

Multiplying Matrices Without Multiplying (ICML-2021)

D. Blalock (MosaicML, MIT CSAIL) and J. Gutttag (MIT CSAIL)

Jungtaek Kim

jtkim@postech.ac.kr

POSTECH

Pohang 37673, Republic of Korea

<https://jungtaek.github.io>

September 24, 2021

Introduction

- ▶ [Blalock and Gutttag, 2021]
- ▶ Matrix multiplication is among the most fundamental subroutines used in machine learning and scientific computing.
- ▶ As a result, there has been a great deal of work on
 1. implementing high-speed matrix multiplication libraries,
 2. designing custom hardware to accelerate multiplication of certain classes of matrices,
 3. speeding up distributed matrix multiplication,
 4. designing efficient Approximate Matrix Multiplication (AMM) algorithms,under various assumptions.
- ▶ In this paper, the authors focus on the AMM task *under the assumptions that the matrices are tall, relatively dense, and resident in a single machine's memory.*

Introduction

- ▶ This setting arises naturally in machine learning and data mining when one has a data matrix A whose rows are samples and a linear operator B one wishes to apply to these samples.
- ▶ B could be a linear classifier, linear regressor, or an embedding matrix, among other possibilities:

$$\mathbf{t} = AB. \quad (1)$$

- ▶ As a concrete example, consider the task of approximating a softmax classifier trained to predict image labels given embeddings derived from a neural network.
- ▶ Here, the rows of A are the embeddings for each image, and the columns of B are the weight vectors for each class.

Approximating Softmax Classifiers

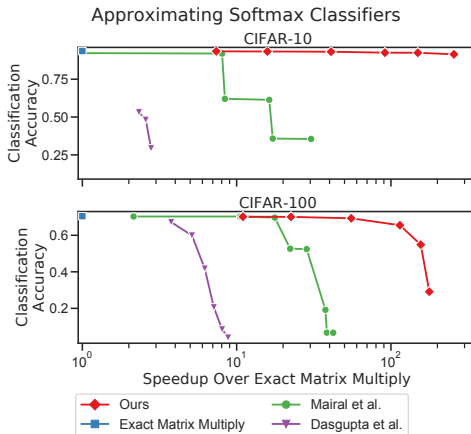


Figure 1: Ours achieves a dramatically better speed-accuracy tradeoff than the best existing methods when approximating two linear classifiers.

Introduction

- ▶ Traditional AMM methods construct matrices $\mathbf{V}_A, \mathbf{V}_B \in \mathbb{R}^{D \times d}$, $d \ll D$ such that

$$\mathbf{AB} \approx (\mathbf{AV}_A)(\mathbf{V}_B^\top \mathbf{B}). \quad (2)$$

- ▶ Often, \mathbf{V}_A and \mathbf{V}_B are sparse, embody some sort of sampling scheme, or have other structure such that these projection operations are faster than a dense matrix multiply.

Multiply-ADDitioN-IESS, a.k.a. MADDNESS

- ▶ It employs a nonlinear preprocessing function and reduces the problem to table lookups.
- ▶ Moreover, in the case that B is known ahead of time, which happens when applying a trained linear model to new data, among other situation, MADDNESS does not require any multiply-add operations.
- ▶ MADDNESS is most closely related to vector quantization methods used for similarity search.
- ▶ However, instead of using an expensive quantization function that requires many multiply-adds, we introduce a family of quantization functions that require no multiply-adds.

Problem Formulation

- ▶ Let $\mathbf{A} \in \mathbb{R}^{N \times D}$ and $\mathbf{B} \in \mathbb{R}^{D \times M}$ be two matrices, with $N \gg D \geq M$.
- ▶ Given a computation time budget τ , this task is to construct three functions $g(\cdot)$, $h(\cdot)$, and $f(\cdot, \cdot)$, along with constants α and β , such that

$$\|\alpha f(g(\mathbf{A}), h(\mathbf{B})) + \beta - \mathbf{AB}\|_F < \varepsilon(\tau) \|\mathbf{AB}\|_F, \quad (3)$$

for as small an error $\varepsilon(\tau)$ possible.

- ▶ The constants α and β are separated from $f(\cdot, \cdot)$ so that $f(\cdot, \cdot)$ can produce low-bitwidth outputs even when the entries of \mathbf{AB} do not fall in this range.
- ▶ We assume the existence of a training set $\tilde{\mathbf{A}}$, whose rows are drawn from the same distribution as the rows of \mathbf{A} .
- ▶ This is a natural assumption in the case that rows of \mathbf{A} represent examples in training data, or structured subsets thereof (such as patches of images).

