CLUSTERING THE SKETCH: DYNAMIC COMPRESSION FOR EMBEDDING TABLES



Henry Tsang



NORMAL Thomas Ahle

Embedding Tables

We address the challenges associated with large categorical feature tables in Recommendation Systems by co-learning a hash projection matrix and a set of embeddings. The hash projections are learned by occasionally decompressing and performing kmeans clustering on the embeddings, and the embeddings are learned by SGD.

The methods can be used in any Neural Network that uses embedding tables, and is today using hashing or quantization.

Table Compression



The Hashing Trick: Each ID is hashed to one location in a table (here with 1000 rows) and it is assigned the embedding vector stored at the location. Many IDs are likely to share the same vector.



Sum Hash Embeddings Each ID is hashed to two rows, one per table, and its embedding vector is assigned to be the sum of those two vectors. Sometimes a separate table is used for weights





Normal Computing



Compositional Embeddings, CE: Each ID is hashed to a location in each of, say, 4 different tables. The 4 vectors stored there are concatenated into the final embedding vector. Given the large number of possible combinations (here 1000^4), it is unlikely that two IDs get assigned the exact same embedding vector, even if they may share each part with some other IDs.

Clustered Compositional Embeddings, CCE: Combining Compositional Embeddings with Sum Hashing allows us to introduce a clustering step, shown in the big central picture on the poster. Each ID gets assigned a vector that is the concatenation of smaller sums. In each of the four blocks, we apply clustering every epoch, setting the results in the green

tables, and replacing the hash functions in the blue tables with new random values.

Our repository github.com/thomasahle/cce implements more than 11 different table compression algorithms.

 $\times k$

Algorithm 2 Sparse CCE for Least Squares

- **Require:** $X \in \mathbb{R}^{n \times d_1}, Y \in \mathbb{R}^{n \times d_2}, k \in \mathbb{N}$ 1: $H_0 = 0 \in \mathbb{R}^{d_1 \times 2k}$ > Initialize assignments 2: $M_0 = 0 \in \mathbb{R}^{2k \times d_2}$ > Initialize codebook
- 3: for $i = 0, 1, \ldots$ do

4:
$$T_i = H_i M_i \in \mathbb{R}^{d_1 \times d_2}$$

5:
$$A = \operatorname{kmeans}(T_i) \in \{0, 1\}^{d_1 \times k}$$

6:
$$C \sim \text{countsketch}() \in \{-1, 0, 1\}^{d_1}$$

7:
$$H_{i+1} = [A \mid C] \in \mathbb{R}^{d_1 \times 2k}$$

8: $M_{i+1} = \arg \min_M \|XH_{i+1}M - Y\|_F^2$

8:
$$M_{i+1} =$$

9: **end for**

Single iteration of CCE

- 1. Starting from a random embedding table, each ID is hashed to a vector in each of 2 small tables.
- 2. During training, the embedding of an ID is taken to be the mean of the two referenced code words.
- 3. After training for an epoch, the vectors for all (or a sample of) the IDs are computed and clustered. This leaves a new small table in which similar IDs are represented by the same vector.
- 4. We can choose to combine the cluster centers with a new random table (and new hash function), after which the process can be repeated for an increasingly better understanding of which ID should be combined.



thomasahle.github.io/cce



BETTER COMPRESSION BY COMBINING HASHING AND CLUSTERING





Theorem 3.1. Assume $d_1 > k > d_2$ and let $T^* \in \mathbb{R}^{d_1 \times d_2}$ be the matrix that minimizes $\|XT^* - Y\|_F^2$, then $T_i = H_i M_i$ from Algorithm 1 exponentially approaches the optimal loss in the sense

$$E[\|XT_i - Y\|_F^2] \le (1 - \rho)^{ik} \|XT^*\|_F^2 + \|XT^* - Y\|_F^2$$

where $\rho = \|X\|_{-2}^2 / \|X\|_F^2 \approx 1/d_1$ is the smallest singular value of X squared divided by the sum of singular values squared.





