

BlockTensorDecompositions.jl: A Unified Constrained Tensor Decomposition Julia Package

Nicholas J. E. Richardson

Department of Mathematics

Noah Marusenko

Department of Computer Sci

Michael P. Friedlander

ence

Departments of Mathematics

and Computer Science

Table of contents

1	Introduction	2
1.1	Related tools	2
1.2	Contributions	2
2	Tensor Decompositions	3
2.1	Notation	3
2.1.a	Sets	3
2.1.b	Vectors, Matrices, and Tensors	3
2.1.c	Products of Tensors	5
2.1.d	Gradients, Norms, and Lipschitz Constants	8
2.2	Common Decompositions	11
2.2.a	Representing Tucker Decompositions	15
2.2.3	Tensor rank	16
3	Computing Decompositions	17
3.1	Optimization Problem	17
3.2	Base algorithm	18
3.2.a	High level code	18
3.2.b	Computing Gradients	19
3.2.c	Computing Lipschitz Step-sizes	23
4	Computational Techniques	25
4.1	For Improving Convergence Speed	25
4.1.a	Sub-block Descent	26
4.1.b	Momentum	29
4.1.c	Empirical Evidence for Sub-Block Descent and Momentum	31
4.2	For Flexibility	33
4.2.a	Convergence Criteria and Stats	33
4.2.b	BlockUpdate Language	34
4.2.c	Constraints	34
5	Partial Projection and Rescaling	34
6	Multi-scale	35
7	Conclusion	35
8	Appendix	35

8.1 Building the Hessian from two gradients	35
Bibliography	36

1 Introduction

- Tenors are useful in many applications
- Need tools for fast and efficient decompositions

For the scientific user, it would be most useful for there to be a single piece of software that can take as input 1) any reasonable type of factorization model and 2) constraints on the individual factors, and produce a factorization. Details like what rank to select, how the constraints should be enforced, and convergence criteria should be handled automatically, but customizable to the knowledgeable user. These are the core specification for BlockTensorDecompositions.jl.

1.1 Related tools

- Packages within Julia
- Other languages
- Hint at why I developed this

Beyond the external usefulness already mentioned, this package offers a playground for fair comparisons of different parameters and options for performing tensor factorizations across various decomposition models. There exist packages for working with tensors in languages like Python (TensorFlow [1], PyTorch [2], and TensorLy [3]), MATLAB (Tensor Toolbox [4]), R (rTensor [5]), and Julia (TensorKit.jl [6], Tullio.jl [7], OMEinsum.jl [8], and TensorDecompositions.jl [9]). But they only provide a groundwork for basic manipulation of tensors and the most common tensor decomposition models and algorithms, and are not equipped to handle arbitrary user defined constraints and factorization models.

Some progress towards building a unified framework has been made [10–12]. But these approaches don't operate on the high dimensional tensor data natively and rely on matricizations of the problem, or only consider nonnegative constraints. They also don't provide an all-in-one package for executing their frameworks.

1.2 Contributions

- Fast and flexible tensor decomposition package
- Framework for creating and performing custom
 - tensor decompositions
 - constrained factorization (the what)
 - iterative updates (the how)
- Implement new “tricks”
 - a (Lipschitz) matrix step size for efficient sub-block updates
 - multi-scaled factorization when tensor entries are discretizations of a continuous function
 - partial projection and rescaling to enforce linear constraints (rather than Euclidean projection)
 - ?? rank detection ??

The main contribution is a description of a fast and flexible tensor decomposition package, along with a public implementation written in Julia: `BlockTensorDecompositions.jl`. This package provides a framework for creating and performing custom tensor decompositions. To the author's knowledge, it is the first package to provide automatic factorization to a large class of constrained tensor decompositions problems, as well as a framework for implementing new constraints and iterative algorithms. This paper also describes three new techniques not found in the literature that empirically converge faster than traditional block-coordinate descent.

2 Tensor Decompositions

- the math section of the paper

This section reviews the notation used throughout the paper and commonly used tensor decompositions.

2.1 Notation

- tensor notation, use MATLAB notation for indexing so subscripts can be used for a sequence of tensors

2.1.a Sets

The set of real number is denoted as \mathbb{R} and its restrictions to nonnegative numbers is denoted as $\mathbb{R}_+ = \mathbb{R}_{\geq 0} = \{x \in \mathbb{R} \mid x \geq 0\}$.

We use $[N] = \{1, 2, \dots, N\} = \{n\}_{n=1}^N$ to denote integers from 1 to N .

Usually, lower case symbols will be used for the running index, and the capitalized letter will be the maximum letter it runs to. This leads to the convenient shorthand $i \in [I]$, $j \in [J]$, etc.

We use a capital delta Δ to denote sets of vectors or higher order tensors where the slices or fibres along a specified dimension sum to 1, i.e. generalized simplexes.

Usually, we use script letters ($\mathcal{A}, \mathcal{B}, \mathcal{C}$, etc.) for other sets.

2.1.b Vectors, Matrices, and Tensors

Vectors are denoted with lowercase letters (x, y , etc.), and matrices and higher order tensors with uppercase letters (commonly A, B, C and X, Y, Z). The order of a tensor is the number of axes it has. We would call vectors “order-1” or “1st order” tensors, and matrices “order-2” or “2nd order” tensors.

To avoid confusion between entries of a vector/matrix/tensor and indexing a list of objects, we use square brackets to denote the former, and subscripts to denote the later. For example, the entry in the i th row and j th column of a matrix $A \in \mathbb{R}$ is $A[i, j]$. This follows MATLAB/Julia notation where $A[i, j]$ points to the entry $A[i, j]$. We contrast this with a list of I objects being denoted as a_1, \dots, a_I , or more compactly, $\{a_i\}$ when it is clear the index $i \in [I]$.

The transpose $A^\top \in \mathbb{R}^{J \times I}$ of a matrix $A \in \mathbb{R}^{I \times J}$ flips entries along the main diagonal: $A^\top[j, i] = A[i, j]$. In Julia, the transpose of a matrix is typed with a single apostrophe A' .

The n -slices, n th mode slices, or mode n slices of an N th order tensor A are denoted with the slice $A[:, \dots, :, i_n, :, \dots, :]$. For a 3rd order tensor A , the 1st, 2nd, and 3rd mode slices $A[i, :, :]$, $A[:, j, :]$, and $A[:, :, k]$ have special names and are called the horizontal, lateral, and frontal slices and are displayed in Figure 1. In Julia, the 1-, 2-, and 3-slices of a third order array A would be `eachslice(A, dims=1)`, `eachslice(A, dims=2)`, and `eachslice(A, dims=3)`.

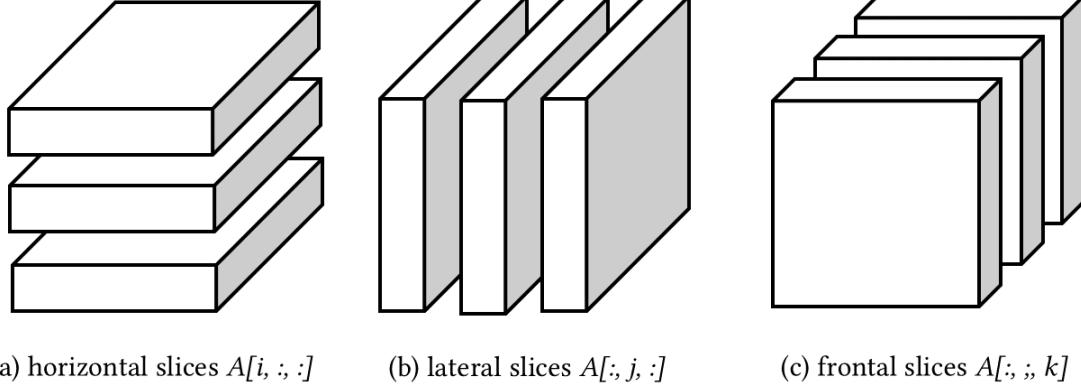


Figure 1: Slices of an order 3 tensor A .

The n -fibres, n th mode fibres, or mode n fibres of an N th order tensor A are denoted $A[i_1, \dots, i_{n-1}, :, i_{n+1}, \dots, i_N]$. For example, the 1-fibres of a matrix M are the column vectors $M[:, j]$, and the 2-fibres are the row vectors $M[i, :]$. For order-3 tensors, the 1st, 2nd, and 3rd mode fibres $A[:, j, k]$, $A[i, :, k]$, and $A[i, j, :]$ are called the vertical/column, horizontal/row, and depth/tube fibres respectively and are displayed in Figure 2. Natively in Julia, the 1-, 2-, and 3-fibres of a third order array A would be `eachslice(A, dims=(2,3))`, `eachslice(A, dims=(1,3))`, and `eachslice(A, dims=(1,2))`. BlockTensorDecomposition.jl defines the function `eachfibre(A; n)` to do exactly this. For example, the 1-fibres of an array A would be `eachfibre(A, n=1)`.

For matrices, the 1-fibres are the same as the 2-slices (and vice versa), but for N th order tensors in general, fibres are always vectors, whereas n -slices are $(N - 1)$ th order tensors.

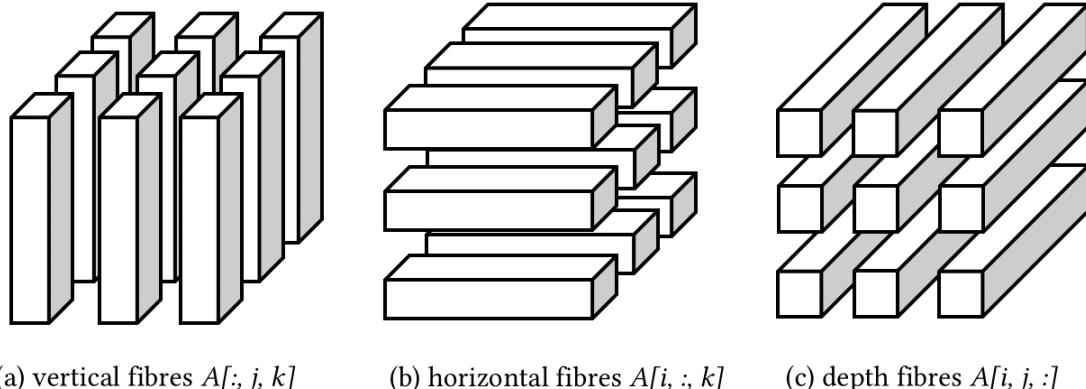


Figure 2: Fibres of an order 3 tensor A .