Background: Product Quantization

- ▶ Product quantization (PQ) [Jégou et al., 2011] is a classic vector quantization algorithm for approximating inner products and Euclidean distances.
- ▶ PQ serves as the basis for nearly all vector quantization methods.
- ▶ The basic intuition behind PQ is that $\mathbf{a}^\top \mathbf{b} \approx \hat{\mathbf{a}}^\top \mathbf{b}$, where $\|\hat{\mathbf{a}} - \mathbf{a}\|$ is small but $\hat{\mathbf{a}}$ has special structure allowing the product to be computed quickly.
- ▶ In somewhat more detail, PQ consists of the following:
 1. Prototype Learning,
 2. Encoding Function $g(\mathbf{a})$,
 3. Table Construction $h(\mathbf{B})$,
 4. Aggregation $f(\cdot, \cdot)$.

Background: Product Quantization

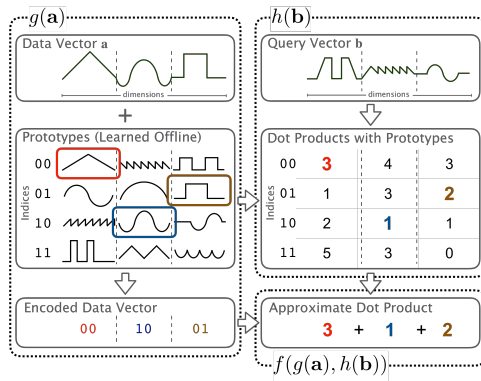


Figure 2: Product Quantization. The $g(\cdot)$ function returns the index of the most similar prototype to the data vector a in each subspace. The $h(\cdot)$ function computes a lookup table of dot products between the query vector b and each prototype in each subspace. The aggregation function $f(\cdot, \cdot)$ sums the table entries corresponding to each index.

Background: Prototype Learning

- ▶ Let $\tilde{\mathbf{A}} \in \mathbb{R}^{N \times D}$ be a training set, K be the number of prototypes per subspace, C be the number of subspaces, and $\{\mathcal{J}^{(c)}\}_{c=1}^C$ be the mutually exclusive and collectively exhaustive sets of indices associated with each subspace.
- ▶ The training-time task of PQ is to learn C sets of prototypes $\mathbf{P}^{(c)} \in \mathbb{R}^{K \times |\mathcal{J}^{(c)}|}$ and assignments $\mathbf{z}^{(c)}$ such that:

$$\sum_{i=1}^N \sum_{c=1}^C \sum_{j \in \mathcal{J}^{(c)}} \left(\tilde{\mathbf{A}}_{ij} - \mathbf{P}_{\mathbf{z}^{(c)}, j}^{(c)} \right)^2, \quad (4)$$

is minimized.

- ▶ It does this by running K -means separately in each subspace $\mathcal{J}^{(c)}$ and using the resulting centroids and assignments to populate $\mathbf{P}^{(c)}$ and $\mathbf{z}^{(c)}$.

Background: Encoding Function

- ▶ Given the learned prototypes, PQ replaces each row \mathbf{a} of \mathbf{A} with the concatenation of its C K -means centroid assignments in each of the C subspaces.
- ▶ Formally,

$$g^{(c)}(\mathbf{a}) = \arg \min_k \sum_{j \in \mathcal{J}^{(c)}} \left(\mathbf{a}_j - \mathbf{P}_{k,j}^{(c)} \right). \quad (5)$$

Background: Table Construction

- Using these same prototypes, PQ constructs a lookup table $h^{(c)}(\mathbf{b}) \in \mathbb{R}^K$ in each of the C subspaces for each column \mathbf{b} of \mathbf{B} , where

$$h^{(c)}(\mathbf{b})_k = \sum_{j \in \mathcal{J}^{(c)}} b_j \mathbf{P}_{k,j}^{(c)}. \quad (6)$$

Background: Aggregation

- Given the encoding of \mathbf{a} and the lookup tables for \mathbf{b} , the product can be approximated as

$$\mathbf{a}^\top \mathbf{b} = \sum_{c=1}^C \mathbf{a}^{(c)\top} \mathbf{b}^{(c)} \approx \sum_{c=1}^C h^{(c)}(\mathbf{b})_k, \quad (7)$$

where $k = g^{(c)}(\mathbf{a})$.

MADNESS

- ▶ PQ and its variants yield a large speedup with $N, M \gg D$.
- ▶ However, we require an algorithm that only needs $N \gg M, D$, a more relaxed scenario common when using linear models and transforms.
- ▶ In this setting, the preprocessing time $g(\mathbf{A})$ can be significant, since ND may be similar to (or even larger than) NM .
- ▶ To address this case, the authors introduce a new $g(\mathbf{A})$ function that yields large speedups even on much smaller matrices.
- ▶ The main idea behind our function is to determine the “most similar” prototype through locality-sensitive hashing [Indyk and Motwani, 1998].

