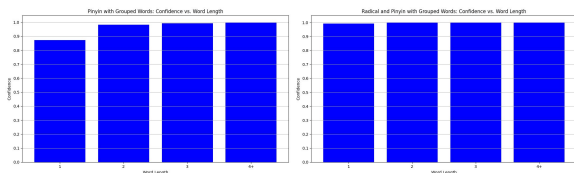
- wōhēnxīhuānbīngqínlín

Option 2: Pinyin with spaces: ambiguous 5% of the time.

- wō hēn xī huān bīng qín lín

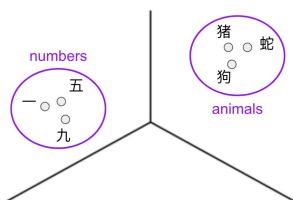
Option 3: **Spaced pinyin with radicals**: ambiguous <1% of the time.

- 戈wō 犭hēn 口xī 欠huān 冫bīng 水qín 氵lín



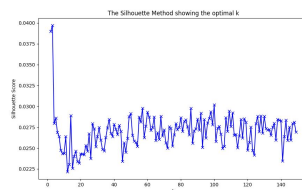
Problem: **radicals do not provide the best semantic categories**. For example, the numbers 1 to 10 are split among seven radicals.

Solution: **apply clustering algorithms to Chinese character embeddings** to create better semantic categories.



## 3. CLUSTERING CONSIDERATIONS

- Reasonable cluster size.** Some clustering algorithms like DBSCAN have no check on cluster size. Others like Gaussian mixture models tend to generate fairly even clusters.
- High intra-cluster similarity and low similarity across different clusters.** We use the **silhouette score** to measure a character's similarity to characters in its own cluster versus those in other clusters.
- Reasonable number of clusters.** We set cluster count to 214 (the number of radicals) for meaningful performance comparisons. However, cluster count can be optimized by plotting the silhouette score of different cluster counts.



Clustering Algorithm	Mean Char. Silhouette	Median Cluster Size	Cluster Size Std. Dev.
Radicals	-0.18	10	65.3
DBSCAN	0.332	2	10.0
Ward hierarchical	0.004	40	29.3
K-means	0.036	40.5	28.0
Gaussian mixture	0.030	42	24.9

## 4. IMPROVING CLUSTERING RESULTS

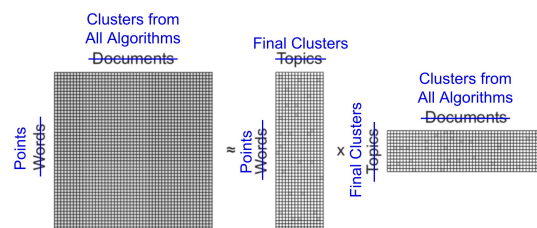
- Removing outliers** (characters that are semantically distant from other characters) before clustering. We developed **geometric ranked outlier score** (measure of outlier status).
- Incorporating human feedback**: we developed an experimental way to revise clustering results based on human feedback. It currently takes very long to converge.
- Consensus clustering**: aggregating the results of multiple clustering algorithms to create better clusters. We discovered a novel method for **conducting consensus clustering with topic modeling algorithms**.