Since we commonly use I as the size of a tensor's dimension, we use id_I to denote the identity tensor of size I (of the appropriate order). When the order is 2, id_I is an $I \times I$ matrix with ones along the main diagonal, and zeros elsewhere. For higher orders N , this is an $\overbrace{I \times \cdots \times I}^{N \text{ times}}$ tensor where $\text{id}_I[i_1, \dots, i_N] = 1$ when $i_1 = \dots = i_N \in [I]$, and is zero otherwise.

`BlockTensorDecomposition.jl` defines `identity_tensor(I, ndims)` to construct id_I .

2.1.c Products of Tensors

Definition 2.1: The outer product \otimes between two tensors $A \in \mathbb{R}^{I_1 \times \cdots \times I_M}$ and $B \in \mathbb{R}^{J_1 \times \cdots \times J_N}$ yields an order $M + N$ tensor $A \otimes B \in \mathbb{R}^{I_1 \times \cdots \times I_M \times J_1 \times \cdots \times J_N}$ that is entry-wise

$$(A \otimes B)[i_1, \dots, i_M, j_1, \dots, j_N] = A[i_1, \dots, i_M]B[j_1, \dots, j_N].$$

TODO Define in `BlockTensorDecomposition.jl`

The Frobenius inner product between two tensors $A, B \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ yields a real number $A \cdot B \in \mathbb{R}$ and is defined as

$$\langle A, B \rangle = A \cdot B = \sum_{i_1=1}^{I_1} \cdots \sum_{i_N=1}^{I_N} A[i_1, \dots, i_N]B[i_1, \dots, i_N].$$

Julia's standard library package `LinearAlgebra` implements the Frobenius inner product with `dot(A, B)` or `A .· B`.

The n -slice dot product \cdot_n between two tensors $A \in \mathbb{R}^{K_1, \dots, K_{n-1}, I, K_{n+1}, \dots, K_N}$ and $B \in \mathbb{R}^{K_1, \dots, K_{n-1}, J, K_{n+1}, \dots, K_N}$ returns a matrix $(A \cdot_n B) \in \mathbb{R}^{I \times J}$ with entries

$$(A \cdot_n B)[i, j] = \sum_{k_1 \dots k_{n-1} k_{n+1} \dots k_N} A[k_1, \dots, k_{n-1}, i, k_{n+1}, \dots, k_N]B[k_1, \dots, k_{n-1}, j, k_{n+1}, \dots, k_N].$$

This product can also be thought of as taking the dot product $(A \cdot_n B)[i, j] = A_i \cdot B_j$ between all pairs of n th order slices of A and B , which exactly how `BlockTensorDecomposition.jl` defines the operation.

```
function slicewise_dot(A::AbstractArray, B::AbstractArray; dims=1)
    C = zeros(size(A, dims), size(B, dims))
    if A === B # use faster routine if they are the same
        return _slicewise_self_dot!(C, A; dims)
    end

    for (i, A_slice) in enumerate(eachslice(A; dims))
        for (j, B_slice) in enumerate(eachslice(B; dims))
            C[i, j] = A_slice .· B_slice
        end
    end
end
```

```

    end
    return C
end

function _slicewise_self_dot!(C, A; dims=1)
    enumerated_A_slices = enumerate(eachslice(A; dims))
    for (i, Ai_slice) in enumerated_A_slices
        for (j, Aj_slice) in enumerated_A_slices
            if i > j
                continue
            else # only compute the upper triangle entries of C
                C[i, j] = Ai_slice * Aj_slice
            end
        end
    end
    return Symmetric(C) # indexing C[2,1] points to the entry in C[1,2]
end

```

BlockTensorDecomposition.jl defines this operation with `slicewise_dot(A, B, n)`. In the special case where $A = B$, a more efficient method that only computes entries where $i \leq j$ is defined since $A \cdot_n A$ is a symmetric matrix.

The n -slice product of a tensor with itself $X \cdot_n X$ should be thought of as a generalization of the Gram matrix $X^\top X$ since it considers the matrix generated by taking the dot product between every n th mode slice, just like how the Gram matrix considers the dot product between every pair of columns.

The n -mode product \times_n between a tensor $A \in \mathbb{R}^{I_1 \times \dots \times I_N}$ and matrix $B \in \mathbb{R}^{I_n \times J}$, returns a tensor $(A \times_n B) \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N}$ with entries

$$(A \times_n B)[i_1, \dots, i_{n-1}, j, i_{n+1}, \dots, i_N] = \sum_{i_n=1}^{I_n} A[i_1, \dots, i_{n-1}, i_n, i_{n+1}, \dots, i_N] B[i_n, j].$$

BlockTensorDecomposition.jl defines this operation with `nmode_product(A, B, n)`.

```

function nmode_product(A::AbstractArray, B::AbstractMatrix, n::Integer)
    # convert the problem to the mode-1 product
    Aperm = swapdims(A, n)
    Cperm = Aperm ×₁ B
    return swapdims(Cperm, n) # swap back
end

function ×₁( A::AbstractArray, B::AbstractMatrix )
    # Turn the 1-mode product into matrix-matrix multiplication
    sizeA = size(A)
    Amat = reshape(A, sizeA[1], :)

```

```

# Initialize the output tensor
C = zeros(size(B, 1), sizeA[2:end]...)
Cmat = reshape(C, size(B, 1), prod(sizeA[2:end]))

# Perform matrix-matrix multiplication Cmat = B*Amat
mul!(Cmat, B, Amat)

return C # Output entries of Cmat in tensor form
end

function swapdims(A::AbstractArray, a::Integer, b::Integer=1)
    # Construct a permutation where a and b are swapped
    # e.g. [4, 2, 3, 1, 5, 6] when a=4 and b=1
    dims = collect(1:ndims(A))
    dims[a] = b; dims[b] = a
    return permutedims(A, dims)
end

```

! Note

If we were only working with a fixed order of tensors, we could have defined \times_1 entry-wise with Tullio.jl. The function definition `tulliox1` below gives an example for order three tensors.

```

function tulliox1(A::AbstractArray{_,3}, B::AbstractMatrix)
    @tullio C[i, j, k] := A[r, j, k] * B[i, r]
    return C
end

```

But we would need a new definition for each ordered tensor, or use Julia's meta programming to write a method for each order at runtime.

The n -mode product and n -slice product can be thought of as opposites of each other. The n -mode product sums over just the n th dimension of the first tensor, whereas the n -slice product sums over all but the n th dimension.

We can extend the n -mode product to sum over multiple indices between two tensors.

The multi-mode product $\times_{1,\dots,n} = \times_{1:n} = \times_{[n]}$ between a tensor $A \in \mathbb{R}^{I_1 \times \dots \times I_N}$ and tensor $B \in \mathbb{R}^{I_1 \times \dots \times I_n}$, returns a tensor $(A \times_{[n]} B) \in \mathbb{R}^{I_{n+1} \times \dots \times I_N}$ with entries

$$(A \times_{[n]} B)[i_{n+1}, \dots, i_N] = \sum_{i_1, \dots, i_n} A[i_1, \dots, i_n, i_{n+1}, \dots, i_N] B[i_1, \dots, i_n].$$

This product contracts the first n indexes of A with every index of B .

More generally, we can contract any number of indexes such as the last n indexes of A with every index of B with $\times_{N-n+1, \dots, N} = \times_{(N-n+1):N} = \times_{[-n]}$,

$$(A \times_{[-n]} B)[i_1, \dots, i_n] = \sum_{i_{n+1}, \dots, i_N} A[i_1, \dots, i_{N-n}, i_{N-n+1}, \dots, i_N] B[i_{N-n+1}, \dots, i_N],$$

or specific indexes. For example, we would define $(A \times_{1,3,5} B) \in \mathbb{R}^{I_2 \times I_4 \times I_6}$ where $A \in \mathbb{R}^{I_1 \times \dots \times I_6}$ and $B \in \mathbb{R}^{I_1 \times I_3 \times I_5}$ to be

$$(A \times_{1,3,5} B)[i_2, i_4, i_6] = \sum_{i_1, i_3, i_5} A[i_1, i_2, i_3, i_4, i_5, i_6] B[i_1, i_3, i_5].$$

When A a *half-symmetric*¹ tensor of order $2N$

$$A[i_1, \dots, i_N, i_{N+1}, \dots, i_{2N}] = A[i_{N+1}, \dots, i_{2N}, i_1, \dots, i_N], \quad (1)$$

we have

$$A \times_{[N]} B = A \times_{[-N]} B$$

for tensors B of order N .

2.1.d Gradients, Norms, and Lipschitz Constants

Definition 2.2 (Gradient): The gradient $\nabla f : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}^{I_1 \times \dots \times I_N}$ of a (differentiable) function $f : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}$ is defined entry-wise for a tensor $A \in \mathbb{R}^{I_1 \times \dots \times I_N}$ by

$$\nabla f(A)[i_1, \dots, i_N] = \frac{\partial f}{\partial A[i_1, \dots, i_N]}(A).$$

Definition 2.3 (Hessian): The Hessian $\nabla^2 f : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}^{(I_1 \times \dots \times I_N)^2}$ of a second-differentiable function $f : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}$ is the gradient of the gradient and is defined for a tensor $A \in \mathbb{R}^{I_1 \times \dots \times I_N}$ entry-wise by

$$\nabla^2 f(A)[i_1, \dots, i_N, j_1, \dots, j_N] = \frac{\partial^2 f}{\partial A[i_1, \dots, i_N] \partial A[j_1, \dots, j_N]}(A).$$

For a tensor input A of order N , the Hessian tensor $\nabla^2 f(A)$ is of order $2N$.

See Section 8.1 for how this definition can be reproduced by performing two gradients $\nabla^2 f = \nabla(\nabla f)$.

¹For example, the Hessian of a scalar function (see Definition 2.3).