Hash Function Family

- ▶ The family of hash functions in this paper is balanced binary regression trees, with each leaf of the tree acting as one hash bucket.
- ▶ The leaf for a vector x is chosen by traversing the tree from the root and moving to the left child if the value x_j at some index j is below a node-specific threshold v , and to the right child otherwise.
- ▶ To enable the use of SIMD instructions, the tree is limited to 16 leaves and all nodes at a given level of the tree are required to split on the same index j .

Algorithm 1 MADNESSHASH

```
1: Input: vector  $\mathbf{x}$ , split indices  $j^1, \dots, j^4$ , split thresh-  
   olds  $v^1, \dots, v^4$   
2:  $i \leftarrow 1$  // node index within level of tree  
3: for  $t \leftarrow 1$  to 4 do  
4:    $v \leftarrow v_i^t$  // lookup split threshold for node  $i$  at level  $t$   
5:    $b \leftarrow x_{j^t} \geq v ? 1 : 0$  // above split threshold?  
6:    $i \leftarrow 2i - 1 + b$  // assign to left or right child  
7: end for  
8: return  $i$ 
```

Learning the Hash Function Parameters

- ▶ The split indices j^1, j^2, j^3, j^4 and split thresholds v_1, v_2, v_3, v_4 are optimized on the training matrix $\tilde{\mathbf{A}}$ using a greedy tree construction algorithm.
- ▶ To describe this algorithm, we introduce the notion of a bucket \mathcal{B}_i^t , which is the set of vectors mapped to node i in level t of the tree.
- ▶ The root of the tree is level 0 and \mathcal{B}_1^0 contains all the vectors.
- ▶ The sum of squared errors (SSE) loss is used:

$$\mathcal{L}(j, \mathcal{B}) = \sum_{\mathbf{x} \in \mathcal{B}} \left(x_j - \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}' \in \mathcal{B}} x'_j \right)^2, \quad (8)$$

$$\mathcal{L}(\mathcal{B}) = \sum_j \mathcal{L}(j, \mathcal{B}). \quad (9)$$

Adding The Next Level to the Hashing Tree

Algorithm 2 Adding The Next Level to the Hashing Tree

```
1: Input: buckets  $\mathcal{B}_1^{t-1}, \dots, \mathcal{B}_{2^{t-1}}^{t-1}$ , training matrix  $\tilde{A}$   
   // greedily choose next split index and thresholds  
2:  $\hat{\mathcal{J}} \leftarrow \text{heuristic\_select\_idxs}(\mathcal{B}_1^{t-1}, \dots, \mathcal{B}_{2^{t-1}}^{t-1})$   
3:  $l^{\min}, j^{\min}, v^{\min} \leftarrow \infty, \text{NaN}, \text{NaN}$   
4: for  $j \in \hat{\mathcal{J}}$  do  
5:    $l \leftarrow 0$            // initialize loss for this index to 0  
6:    $v \leftarrow []$        // empty list of split thresholds  
7:   for  $i \leftarrow 1$  to  $2^{t-1}$  do  
8:      $v_i, l_i \leftarrow \text{optimal\_split\_threshold}(j, \mathcal{B}_i^{t-1})$   
9:      $\text{append}(v, v_i)$  // append threshold for bucket  $i$   
10:     $l \leftarrow l + l_i$  // accumulate loss from bucket  $i$   
11:   end for  
12:   if  $l < l^{\min}$  then  
13:      $l^{\min} \leftarrow l, j^{\min} \leftarrow j, v^{\min} \leftarrow v$  // new best split  
14:   end if  
15: end for  
   // create new buckets using chosen split  
16:  $\mathcal{B} \leftarrow []$   
17: for  $i \leftarrow 1$  to  $2^{t-1}$  do  
18:    $\mathcal{B}_{\text{below}}, \mathcal{B}_{\text{above}} \leftarrow \text{apply\_split}(v_i^{\min}, \mathcal{B}_i^{t-1})$   
19:    $\text{append}(\mathcal{B}, \mathcal{B}_{\text{below}})$   
20:    $\text{append}(\mathcal{B}, \mathcal{B}_{\text{above}})$   
21: end for  
22: return  $\mathcal{B}, l^{\min}, j^{\min}, v^{\min}$ 
```