Geometric Ranked Outlier Score									
r = 0		r = 0.5		r = 1		r = 2		r = ∞	
HI-scores	Lo-scores	HI-scores	Lo-scores	HI-scores	Lo-scores	HI-scores	Lo-scores	HI-scores	Lo-scores
猪 0.544	猪 0.821	五 0.549	五 0.858	猪 0.888	猪 0.584	猪 1.065	猪 0.887	猪 1.069	猪 0.889
五 0.539	五 0.821	猪 0.549	猪 0.862	五 0.791	五 0.583	五 1.039	五 0.888	五 1.069	五 0.811
衣 0.532	衣 0.824	猪 0.542	猪 0.864	猪 0.787	猪 0.582	猪 1.027	猪 0.81	猪 1.063	猪 0.814
蛇 0.523	蛇 0.824	猪 0.537	猪 0.865	猪 0.785	猪 0.582	猪 1.017	猪 0.811	猪 1.063	猪 0.815
猪 0.51	猪 0.849	猪 0.534	猪 0.868	猪 0.78	猪 0.581	猪 1.014	猪 0.811	猪 1.04	猪 0.817
猪 0.51	猪 0.849	猪 0.521	猪 0.868	猪 0.777	猪 0.58	猪 1.008	猪 0.812	猪 1.04	猪 0.818
猪 0.506	猪 0.849	猪 0.521	猪 0.872	猪 0.776	猪 0.578	猪 1.007	猪 0.814	猪 1.029	猪 0.819
猪 0.506	猪 0.849	猪 0.514	猪 0.879	猪 0.774	猪 0.578	猪 1.004	猪 0.815	猪 1.026	猪 0.82
猪 0.505	猪 0.853	猪 0.511	猪 0.881	猪 0.772	猪 0.578	猪 1.002	猪 0.816	猪 1.026	猪 0.823
猪 0.501	猪 0.853	猪 0.512	猪 0.885	猪 0.772	猪 0.577	猪 0.998	猪 0.818	猪 1.022	猪 0.827
猪 0.498	猪 0.854	猪 0.511	猪 0.887	猪 0.772	猪 0.577	猪 0.995	猪 0.819	猪 1.021	猪 0.827
猪 0.497	猪 0.854	猪 0.51	猪 0.888	猪 0.771	猪 0.576	猪 0.995	猪 0.82	猪 1.021	猪 0.829
猪 0.496	猪 0.854	猪 0.508	猪 0.889	猪 0.771	猪 0.575	猪 0.994	猪 0.823	猪 1.02	猪 0.829
猪 0.489	猪 0.857	猪 0.507	猪 0.89	猪 0.768	猪 0.575	猪 0.993	猪 0.824	猪 1.019	猪 0.829
猪 0.486	猪 0.857	猪 0.506	猪 0.89	猪 0.765	猪 0.571	猪 0.993	猪 0.824	猪 1.019	猪 0.83
猪 0.486	猪 0.858	猪 0.506	猪 0.89	猪 0.765	猪 0.571	猪 0.992	猪 0.824	猪 1.019	猪 0.831
猪 0.485	猪 0.858	猪 0.502	猪 0.89	猪 0.764	猪 0.567	猪 0.992	猪 0.824	猪 1.015	猪 0.831
猪 0.484	猪 0.859	猪 0.499	猪 0.89	猪 0.763	猪 0.565	猪 0.992	猪 0.825	猪 1.015	猪 0.831
猪 0.484	猪 0.862	猪 0.493	猪 0.89	猪 0.763	猪 0.561	猪 0.991	猪 0.826	猪 1.014	猪 0.832
猪 0.482	猪 0.862	猪 0.491	猪 0.89	猪 0.762	猪 0.561	猪 0.99	猪 0.826	猪 1.011	猪 0.833

Geometric ranked outlier score takes one parameter,  $r$ , that controls the extent to which closer points, as opposed to farther points, influence a point's outlier status.

Original	Iteration 1
Cluster 1: 猫狗	Cluster 1: 猫狗狼狐鹰
Cluster 2: 狮虎	Cluster 2: 狮虎
Cluster 3: 鸡鸭	Cluster 3: 鸡鸭
Iteration 2	Iteration 3
Cluster 1: 猫鹰	Cluster 1: 猫狮虎
Cluster 2: 狮虎	Cluster 2: 狗狼狐
Cluster 3: 狗狼狐鸡鸭	Cluster 3: 鸡鸭鹰

Our experimental method allows humans to specify characters that do not belong in a cluster, and re-clusters accordingly. However, convergence takes a long time.



Consensus clustering can be converted into a topic modeling problem. Instead of assigning words in documents to topics, we assign points in clusters from all clustering algorithms to final clusters. Image adapted from Andrea L. Bertozzi.

## 5. FUTURE RESEARCH

Our discovery of a novel method for conducting consensus clustering using topic modeling algorithms could be useful for broader clustering tasks. Further investigation is pending.