Definition 2.4: The Frobenius norm of a tensor A is the square root of its dot product with itself

$$\|A\|_F = \sqrt{\langle A, A \rangle}.$$

For vectors v , this is equivalent to the (Euclidean) 2-norm

$$\|v\|_F = \|v\|_2 = \sqrt{\langle v, v \rangle}.$$

For matrices M , the $(2 \rightarrow 2)$ operator norm is defined as

$$\|M\|_{\text{op}} = \sup_{\|v\|_2=1} \|Mv\|_2 = \sigma_1(M)$$

where $\sigma_1(M)$ is the largest singular value of M .

For tensors T , the operator-norm is ambiguous since there are multiple ways we can treat tensors as function on other tensors. There is a canonical way to do this for vectors $x \mapsto v^\top x$ and matrices $x \mapsto Mx$, but not tensors. In the case of the Hessian tensor $\nabla^2 f(A) \in \mathbb{R}^{(I_1 \times \cdots \times I_N)^2}$ evaluated at $A \in \mathbb{R}^{I_1 \times \cdots \times I_N}$, it is natural to consider the function $X \mapsto \nabla^2 f(A) \times_{[N]} X$ for $X \in \mathbb{R}^{I_1 \times \cdots \times I_N}$. This gives us our definition of the operator norm on tensors.

Definition 2.5 (Operator Norm): The operator norm of a half-symmetric tensor $A \in \mathbb{R}^{(I_1 \times \cdots \times I_N)^2}$ (Equation 1) is defined as

$$\|A\|_{\text{op}} = \sup_{\|X\|_F=1} \|A \times_{[N]} X\|_F. \quad (2)$$

In Equation 2, $X \in \mathbb{R}^{I_1 \times \cdots \times I_N}$. Note that this definition agrees with the usual operator norm on matrices when $N = 1$.

Theorem 2.1 (Norm of an outer product): Let $T = A_1 \otimes \cdots \otimes A_N \in \mathbb{R}^{(I_1 \times \cdots \times I_N)^2}$ where $A_n \in \mathbb{R}^{I_n \times I_n}$ are symmetric matrices.

Then T is half-symmetric and

$$\|T\|_{\text{op}} = \prod_{n=1}^N \|A_n\|_{\text{op}}.$$

 Warning

According to how the outer product \otimes is defined in Definition 2.1, the product $A_1 \otimes \cdots \otimes A_N$ shown in Theorem 2.1 is really an element of $\mathbb{R}^{I_1 \times I_1 \times \cdots \times I_N \times I_N}$. Note how the indexes are ordered differently than an element of $\mathbb{R}^{(I_1 \times \cdots \times I_N)^2} = \mathbb{R}^{I_1 \times \cdots \times I_N \times I_1 \times \cdots \times I_N}$. Correcting for this with explicit notation becomes cumbersome and would require tensor transposes, a new definition of an outer product, or reordering of indexes in the definition of a half-symmetric tensor. These can have knock-on effects to the definition of the Hessian, multi-mode product, and the operator norm.

To avoid the headache, the equality

$$T = A_1 \otimes \cdots \otimes A_N$$

in Theorem 2.1 should be thought of as the following entry-wise equation

$$T[i_1, \dots, i_N, j_1, \dots, j_N] = A_1[i_1, j_1] \cdots A_N[i_N, j_N]. \quad (3)$$

With the outer product understood as Equation 3, the results of Theorem 2.1 that T is half-symmetric and its operator norm is the product of the operator norms of the constituent matrices is true.

Definition 2.6 (Lipschitz Function): A function $g : \mathbb{R}^{I_1 \times \cdots \times I_N} \rightarrow \mathbb{R}^{I_1 \times \cdots \times I_N}$ is L -Lipschitz when

$$\|g(A) - g(B)\|_F \leq L\|A - B\|_F, \quad \forall A, B \in \mathbb{R}^{I_1 \times \cdots \times I_N}.$$

We call the smallest such L the Lipschitz constant of g .

Definition 2.7 (Smooth Function): A differentiable function $f : \mathbb{R}^{I_1 \times \cdots \times I_N} \rightarrow \mathbb{R}$ is L -smooth when its gradient $g = \nabla f$ is L -Lipschitz.

Theorem 2.2 (Quadratic Smoothness): Let $f : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}$ be a quadratic function

$$f(X) = \frac{1}{2}A(X, X) + B(X) + C$$

of X with bilinear function $A : (\mathbb{R}^{I_1 \times \dots \times I_N})^2 \rightarrow \mathbb{R}$, linear function $B : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}$, and constant $C \in \mathbb{R}^{I_1 \times \dots \times I_N}$.

Then:

1. the Hessian $\nabla^2 f$ is a constant function that evaluates to $\nabla^2 f(X) = D$ at every point X for some $D \in \mathbb{R}^{(I_1 \times \dots \times I_N)^2}$,
2. the tensor D only depends on the bilinear function A (and not on B and C), and
3. the quadratic function f is L -smooth with constant

$$L = \|D\|_{\text{op}}.$$

2.2 Common Decompositions

- Extensions of PCA/ICA/NMF to higher dimensions
- talk about the most popular Tucker, Tucker-n, CP
- other decompositions
 - high order SVD (see Kolda and Bader)
 - HOSVD (see Kolda, Shifted power method for computing tensor eigenpairs)

A tensor decomposition is a factorization of a tensor into multiple (usually smaller) tensors, that can be recombined into the original tensor. To make a common interface for decompositions, we make an abstract subtype of Julia's `AbstractArray`, and subtype `AbstractDecomposition` for our concrete tensor decompositions.

```
abstract type AbstractDecomposition{T, N} <: AbstractArray{T, N} end
```

Computationally, we can think of a generic decomposition as storing factors (A, B, C, \dots) and operations ($\times_a, \times_b, \dots$) for combining them. This is what we do in `BlockTensorDecomposition.jl`.

```
struct GenericDecomposition{T, N} <: AbstractDecomposition{T, N}
    factors::Tuple{Vararg{AbstractArray{T}}}} # e.g. (A, B, C)
    contractions::Tuple{Vararg{Function}}}} # e.g. (×₁, ×₂)
end
# Y = A ×₁ B ×₂ C
array(G::GenericDecomposition) = multifoldl(contractions(G), factors(G))
```

The function `multifoldl` applies the given operations between each factor, from left to right.

```

function multifoldl(ops, args)
    @assert (length(ops) + 1) == length(args)
    x, xs... = args
    for (op, arg) in zip(ops, xs)
        x = op(x, arg)
    end
    return x
end

```

Different types of decompositions define different operations, and different “ranks” of the same decomposition specific the sizes of the factors used.

A commonly used family of decompositions can be derived from the Tucker decomposition.

Definition 2.8: A rank- (R_1, \dots, R_N) Tucker decomposition of a tensor $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$ produces N matrices $A_n \in \mathbb{R}^{I_n \times R_n}$, $n \in [N]$, and core tensor $B \in \mathbb{R}^{R_1 \times \dots \times R_N}$ such that

$$Y[i_1, \dots, i_N] = \sum_{r_1=1}^{R_1} \dots \sum_{r_N=1}^{R_N} A_1[i_1, r_1] \cdots A_r[i_N, r_N] B[r_1, \dots, r_N] \quad (4)$$

entry-wise. More compactly, this decomposition can be written using the n -mode product, or with double brackets

$$Y = B \times_1 A_1 \times_2 \dots \times_N A_N = B \bigtimes_n A_n = \llbracket B; A_1, \dots, A_N \rrbracket. \quad (5)$$

The *Tucker Product* defined by Equation 5 is implemented in BlockTensorDecomposition.jl with `tuckerproduct(B, (A1, ..., AN))` and computes

$$B \bigtimes_n A_n = \llbracket B; A_1, \dots, A_N \rrbracket.$$

It can also optionally “exclude” one of the matrix factors with the call `tuckerproduct(B, (A1, ..., AN); exclude=n)` to compute

$$B \bigtimes_{m \neq n} A_m = \llbracket B; A_1, \dots, A_{n-1}, \text{id}_{R_n}, A_{n+1}, \dots, A_N \rrbracket.$$

```

function tuckerproduct(core, matrices; exclude=nothing)
    N = ndims(core)
    if isnothing(exclude)
        return multifoldl(tucker_contractions(N), (core, matrices...))
    else
        return multifoldl(getnotindex(tucker_contractions(N), exclude), (core,
getnotindex(matrices, exclude)...))
    end
end

```

```

    end
end

tucker_contractions(N) = Tuple((B, A) -> nmode_product(B, A, n) for n in 1:N)

```

Sometimes we write $A_0 = B$ to ease notation, and suggest the “zeroth” factor of the tucker decomposition is the core tensor B . In the special case when $N = 3$, we can visualize Tucker decomposition as multiplying the core tensor by matrices on all three sides as shown in Figure 3.

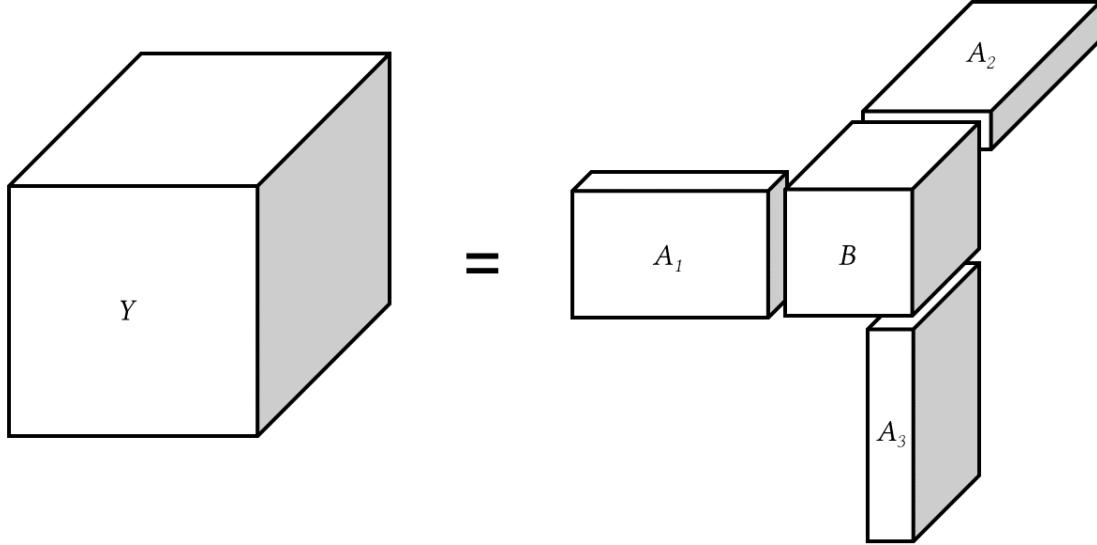


Figure 3: Tucker factorization of a 3rd order tensor Y .