Optimizing the Prototypes

- ▶ Instead of suggesting prototypes or table optimizations based on knowledge of B , the authors introduce a method to choose prototypes such that \tilde{A} can be reconstructed from its prototypes with as little squared error as possible.
- ▶ Let $P \in \mathbb{R}^{KC \times D}$ be a matrix whose diagonal blocks of size $K \times |\mathcal{J}^{(c)}|$ consist of the K learned prototypes in each subspace c .
- ▶ The training matrix \tilde{A} can be approximately reconstructed as $\tilde{A} = GP$, where G serves to select the appropriate prototype in each subspace.
- ▶ Rows of G are formed by concatenating the one-hot encoded representations of each assignment for the corresponding row of \tilde{A} .
- ▶ For example if a row were assigned prototypes $\langle 3 \ 1 \ 2 \rangle$ with $K = 4$, $C = 3$, its row in G would be $\langle 0010 \ 1000 \ 0100 \rangle \in \mathbb{R}^{12}$.
- ▶ Finally, it is solved with ridge regression: $P = (G^\top G + \lambda I)^{-1} G^\top \tilde{A}$.

Fast 8-bit Aggregation

- ▶ Let $\mathbf{T} \in \mathbb{R}^{M \times C \times K}$ be the tensor of lookup tables for all M columns of \mathbf{B} .
- ▶ Given the encodings \mathbf{G} , the function $f(\cdot, \cdot)$ is defined as

$$f(g(\mathbf{A}), h(\mathbf{B}))_{n,m} = \sum_{c=1}^C \mathbf{T}_{m,c,k}, \quad (10)$$

where $k = g^{(c)}(\mathbf{a}_n)$.

- ▶ We propose an alternative that sacrifices a small amount of accuracy for a significant increase in speed.
- ▶ Instead of using addition instructions, we use averaging instructions, such as `vpavgb` on x86 or `vrhadd` on ARM.

Theoretical Guarantees

Theorem 4.1 (Generalization Error of MADDNESS). *Let \mathcal{D} be a probability distribution over \mathbb{R}^D and suppose that MADDNESS is trained on a matrix $\tilde{\mathbf{A}} \in \mathbb{R}^{N \times D}$ whose rows are drawn independently from \mathcal{D} and with maximum singular value bounded by σ_A . Let C be the number of codebooks used by MADDNESS and $\lambda > 0$ be the regularization parameter used in the ridge regression step. Then for any $\mathbf{b} \in \mathbb{R}^D$, any $\mathbf{a} \sim \mathcal{D}$, and any $0 < \delta < 1$, we have with probability at least $1 - \delta$ that*

$$\mathbb{E}_{\mathcal{D}}[\mathcal{L}(\mathbf{a}, \mathbf{b})] \leq \mathbb{E}_{\tilde{\mathbf{A}}}[\mathcal{L}(\mathbf{a}, \mathbf{b})] + \frac{C\sigma_A\|\mathbf{b}\|_2}{2\sqrt{\lambda}} \left(\frac{1}{256} + \frac{8 + \sqrt{\nu(C, D, \delta)}}{\sqrt{2n}} \right) \quad (10)$$

where $\mathcal{L}(\mathbf{a}, \mathbf{b}) \triangleq |\mathbf{a}^\top \mathbf{b} - \alpha f(g(\mathbf{a}), h(\mathbf{b})) - \beta|$, α is the scale used for quantizing the lookup tables, β is the constants used in quantizing the lookup tables plus the debiasing constant of Section 4.4, and

$$\nu(C, D, \delta) \triangleq C(4 \lceil \log_2(D) \rceil + 256) \log 2 - \log \delta. \quad (11)$$

How Fast is MADDNESS?

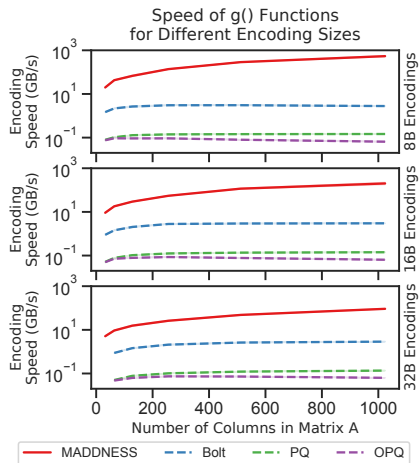


Figure 3: MADDNESS encodes the A matrix orders of magnitude more quickly than existing vector quantization methods.

How Fast is MADDNESS?