Setting all the matrices of a Tucker decomposition to the identity matrix but the first gives the Tucker-1 decomposition.

Definition 2.9: A rank- R Tucker-1 decomposition of a tensor $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$ produces a matrix $A \in \mathbb{R}^{I_1 \times R}$, and core tensor $B \in \mathbb{R}^{R \times I_2 \times \dots \times I_N}$ such that

$$Y[i_1, \dots, i_N] = \sum_{r=1}^R A[i_1, r] B[r, i_2, \dots, i_N] \quad (6)$$

entry-wise or more compactly,

$$Y = AB = B \times_1 A = \llbracket B; A \rrbracket.$$

Note we extend the usual definition of matrix-matrix multiplication

$$(AB)[i, j] = \sum_{r=1}^R A[i, r] B[r, j]$$

to tensors B in the compact notation for Tucker-1 decomposition $Y = AB$.

More generally, any number of matrices can be set to the identity matrix giving the Tucker- n decomposition.

Definition 2.10: A rank- (R_1, \dots, R_n) Tucker- n decomposition of a tensor $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$ produces n matrices A_1, \dots, A_n , and core tensor $B \in \mathbb{R}^{R_1 \times \dots \times R_n \times I_{n+1} \times \dots \times I_N}$ such that

$$Y[i_1, \dots, i_N] = \sum_{r_1=1}^{R_1} \dots \sum_{r_N=1}^{R_N} A_1[i_1, r_1] \cdots A_n[i_N, r_n] B[r_1, \dots, r_n, i_{n+1}, \dots, i_N] \quad (7)$$

entry-wise, or compactly written in the following three ways,

$$\begin{aligned} Y &= B \times_1 A_1 \times_2 \dots \times_n A_n \times_{n+1} \text{id}_{I_{n+1}} \times_{n+2} \dots \times_N \text{id}_{I_N} \\ Y &= B \times_1 A_1 \times_2 \dots \times_n A_n \\ Y &= [\![B; A_1, \dots, A_n]\!]. \end{aligned}$$

Lastly, if we set the core tensor B to the identity tensor id_R , we obtain the **canonical decomposition/parallel factors model** (CANDECOMP/PARAFAC or CP for short).

Definition 2.11: A rank- R CP decomposition of a tensor $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$ produces N matrices $A_n \in \mathbb{R}^{I_n \times R}$, such that

$$Y[i_1, \dots, i_N] = \sum_{r=1}^R A_1[i_1, r] \cdots A_r[i_N, r] \quad (8)$$

entry-wise. More compactly, this decomposition can be written using the n -mode product, or with double brackets

$$Y = \text{id}_R \times_1 A_1 \times_2 \dots \times_N A_N = \text{id}_R \bigtimes_n A_n = [\![A_1, \dots, A_N]\!].$$

Note CP decomposition is sometimes referred to as Kruskal decomposition, and requires the core only be diagonal (and not necessarily identity) and the factors A_n have normalized columns $\|A_n[:, r]\|_2 = 1$.

Other factorization models are used that combine aspects of CP and Tucker decomposition [13], are specialized for order 3 tensors [14, 15], or provide alternate decomposition models entirely like tensor-trains [16]. But the (full) Tucker, and its special cases Tucker- n , and CP decomposition are most commonly used extensions of the low-rank matrix factorization to tensors. These factorizations are summarized in Table 1.

Table 1: Summary of common tensor factorizations. Here, N is the order of the factorized tensor.

Name	Bracket Notation	n -mode Product	Entry-wise
Tucker	$\llbracket A_0; A_1, \dots, A_N \rrbracket$	$A_0 \times_1 A_1 \times_2 \dots \times_N A_N$	Equation 4
Tucker-1	$\llbracket A_0; A_1 \rrbracket$	$A_0 \times_1 A_1$	Equation 6
Tucker- n	$\llbracket A_0; A_1, \dots, A_n \rrbracket$	$A_0 \times_1 A_1 \times_2 \dots \times_n A_n$	Equation 7
CP	$\llbracket A_1, \dots, A_N \rrbracket$	$\text{id}_R \times_1 A_1 \times_2 \dots \times_N A_N$	Equation 8

TODO add discussion on other decompositions - high order SVD (see Kolda and Bader) - HOSVD (see Kolda, Shifted power method for computing tensor eigenpairs)

Tensor decompositions are not necessarily unique. It should be clear that scaling one factor by $x \neq 0$ and dividing another by x yields the same original tensor. Furthermore, fibres and slices can be permuted without affecting the the original tensor. Up to these manipulations, for a fixed rank, there exist criteria that ensures their decompositions are unique [13, 17, 18].

2.2.a Representing Tucker Decompositions

There are implemented in `BlockTensorDecomposition.jl` and can be called, for a third order tensor, with `Tucker((B, A1, A2, A3))`, `Tucker1((B, A1))`, and `CPDecomposition((A1, A2, A3))`. These Julia structs store the tensor in its factored form. We could define the contractions for these types and use the common interface provided by `array`, but it turns out we can reconstruct the whole tensor more efficiently. If the recombined tensor or particular entries are requested, Julia dispatches on the type of decomposition and calls a particular method of `array` or `getindex`. The implementations for efficient array construction and index access are provided below.

```
array(T::Tucker) = multifoldl(tucker_contractions(ndims(T)), factors(T))
tucker_contractions(N) = Tuple((G, A) -> nmode_product(G, A, n) for n in 1:N)
```

TODO add `getindex` method for Tucker type

```
function array(T::Tucker1)
    B, A = factors(T)
    return B ×1 A
end

function getindex(T::Tucker1, I::Vararg{Int})
    B, A = factors(T)
    i, J... = I # (i, J) = (I[1], I[begin+1:end])
    return (@view A[i, :]) ∙ view(B, :, J...)
end
```

```
array(CPD::CPDecomposition) =
    mapreduce(vector_outer, +, zip((eachcol.(factors(CPD))))...))
```

```

vector_outer(v) = reshape(kron(reverse(v)...),length.(v))

getindex(CPD::CPD{Vararg{Int}}) =
    sum(reduce(.*, (@view f[i,:]) for (f,i) in zip(factors(CPD), I)))

```

2.3 Tensor rank

- tensor rank
- constrained rank (nonnegative etc.)

The rank of a matrix $Y \in \mathbb{R}^{I \times J}$ can be defined as the smallest $R \in \mathbb{Z}_+$ such that there exists an exact factorization $Y = AB$ for some $A \in \mathbb{R}^{I \times R}$ and $B \in \mathbb{R}^{R \times J}$.

Although this can be extended to higher order tensors, we must specify under which factorization model we are using. For example, the *CP-rank* R of a tensor Y is the smallest such R that omits an exact CP decomposition of Y .

Definition 2.12: The CP rank of a tensor $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$ is the smallest R such that there exist factors $A_n \in \mathbb{R}^{I_n \times R}$ and $Y = [A_1, \dots, A_N]$,

$$\text{rank}_{\text{CP}}(Y) = \min\{R \mid \exists A_n \in \mathbb{R}^{I_n \times R}, n \in [N] \quad \text{s.t.} \quad Y = [A_1, \dots, A_N]\}.$$

In a similar way, we can define the *Tucker-1-rank* R .

Definition 2.13: The Tucker-1 rank of a tensor $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$ is the smallest R such that there exist factors $A \in \mathbb{R}^{I_1 \times R}$ and $B \in \mathbb{R}^{R \times I_2 \times \dots \times I_N}$ where $Y = AB$

$$\text{rank}_{\text{Tucker-1}}(Y) = \min\{R \mid \exists A \in \mathbb{R}^{I_1 \times R}, B \in \mathbb{R}^{R \times I_2 \times \dots \times I_N} \quad \text{s.t.} \quad Y = AB\}$$

For the Tucker and Tucker- n decompositions, we instead call a particular factorization **a** rank- (R_1, \dots, R_N) Tucker factorization or **a** rank- (R_1, \dots, R_n) Tucker- n factorization, rather than **the** CP- or Tucker-1-rank of a tensor or **the** rank of a matrix.

One reason CP and Tucker-1 only need a single rank R can be explained by considering the case when the order of the tensor $N = 2$ (matrices). The two factorizations become equivalent and are equal to low-rank R matrix factorization $Y = AB$. In fact, Tucker-1 is always equivalent to a low-rank matrix factorization, if you consider a flattening of the tensor to arrange the entries as a matrix.

The idea of tensor rank can be generalized further to constrained rank. These are the smallest rank R such that the factors in the decomposition obey the given set of constraints.

For example, the nonnegative Tucker-1 rank is defined as

$$\text{rank}_{\text{Tucker-1}}^+(Y) = \min \left\{ R \mid \exists A_n \in \mathbb{R}_+^{I_n \times R}, B \in \mathbb{R}_+^{R \times I_2 \times \dots \times I_N} \quad \text{s.t.} \quad Y = AB \right\}.$$

More restrictive constraints increase the rank of the tensor since there is less freedom in selecting the factors.

Most tensor decomposition algorithms require the rank as input [CITE] since calculating the rank of the tensor can be NP-hard in general [19]. For applications where the rank is not known a priori, a common strategy is to attempt a decomposition for a variety of ranks, and select the model with smallest rank that still achieves good fit between the factorization and the original tensor.

3 Computing Decompositions

- Given a data tensor and a model, how do we fit the model?

Many tensor decompositions algorithms exist in the literature. Usually, they cyclically (or in a random order) update factors until their reconstruction satisfies some convergence criterion. The base algorithm described in Section 3.2 provides flexible framework for wide class of constrained tensor factorization problems. This framework was selected based on empirical observations where it outperforms other similar algorithms, and has also been observed in the literature [10].

3.1 Optimization Problem

- Least squares (can use KL, 1 norm, etc.)