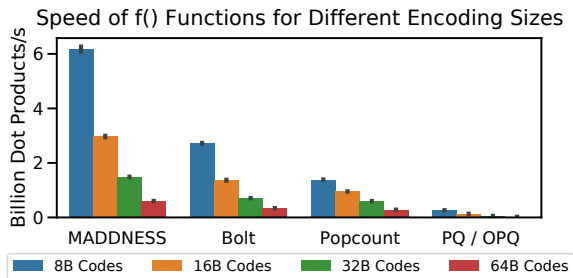


Figure 4: Given the preprocessed matrices, MADDNESS computes the approximate output twice as fast as the fastest existing method.

Softmax Classifier

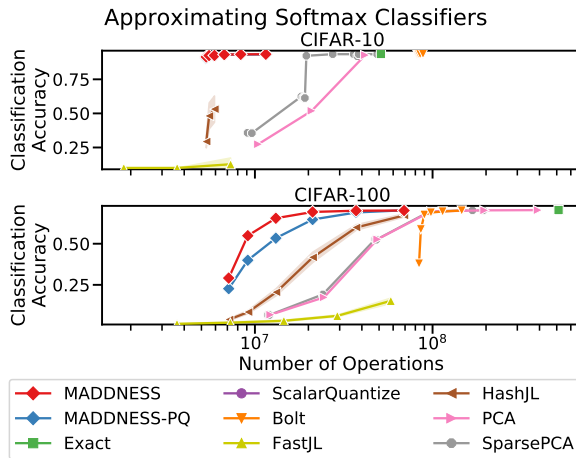


Figure 5: MADDNESS achieves a far better speed- accuracy tradeoff than any existing method when approximating two softmax classifiers.

Kernel-Based Classification

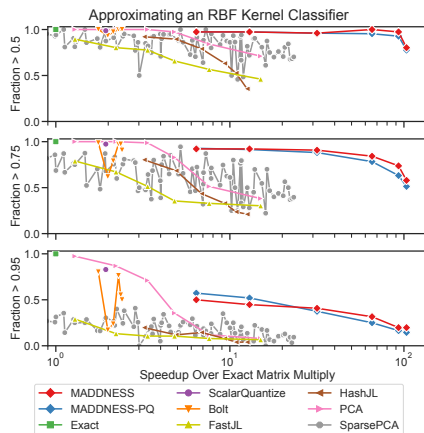


Figure 6: Fraction of UCR datasets for which each method preserves a given fraction of the original accuracy. MADDNESS enables much greater speedups for a given level of accuracy degradation.

Image Filtering

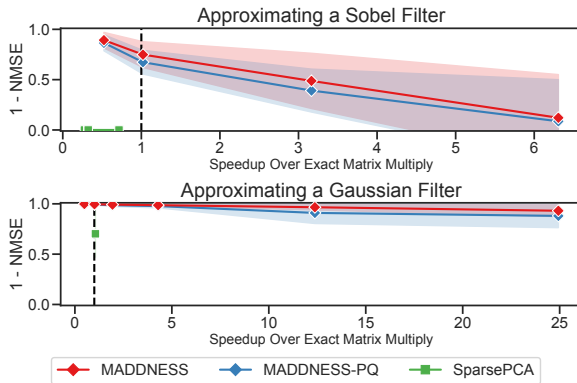


Figure 7: Despite there being only two columns in the matrix B , MADDNESS still achieves a significant speedup with reasonable accuracy. Methods that are Pareto dominated by exact matrix multiplication on both tasks are not shown; this includes all methods but MADDNESS and SparsePCA.

Discussion and Conclusion

- ▶ MADDNESS achieves order-of-magnitude speedups compared to existing AMM methods
- ▶ It achieves up to two-order-of-magnitude speedups compared to the dense baseline.
- ▶ It also compresses matrices by up to three orders of magnitude.
- ▶ These results are evaluated on a CPU, and are obtainable only when there is a training set for one matrix.
- ▶ The authors also claim superior performance only when one matrix is larger than the other, and both matrices are tall—the regime wherein our extremely fast (but less accurate) encoding function is beneficial.

Discussion and Conclusion

- ▶ The authors have not demonstrated results using GPUs or other accelerators.
- ▶ They also have not evaluated a multi-CPU-threaded extension of this algorithm.
- ▶ They have not demonstrated results using convolutional layers in neural networks, or results accelerating full networks.

References I

- D. Blalock and J. Gutttag. Multiplying matrices without multiplying. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 992–1004, Virtual, 2021.
- P. Indyk and R. Motwani. Approximate nearest neighbor: Towards removing the curse of dimensionality. In *Proceedings of the Annual ACM Symposium on Theory of Computing (STOC)*, pages 604–613, 1998.
- H. Jégou, M. Douze, and C. Schmid. Product quantization for nearest neighbor search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128, 2011.