Ideally, we would be given a data tensor Y and decomposition model, and compute an exact factorization of Y into its factors. Because there is often measurement, numerical, or modeling error, an exact factorization of Y for a particular rank may not exist. To over come this, we instead try to fit the model to the data. Let X be the reconstruction of factors A_1, \dots, A_N according to some decomposition for a fixed rank. We assume we know the size of the factors A_1, \dots, A_N and how they are combined to produce a tensor the same size of Y , i.e. the map $g : (A_1, \dots, A_N) \mapsto X$.

There are many loss functions that can be used to determine how close the model X is to the data Y . In principle, any distance or divergence $d(Y, X)$ could be used. We use the L_2 loss or least-squares distance between the tensors $\|X - Y\|_F^2$, but other losses are used for tensor decomposition in practice such as the KL divergence [CITE].

The main optimization we must solve is now given.

Definition 3.1: The constrained least-squares tensor factorization problem is to solve

$$\min_{A_1, \dots, A_N} \frac{1}{2} \|g(A_1, \dots, A_N) - Y\|_F^2 \quad \text{s.t.} \quad (A_1, \dots, A_N) \in \mathcal{C}_1 \times \dots \times \mathcal{C}_N \quad (9)$$

for a given data tensor Y , constraints $\mathcal{C}_1, \dots, \mathcal{C}_N$, and decomposition model g with fixed rank.

Note the problem would have the same solutions as simply using the objective $\|g(A_1, \dots, A_N) - Y\|$ without squaring and dividing by 2. We define the objective in Equation 9 to make computing the function value and gradients faster.

3.2 Base algorithm

- Use Block Coordinate Descent / Alternating Proximal Descent
 - do *not* use alternating least squares (slower for unconstrained problems, no closed form update for general constrained problems)

Let $f(A_1, \dots, A_N) := \frac{1}{2} \|g(A_1, \dots, A_N) - Y\|_F^2$ be the objective function we wish to minimize in Equation 9. Following Xu and Yin [10], the general approach we take to minimize f is to apply block coordinate descent using each factor as a different block. Let A_n^t be the t th iteration of the n th factor, and let

$$f_n^t(A_n) := \frac{1}{2} \|g(A_1^{t+1}, \dots, A_{n-1}^{t+1}, A_n, A_{n+1}^t, \dots, A_N^t) - Y\|_F^2$$

be the (partially updated) objective function at iteration t for factor n .

Given initial factors A_1^0, \dots, A_N^0 , we cycle through the factors $n \in [N]$ and perform the update

$$A_n^{t+1} \leftarrow \arg \min_{A_n \in \mathcal{C}_n} \langle \nabla f_n^t(A_n^t), A_n - A_n^t \rangle + \frac{L_n^t}{2} \|A_n - A_n^t\|_F^2,$$

for $t = 1, 2, \dots$ until some convergence criterion is satisfied (see Section 4.2.a).

This implicit update has the *projected gradient descent* closed form solution for convex constraints \mathcal{C}_n ,

$$A_n^{t+1} \leftarrow P_{\mathcal{C}_n} \left(A_n^t - \frac{1}{L_n^t} \nabla f_n^t(A_n^t) \right). \quad (10)$$

We typically choose L_n^t to be the Lipschitz constant of ∇f_n^t , since it is a sufficient condition to guarantee $f_n^t(A_n^{t+1}) \leq f_n^t(A_n^t)$, but other step sizes can be used in theory [20 (Sec. 1.2.3)].

?ASIDE? To write ∇f_n^t , we have assumed (block) differentiability of the decomposition model g . In practice, most decompositions are “block-linear” (freeze all factors but one and you have a linear function) and in rare cases are “block-affine”. “block-affine” is enough to ensure f_n^t is convex (i.e. f is “block-convex”) so the updates Equation 10 converge to a Nash equilibrium (block minimizer).

3.2.a High level code

To ensure the code stays flexible, the main algorithm of `BlockTensorDecomposition.jl`, `factorize`, is defined at a very high level.

```
factorize(Y; kwargs...) =
    _factorize(Y; (default_kwargs(Y; kwargs...))...)
```

```

"""
Inner level function once keyword arguments are set
"""

function _factorize(Y; kwargs...)
    decomposition, previous, updateprevious!, parameters, updateparameters!, update!,
    stats_data, getstats, converged, kwargs = initialize(Y, kwargs)

    while !converged(stats_data; kwargs...)
        # Usually one cycle of updates through each factor in the decomposition
        update!(decomposition; parameters...)

        # This could be the next stepsize or other info used by update!
        updateparameters!(parameters, decomposition, previous)

        push!(stats_data,
              getstats(decomposition, Y, previous, parameters, stats_data))

        # Update one or two previous iterates. For example, used for momentum
        updateprevious!(previous, parameters, decomposition)
    end

    kwargs = postprocess!(decomposition, Y, previous, parameters, stats_data,
                          updateparameters!, getstats, kwargs)

    return decomposition, stats_data, kwargs
end

```

The magic of the code is in defining the functions at runtime for a particular decomposition requested, from a reasonable set of default keyword arguments. This is discussed further in Section 4.2.

3.2.b Computing Gradients

- Use Auto diff generally
- But hand-crafted gradients and Lipschitz calculations *can* be faster (e.g. symmetrized slice-wise dot product)

Generally, we can use automatic differentiation on f to compute gradients. Some care needs to be taken otherwise the forward or backwards pass will have to be recompiled every iteration since the factors are updated every iteration.

But for Tucker decompositions, we can compute gradients faster than what an automatic differentiation scheme would give, by taking advantage of symmetry and other computational shortcuts.

Starting with the Tucker-1 decomposition (Definition 2.9), we would like to compute $\nabla_B f(B, A)$ and $\nabla_A f(B, A)$ for $f(B, A) = \frac{1}{2} \|AB - Y\|_F^2$ for a given input Y . We have the gradient

$$\nabla_B f(B, A) = A^\top (AB - Y) = (B \times_1 A - Y) \times_1 A^\top \quad (11)$$

by chain rule, but it is more efficient to calculate the gradient as

$$\nabla_B f(B, A) = (A^\top A)B - A^\top Y = B \times_1 (A^\top A) - Y \times_1 A^\top. \quad (12)$$

²For $A \in \mathbb{R}^{I \times R}$, $B \in \mathbb{R}^{R \times J \times K}$, and $Y \in \mathbb{R}^{I \times J \times K}$, Equation 11 requires

$$\underbrace{2IJKR}_{AB-Y} + \underbrace{IJK(2I-1)}_{A^\top(AB-Y)} \sim 2IJKR + 2I^2JK$$

floating point operations (FLOPS) whereas Equation 12 only uses

$$\underbrace{\frac{R(R+1)}{2}(2I-1)}_{A^\top A} + \underbrace{RJK(2I-1)}_{A^\top Y} + \underbrace{2R^2JK}_{(A^\top A)B-(A^\top Y)} \sim 2IJKR + 2R^2JK + IR^2$$

FLOPS³. So for small ranks $R \ll I$, Equation 12 is cheaper.

A similar story can be said about $\nabla_A f(B, A)$ which is most efficiently computed as

$$\nabla_A f(B, A) = A(B \cdot_1 B) - Y \cdot_1 B.$$

! Note

For the family of Tucker decompositions, the objective function f is “block-quadratic” with respect to the factors. This means the gradient with respect to a factor is an affine function of that factor. This is exactly what we see in Equation 12 where B is multiplied by the “slope” $A^\top A$ plus a shift of $-Y \times_1 A^\top$.

The associated implementation with `BlockTensorDecomposition.jl` is shown below. We define a `make_gradient` which takes the decomposition, factor index n , and data tensor Y , and creates a function that computes the gradient for the same type of decomposition. This lets us manipulate the function that computes the gradient, rather than just the computed gradient.

```
function make_gradient(T::Tucker1, n::Integer, Y::AbstractArray; objective::L2,
kwargs...)
    if n==0 # the core is the zeroth factor
        function gradient0(X::Tucker1; kwargs...)
            (B, A) = factors(X)
            AA = A'A
            YA = Y×_1 A
            grad = B×_1 AA - YA
            return grad
        end
    end
end
```

²Seeing Equation 11 and Equation 12 written using the 1-mode product shows how it is “backwards” to normal matrix-matrix multiplication.

³Note we have the smaller factor $R(R+1)/2$ and not the expected R^2 number of entries needed to compute $A^\top A$. The product is a symmetric matrix so only the upper or lower triangle of entries needs to be computed.

```

    return gradient0
elseif n==1 # the matrix is the first factor
    function gradient1(X::Tucker1; kwargs...)
        (B, A) = factors(X)
        BB = slicewise_dot(B, B)
        YB = slicewise_dot(Y, B)
        grad = A*BB - YB
        return grad
    end
    return gradient1
else
    error("No $(n)th factor in Tucker1")
end
end

```

Similarly, we also have special methods for the Tucker and CP Decomposition.

The gradient with respect to the core for a full Tucker factorization is

$$\nabla_B f(B, A_1, \dots, A_N) = B \bigtimes_n A_n^\top A_n - Y \bigtimes_n A_n^\top,$$

and the gradient with respect to the matrix factor A_n is

$$\nabla_{A_n} f(B, A_1, \dots, A_N) = A_n (\tilde{X}_n \cdot_n \tilde{X}_n) - Y \cdot_n \tilde{X}_n$$

where

$$\tilde{X}_n = \left(B \bigtimes_{m \neq n} A_m \right) = [B; A_1, \dots, A_{n-1}, \text{id}_{R_n}, A_{n+1}, \dots, A_N].$$

```

function make_gradient(T::Tucker, n::Integer, Y::AbstractArray; objective::L2,
kwargs...)
    N = ndims(T)
    if n==0 # the core is the zeroth factor
        function gradient_core(X::AbstractTucker; kwargs...)
            B = core(X)
            matrices = matrix_factors(X)
            gram_matrices = map(A -> A'A, matrices) # gram matrices AA = A'A,
                                                # BB = B'B, ...
            grad = tuckerproduct(B, gram_matrices)
            - tuckerproduct(Y, adjoint.(matrices))
            return grad
        end
        return gradient_core
    elseif n in 1:N # the matrix factors start at m=1
        function gradient_matrix(X::AbstractTucker; kwargs...)

```

```

        B = core(X)
        matrices = matrix_factors(X)
        A_n = factor(X, n)
        X_n = tuckerproduct(B, matrices; exclude=n)
        grad = A_n * slicewise_dot(X_n, X_n; dims=n)
            - slicewise_dot(Y, X_n; dims=n)
        return grad
    end
    return gradient_matrix

else
    error("No $(n)th factor in Tucker")
end
end

```

For the CP Decomposition, we can simply treat the core as $B = \text{id}_R$ and compute the gradient with respect to the matrix factors similarly to the Tucker decomposition:

$$\nabla_{A_n} f(A_1, \dots, A_N) = A_n (\tilde{X}_n \cdot_n \tilde{X}_n) - Y \cdot_n \tilde{X}_n$$

where

$$\tilde{X}_n = \left(\text{id}_R \bigtimes_{m \neq n} A_m \right) = [\text{id}_R; A_1, \dots, A_{n-1}, \text{id}_R, A_{n+1}, \dots, A_N].$$

```

function make_gradient(T::CPDecomposition, n::Integer, Y::AbstractArray;
objective::L2, kwargs...)
    N = ndims(T)
    if n in 1:N # the matrix factors start at m=1
        function gradient_matrix(X::AbstractTucker; kwargs...)
            B = core(X)
            matrices = matrix_factors(X)
            A_n = factor(X, n)
            X_n = tuckerproduct(B, matrices; exclude=n)
            grad = A_n * slicewise_dot(X_n, X_n; dims=n)
                - slicewise_dot(Y, X_n; dims=n)
            return grad
        end
        return gradient_matrix
    else
        error("No $(n)th factor in Tucker")
    end
end

```

3.2.c Computing Lipschitz Step-sizes

Similar to automatic differentiation, there exist “automatic Lipschitz” calculations to upper bound the Lipschitz constant of a function [21].

For the family of Tucker decompositions, we can compute the Lipschitz constants of the gradient efficiently similar to how we compute the gradient in Section 3.2.b with the following corollaries of Theorem 2.2.

Corollary 3.1: Let $B \in \mathbb{R}^{R_1 \times \dots \times R_N}$, $A_m \in \mathbb{R}^{I_m \times R_m}$, and $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$. The function

$$f(A) = \frac{1}{2} \| [B; A_1, \dots, A_{n-1}, A, A_{n+1}, \dots, A_N] - Y \|_F^2$$

is quadratic, and L -smooth with constant

$$L_{A_n} = \| \tilde{X}_n \cdot_n \tilde{X}_n \|_{\text{op}}$$

where

$$\tilde{X}_n = \left(B \bigtimes_{m \neq n} A_m \right) = \llbracket B; A_1, \dots, A_{n-1}, \text{id}_{R_n}, A_{n+1}, \dots, A_N \rrbracket.$$

Proof. The result follows from Theorem 2.1 and Theorem 2.2.

Corollary 3.2: Let $A_n \in \mathbb{R}^{I_n \times R_n}$, and $Y \in \mathbb{R}^{I_1 \times \dots \times I_N}$. The function

$$f(B) = \frac{1}{2} \| [B; A_1, \dots, A_N] - Y \|_F^2$$

is quadratic, and L -smooth with constant

$$L_B = \prod_{n=1}^N \| A_n^\top A_n \|_{\text{op}}.$$

Proof. The result follows from Theorem 2.1 and Theorem 2.2.

This yields the following efficient implementations.

! Note

It is tempting to use the identity $\|A^\top A\|_{\text{op}} = \|A\|_{\text{op}}^2$ to calculate the Lipschitz constant without forming $A^\top A$. For tall dense matrices, using this identity is slower and more memory intensive as of Julia 1.11.2. See [LinearAlgebra.jl issue 1185](#) on Github.

```
function make_lipschitz(T::Tucker, n::Integer, Y::AbstractArray; objective::L2,
kwargs...)
    N = ndims(T)
    if n==0 # the core is the zeroth factor
        function lipschitz_core(X::AbstractTucker; kwargs...)
            return prod(A -> opnorm(A'A), matrix_factors(X))
        end
        return lipschitz_core
    elseif n in 1:N # the matrix is the zeroth factor
        function lipschitz_matrix(X::AbstractTucker; kwargs...)
            B = core(X)
            matrices = matrix_factors(X)
            Ī_n = tuckerproduct(B, matrices; exclude=n)
            return opnorm(slicewise_dot(Ī_n, Ī_n; dims=n))
        end
        return lipschitz_matrix
    else
        error("No $(n)th factor in Tucker")
    end
end
```

In the case of Tucker decomposition, the Lipschitz constants simplify to

$$L_B = \|A^\top A\|_{\text{op}}, \quad L_B = \|B \cdot_1 B\|_{\text{op}}.$$

```
function      make_lipschitz(T::Tucker1,      n::Integer,      Y::AbstractArray;
objective::L2, kwargs...)
    if n==0 # the core is the zeroth factor
        function lipschitz0(X::Tucker1; kwargs...)
            A = matrix_factor(X, 1)
            return opnorm(A'A)
        end
        return lipschitz0
    elseif n==1 # the matrix is the zeroth factor
        function lipschitz1(X::Tucker1; kwargs...)
            B = core(X)
```

```

        return opnorm(slicewise_dot(B, B))
    end
    return lipschitz1

else
    error("No $(n)th factor in Tucker1")
end
end

```

Lastly, for CP decomposition, the Lipschitz constants for the matrices can be calculated similarly to the Tucker decomposition.

```

function make_lipschitz(T::CPDecomposition, n::Integer, Y::AbstractArray;
objective::L2, kwargs...)
    N = ndims(T)
    if n in 1:N
        function lipschitz_matrix(X::CPDecomposition; kwargs...)
            id = core(X)
            matrices = matrix_factors(X)
            X̃_n = tuckerproduct(id, matrices; exclude=n)
            return opnorm(slicewise_dot(X̃_n, X̃_n; dims=n))
        end
        return lipschitz_matrix
    else
        error("No $(n)th factor in CPDecomposition")
    end
end

```

4 Computational Techniques

- As stated, algorithm works
- But can be slow, especially for constrained or large problems

As stated, the algorithm described in Section 3.2 works. It will converge to a solution to our optimization problem and factorize the input tensor. It is worth discussing how the algorithm can be modified to improve convergence to maintain quick convergence for large problems, and what sort of architectural methods are used to allow for maximum flexibility, without over engineering the package.

4.1 For Improving Convergence Speed

There are a few techniques used to assist convergence. Two ideas that are well studied are discussed in this section. They are 1) breaking up the updates into smaller blocks, and 2) using momentum or acceleration. What is perhaps novel is considering the synergy between these two ideas.

Two more techniques are implemented in `BlockTensorDecomposition.jl` to improve convergence. To the authors knowledge, these are new to tensor factorization, but may or may not be applicable depending on the exact factorization problem or data being studied. For these reasons, these other techniques are discussed separately in Section 5 and Section 6.

4.1.a Sub-block Descent

- Use smaller blocks, but descent in parallel (sub-blocks don't wait for other sub-blocks)
- Can perform this efficiently with a "matrix step-size"

When using block coordinate descent as in Section 3.2, it is natural to treat each factor as its own block. This requires the fewest blocks while ensuring the objective is still convex with respect to each block. We could just as easily use smaller blocks.

In the case of Tucker decomposition, one modification of the update shown in Equation 10 would be to update each column $a_{n,r}$, $r = 1, \dots, R_n$ of the matrix A_n separately. This would be suitable if the constraint that $A_N \in \mathcal{C}_n$ can be broken up further into the constraints $a_{n,r} \in \mathcal{C}_{n,r}$. This is shown in the following update scheme:

$$a_{n,r}^{t+1} \leftarrow P_{\mathcal{C}_{n,r}} \left(a_{n,r}^t - \frac{1}{L_{n,r}^t} \nabla f_{n,r}^t(a_{n,r}^t) \right), \quad (13)$$

where $f_{n,r}^t(a) = \frac{1}{2} \| [B; A_1^{t+1}, \dots, A_{n-1}^{t+1}, A_{n,r}(a), A_{n+1}^t, \dots, A_N^t] - Y \|_F^2$ and

$$A_{n,r}^t(a) = \begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ a_{n,1}^{t+1} & \cdots & a_{n,r-1}^{t+1} & a & a_{n,r+1}^t & \cdots & a_{n,R_n}^t \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & & \end{bmatrix}. \quad (14)$$

In theory, the block update shown in Equation 13 should be a bit more expensive than using the larger blocks on the matrices A shown in Equation 10, since the gradient needs to be recomputed R_n times for each matrix block n , rather than only computing the gradient once per block n . To get around this, we use the fact that $\nabla f_{n,r}^t(a)$ is the r th column from the gradient $\nabla f_n^t(A)$ where $f_n^t(A) = \frac{1}{2} \| [B; A_1^{t+1}, \dots, A_{n-1}^{t+1}, A, A_{n+1}^t, \dots, A_N^t] - Y \|_F^2$. So we can approximate Equation 13 by first calculating the gradient ∇f_n^t at

$$\hat{A}_n^t = \begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ a_{n,1}^t & \cdots & a_{n,r-1}^t & a_{n,r}^t & a_{n,r+1}^t & \cdots & a_{n,R_n}^t \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & & \end{bmatrix}, \quad (15)$$

and then updating each sub-block r according to

$$a_{n,r}^{t+1} \leftarrow P_{\mathcal{C}_{n,r}} \left(a_{n,r}^t - \frac{1}{L_{n,r}^t} \nabla f_n^t(\hat{A}_n^t) \right). \quad (16)$$

Note the difference between Equation 14 and Equation 15 is that we don't use the most recent columns $a_{n,j}$ for $j < r$ in Equation 15.

The update given in Equation 16 can be merged back to an update on the whole block A_n

$$A_n^{t+1} \leftarrow P_{\mathcal{C}_n} \left(A_n^t - \nabla f_n^t(A_n^t) (\hat{L}_n^t)^{-1} \right) \quad (17)$$

where we have the $R_n \times R_n$ diagonal "Lipschitz Matrix"

$$\hat{L}_n^t = \begin{bmatrix} L_{n,1}^t & 0 & 0 \\ 0 & L_{n,2}^t & 0 \\ & \ddots & \vdots \\ 0 & 0 & \dots & L_{n,R_n}^t \end{bmatrix}.$$

It is not too hard to show that the Lipschitz $L_{n,r}^t$ for ∇f_n^t is the Euclidean norm of the r th column⁴ of the matrix $\tilde{X}_n \cdot_n \tilde{X}_n$ from Corollary 3.1,

$$L_{n,r}^t = \|(\tilde{X}_n \cdot_n \tilde{X}_n)[:,r]\|_2.$$

This leads to the following efficient calculation of the Lipschitz matrix in Julia.

```
function diagonal_lipschitz_matrix(T::Tucker, n::Int; kwargs...)
    B = core(T)
    matrices = matrix_factors(T)
    X_n = tuckerproduct(B, matrices; exclude=n)
    return Diagonal_col_norm(slicewise_dot(X_n, X_n; dims=n))
end

Diagonal_col_norm(X) = Diagonal(norm.(eachcol(X)))
```

We can now compare the merged sub-block update Equation 17 to the standard projected gradient descent update shown in Equation 10. The difference is that we calculate a "matrix step-size" $\hat{L}_n^t \in \mathbb{R}^{R_n \times R_n}$ rather than a scalar $L_n^t \in \mathbb{R}$. In practice, this leads to an improvement in convergence speed for two reasons.

First, computing the matrix \hat{L}_n^t often faster than the scalar $L_n^t = \|\tilde{X}_n \cdot_n \tilde{X}_n\|_{\text{op}}$. The former only requires calculating the Euclidean norm of R vectors for a total cost of $2R_n^2$ floating point operations (FLOPs), whereas the latter requires the top eigenvalue of $\tilde{X}_n \cdot_n \tilde{X}_n$. This is usually done with a power method or truncated SVD which can be costlier than the flat rate of $2R_n^2$ FLOPs.

⁴We could have used the r th row of $\tilde{X}_n \cdot_n \tilde{X}_n$ since this matrix is symmetric. Since Julia store matrices in column-major order, many operations that perform column-wise are more efficient than their equivalent row-wise operation.

Secondly, using the matrix \hat{L}_n^t means columns where $L_{n,r}^t$ is small can take larger descent steps. This is because the largest singular value of $\tilde{X}_n \cdot_n \tilde{X}_n$ is an upper bound on the Euclidean norm of each column: $L_{n,r}^t \leq L_n^t$. Using the scalar Lipschitz L_n^t is equivalent to the diagonal matrix

$$D = \begin{bmatrix} L_n^t & 0 & 0 \\ 0 & L_n^t & 0 \\ & \ddots & \vdots \\ 0 & 0 & \dots & L_n^t \end{bmatrix}$$

in the merged sub-block update shown in Equation 17. So each column of A_n is forced to use the worst case (largest) singular value of $\tilde{X}_n \cdot_n \tilde{X}_n$. In this way, the matrix \hat{L}_n^t acts like a cheap approximate Hessian as if we were doing a quasi-Newton update with step-size 1.

For completeness, we can perform the same merged sub-block update to update the core B . In this case, we obtain the more complicated “Lipschitz tensor” $\hat{L}_B^t \in \mathbb{R}^{(R_1 \times \dots \times R_N)^2}$ defined by

$$\hat{L}_B^t = \hat{L}_{B,1}^t \otimes \dots \otimes \hat{L}_{B,N}^t$$

where each matrix $\hat{L}_{B,n}^t \in \mathbb{R}^{R_n \times R_n}$ is diagonal with non-zero entries

$$L_{B,n}^t[r, r] = \|((A_n^t)^\top A_n^t)[:, r]\|_2.$$

The merged sub-block update for the core becomes

$$B^{t+1} \leftarrow P_{\mathcal{C}_B} \left(B^t - \nabla f_0^t(B^t) \times_B (\hat{L}_B^t)^{-1} \right) \quad (18)$$

with the multiplication

$$\begin{aligned} \nabla f_0^t(B^t) \times_B (\hat{L}_B^t)^{-1} &= \nabla f_0^t(B^t) \bigtimes_n (\hat{L}_{B,n}^t)^{-1} \\ &= [\nabla f_0^t(B^t); (\hat{L}_{B,1}^t)^{-1}, \dots, (\hat{L}_{B,N}^t)^{-1}]. \end{aligned} \quad (19)$$

This should be thought of as normalizing each dimension of the tensor $\nabla f_0^t(B^t)$ so that we can take a unit step-size.

TODO use the multi-mode product in stead of defining a new multiplication \times_B here. I think I'll have the same tensor product issue as before where the indices.

Putting the core and matrices Lipschitz calculations together gives us the following Julia code. Note we store \hat{L}_B^t in factored form as a tuple of diagonal matrices to save space and computation.

```
function make_block_lipschitz(T::AbstractTucker, n::Integer, Y::AbstractArray;
objective::L2, kwargs...)
    N = ndims(T)
    if n==0 # the core is the zeroth factor
        function lipschitz_core(X::AbstractTucker; kwargs...)
            # Implementation of lipschitz_core
        end
    else
        # Implementation for higher dimensions
    end
end
```

```

    return map(A -> Diagonal_col_norm(A'A), matrix_factors(X))
end # Returns a tuple of diagonal matrices
return lipschitz_core

elseif n in 1:N
    function lipschitz_matrix(X::AbstractTucker; kwargs...)
        matrices = matrix_factors(X)
        X̃_n = tuckerproduct(core(X), matrices; exclude=n)
        return Diagonal_col_norm(slicewise_dot(X̃_n, X̃_n; dims=n))
    end
    return lipschitz_matrix
else
    error("No $(n)th factor in Tucker")
end
end

```

4.1.b Momentum

- This one is standard
- Use something similar to [10]
- This is compatible with sub-block descent with appropriately defined matrix operations

In practice, we find that extrapolating the iterate based on the prior iterate

$$\hat{A}_n^t \leftarrow A_n^t + \omega_n^t (A_n^t - A_n^{t-1}) \quad (20)$$

for some amount of extrapolation $\omega_n^t \geq 0$ before applying the update Equation 17 greatly improves the speed of descent. This can be thought of as a type of momentum where we continue to move in directions that showed a lot of improvement during the last iteration.

Our selection for ω_n^t follows Xu and Yin's method for block coordinate descent [10], which is itself inspired by Tseng and Yun's coordinate gradient descent method [22].

Given a parameter⁵ $\delta \in [0, 1)$, we define the momentum parameters and τ^t and ω_n^t according to the following updates

$$\begin{aligned}
\tau^0 &= 1 \\
\tau^{t+1} &\leftarrow \frac{1}{2} \left(1 + \sqrt{1 + 4(\tau^t)^2} \right) \\
\hat{\omega}^t &\leftarrow \frac{\tau^t - 1}{\tau^{t+1}} \\
\omega_n^t &\leftarrow \min \left(\hat{\omega}^t, \delta \sqrt{\hat{L}_n^{t-1} (\hat{L}_n^t)^{-1}} \right).
\end{aligned} \quad (21)$$

⁵Usually we pick a number close to 1. For example, we use the default $\delta = 0.9999$.

TODO notation is going to get confusing. We use hat/not hat L for the scalar vs matrix/tensor version. But we use hat/not hat ω for the ideal vs clamped momentum.

What is novel with our approach is that we perform this momentum on the Lipschitz matrices and tensors \hat{L}_n^t rather than scalar Lipschitz constant L_n^t . In this way, we should interpret the operations shown in Equation 21 as operating element-wise. This also means the momentum parameter ω_n^t is a matrix or tensor and takes the same shape as \hat{L}_n^t .

In order to perform Equation 20, we use the equivalent but more efficient formulation

$$\hat{A}_n^t \leftarrow A_n^t (\text{id}_{R_n} + \omega_n^t) - A_n^{t-1} \omega_n^t.$$

```
function (U::MomentumUpdate)(X::T; X_last::T, ω, δ, kwargs...) where T
    n = U.n

    L = U.lipschitz(X; kwargs...)
    L_last = U.lipschitz(X_last; kwargs...)
    ω = min.(ω, δ .* √(L_last/L))

    A, A_last = factor(X, n), factor(X_last, n)

    A .= U.combine(A, id + ω)
    A .-= U.combine(A_last, ω)
end
```

In the code above, the momentum stores the factor n it acts on, how to compute the Lipschitz constant, matrix, or tensor, and how to combine (multiply) the constant with the factor. In the case of matrix factors A_n in a Tucker decomposition, this is simply right matrix-matrix multiplication. The core factor B uses \times_B as described in Equation 19.

The parameters of the momentum update are handled separately. This is to treat the momentum update as “apply this update with these parameters”. The parameters τ^t and $\hat{\omega}^t$ are updated by the following function that keeps track of all parameters needed to perform the iteration. In this case, we keep track of what iteration t we are at, the previous iterate, and a few options for the order in which to cycle through and update the blocks.

```
update_τ(τ) = 0.5*(1 + sqrt(1 + 4*τ^2))

function initialize_parameters(decomposition, Y, previous; momentum::Bool,
random_order, recursive_random_order, kwargs...)
    # parameters for the update step are symbol => value pairs
    # held in a dictionary since we may mutate these, e.g. the step-size
    parameters = Dict{Symbol, Any}()

    # General Looping
    parameters[:iteration] = 0
    parameters[:X_last] = previous[begin] # Last iterate
```

```

parameters[:random_order] = random_order
parameters[:recursive_random_order] = recursive_random_order

# Momentum
if momentum
    parameters[:τ_last] = float(1) # need this field to hold Floats, not
Ints
    parameters[:τ] = update_τ(float(1))
    parameters[:ω] = (parameters[:τ_last] - 1) / parameters[:τ]
    parameters[:δ] = kwargs[:δ]
end

function updateparameters!(parameters, decomposition, previous)
    parameters[:iteration] += 1
    # parameters[:x_last] = previous[begin]
# This is commented since parameters[:x_last] already points to previous[begin]

    if momentum
        parameters[:τ_last] = parameters[:τ]
        parameters[:τ] = update_τ(parameters[:τ_last])
        parameters[:ω] = (parameters[:τ_last] - 1) / parameters[:τ]
    end
end

return parameters, updateparameters!
end

```

4.1.c Empirical Evidence for Sub-Block Descent and Momentum

To showcase that the combination of these two tricks can speed up convergence, we will benchmark them by factorizing a random 10×10 tensor (a matrix) with rank 3. The Julia code is shown below, and the results are shown in Table 2.

```

using BenchmarkTools
using BlockTensorDecomposition

fact = BlockTensorDecomposition.factorize

options = (
    :rank => 3,
    :tolerance => (1, 0.03),
    :converged => (GradientNNCone, RelativeError),
    :δ => 0.9,
)

n_subblock_n_momentum(Y) = fact(Y;
    do_subblock_updates=false,
    momentum=false,
)

```

```

    options...
)

y_subblock_n_momentum(Y) = fact(Y;
    do_subblock_updates=true,
    momentum=false,
    options...
)

n_subblock_y_momentum(Y) = fact(Y;
    do_subblock_updates=false,
    momentum=true,
    options...
)

y_subblock_y_momentum(Y) = fact(Y;
    do_subblock_updates=true,
    momentum=true,
    options...
)

I, J = 10, 10
R = 3

@benchmark n_subblock_n_momentum(Y) setup=(Y=Tucker1((I, J), R))
@benchmark n_subblock_y_momentum(Y) setup=(Y=Tucker1((I, J), R))
@benchmark y_subblock_n_momentum(Y) setup=(Y=Tucker1((I, J), R))
@benchmark y_subblock_y_momentum(Y) setup=(Y=Tucker1((I, J), R))

performance_increase(old, new) = (old - new) / new * 100

```

The code `Tucker1((I, J), R)` produces a random $I \times J$ rank- R matrix by generating two matrices $A \in \mathbb{R}^{I \times R}$ and $B \in \mathbb{R}^{R \times J}$ with standard normal entries, and multiplies them together.

Table 2: Summary of median times to factorize a random 10×10 rank-3 matrix under different methods. The performance increase is given by the formula $(\text{old} - \text{new}) / \text{new}$.

	No Momentum	Yes Momentum
No Sub-Block	48.843 ms	45.738 ms (6.7887% faster)
Yes Sub-Block	27.473 ms (77.785% faster)	24.350 ms (100.59% faster)

In Table 2, you can see that having both sub-block descent and momentum yields the fastest factorization. Moreover, the performance increase is *more* than simply the performance increases

obtained by exclusively sub-block or momentum alone.⁶ This suggests that there is synergy with these two methods and are best used together.

TODO Repeat this experiment on a less trivial factorization. What I've done above can be done with an SVD in less time. Ideally use a $10 \times 10 \times 10$ Tucker decomposition with rank $2 \times 3 \times 4$.

4.2 For Flexibility

- there are a number of software engineering techniques used
- these help flexibility for hot swapping and a language for making custom...
 - convergence criterion (and having multiple stopping conditions)
 - probing info during the iterations (stats collected at the end)
 - having multiple constraints and ways to enforce them
 - cyclically or partially randomly or fully randomly update factors
- smart enough to apply these in a reasonable order

There are a number of software engineering techniques used to ensure BlockTensorDecomposition.jl is flexible and applicable to a wide range of problems. These enable key algorithmic choices to be hot-swapped and easily compared with each other.

4.2.a Convergence Criteria and Stats

- Can request info about any factor at each outer iteration
- any subset of stats can be the convergence criteria

Some iterative algorithms produce the exact solution of a problem after a finite number of iterations. Generalized minimal residual method (GMRES) is a good example of this [TODO cite!]. Our algorithm, like many others, only converges to the exact solution in the limit as the number of iterations grow. Since we would like a solution in finite time, we must halt the algorithm early.

In finite precision, we can halt the algorithm if we can guarantee the solution is accurate to machine precision. This can often be too strict if convergence is not at a fast enough rate. Furthermore, depending on *why* we are decomposing a tensor, we may want different stats to be within a given a tolerance. BlockTensorDecomposition.jl attempts to solve this issue by defining some standard criteria that can be used to halt the algorithm. These are subtypes of the abstract type `AbstractStat` and are listed below.

```
# X is the decomposition model e.g. Tucker((B, A1, A2, A2))
# Y is the input tensor we want to decompose
# norm(X) is the Frobenius norm of X

GradientNorm # norm(∇f(X))
GradientNNCone # sqrt(sum(norm(∇f(Ai)[Ai .> 0 .| ∇f(Ai) .< 0])^2 for Ai in
factors(X)))
ObjectiveValue # 1/2 norm(X - Y)^2
ObjectiveRatio # norm(X_last - Y)^2 / norm(X_current - Y)^2
```

⁶The expected performance increase if sub-block descent and momentum where independence would be $(1 + 0.067887)(1 + 0.77785) = 1.89854$ or only 89.854% faster.

```

RelativeError # norm(X - Y) / norm(Y)
IterateNormDiff # norm(X_current - X_last)
IterateRelativeDiff # norm(X_current - X_last) / norm(X_last)

```

Most of these are self explanatory, except perhaps `GradientNNCone`. When performing first order unconstrained optimization, we usually slow down progress when the norm of the gradient is small. In the limit, we expect to converge to a stationary point where the gradient is zero. When performing optimization under nonnegative constraints, and even further restrictions like simplex constraints, it makes more sense to ignore entries of the gradient where X is negative and the gradient is positive since those coordinates of X will not change after they are projected back to the nonnegative constraint.

TODO add theory of negative gradient is in the normal cone of the constraint?

As many or as few of these stats can be used in the call to `factorize`. Their tolerances can also be set independently of each other. The algorithm stops iterating when at least one of the criteria has a value less than its tolerance. As a fail safe, we also define one more type, `Iteration`, which is always active and halts the algorithm when that many iterations have past.

The `AbstractStat` type can also be subtyped to create custom stats that may be used to probe the iterates and diagnose issues.

```

EuclidianStepSize
EuclidianLipschitz
FactorNorms

```

These stats are recorded every iteration in a `DataFrame` and are one of the returns of `factorize`.

Finally, there are two auxiliary stats `PrintStats` and `DisplayDecomposition` which can be used to print all stats or the current decomposition each iteration.

4.2.b BlockUpdate Language

- construct the updates as a list of updates
- very functional programming
- can apply them in sequence or in a random order (or partially random)

4.2.c Constraints

- one type of update (other than the typical GD update)
- can combine them with composition
 - which is different than projecting onto their intersection!
- Constraint updates combine the constraint with how they are enforced
 - need to go together since there are multiple ways to enforce them e.g. simplex (see next section)

5 Partial Projection and Rescaling

- for bounded linear constraints

- ▶ first project
- ▶ then rescale to enforce linear constraints
- faster to execute than a projection
- often does not loose progress because of the rescaling (decomposition dependent)

6 Multi-scale

- use a coarse discretization along continuous dimensions
- factorize
- linearly interpolate decomposition to warm start larger decompositions

7 Conclusion

- all-in-one package
- provide a playground to invent new decompositions
- like auto-diff for factorizations

8 Appendix

8.1 Building the Hessian from two gradients

To build the Hessian from the definition of the gradient, we first extend the gradient to tensor-valued functions. For a function $F : \mathbb{R}^{J_1 \times \dots \times J_N} \rightarrow \mathbb{R}^{I_1 \times \dots \times I_M}$ where

$$F(X) = [f_{i_1 \dots i_M}(X)]$$

is a tensor of scalar functions $f_{i_1 \dots i_M} : T \rightarrow \mathbb{R}$, the gradient of F at X is defined entry-wise as

$$\begin{aligned} \nabla F(X)[i_1, \dots, i_M][j_1, \dots, j_N] &= \nabla f_{i_1 \dots i_M}(X)[j_1, \dots, j_N] \\ &= \frac{\partial f_{i_1 \dots i_M}}{\partial X[j_1, \dots, j_N]}(X). \end{aligned}$$

This treats the gradient at X as a tensor of tensors $\nabla F(X) \in (\mathbb{R}^{I_1 \times \dots \times I_M})^{J_1 \times \dots \times J_N}$. This is naturally isomorphic to a tensor of order $M + N$ with entries

$$\nabla F(X)[i_1, \dots, i_M, j_1, \dots, j_N] = \frac{\partial f_{i_1 \dots i_M}}{\partial X[j_1, \dots, j_N]}(X). \quad (22)$$

So we conclude that the gradient at X of a tensor-valued function F is $\nabla F : \mathbb{R}^{J_1 \times \dots \times J_N} \rightarrow \mathbb{R}^{I_1 \times \dots \times I_M \times J_1 \times \dots \times J_N}$ is given by Equation 22.

We can define the Hessian of a scalar function $f : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}$ at X as $\nabla^2 f(X) = \nabla(\nabla f)(X) : \mathbb{R}^{I_1 \times \dots \times I_N} \rightarrow \mathbb{R}^{(I_1 \times \dots \times I_N)^2}$. The inner nabla ∇ is the gradient of the scalar function f , and the outer nabla ∇ is the gradient of the tensor-valued function ∇f .

This means

$$\nabla^2 f(A)[i_1, \dots, i_N, j_1, \dots, j_N] = \frac{\partial^2 f}{\partial A[j_1, \dots, j_N] \partial A[i_1, \dots, i_N]}(A),$$

but if the function has continuous second derivatives, we can perform the partial derivatives in either order

$$\frac{\partial^2 f}{\partial A[j_1, \dots, j_N] \partial A[i_1, \dots, i_N]}(A) = \frac{\partial^2 f}{\partial A[i_1, \dots, i_N] \partial A[j_1, \dots, j_N]}(A).$$

Bibliography

- [1] Martín Abadi *et al.*, “TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems,” 2015. <https://www.tensorflow.org/>
- [2] J. Ansel *et al.*, “PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation,” in *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, Apr. 2024, vol. 2, pp. 929–947. doi: 10.1145/3620665.3640366.
- [3] J. Kossaifi, Y. Panagakis, A. Anandkumar, and M. Pantic, “TensorLy: Tensor Learning in Python,” *Journal of Machine Learning Research*, vol. 20, no. 26, pp. 1–6, 2019, Accessed: Aug. 16, 2024. [Online]. Available: <http://jmlr.org/papers/v20/18-277.html>
- [4] B. W. Bader and T. G. Kolda, “Tensor Toolbox for MATLAB.” Sep. 2023.
- [5] J. Li, J. Bien, and M. T. Wells, “rTensor: An R Package for Multidimensional Array (Tensor) Unfolding, Multiplication, and Decomposition,” *Journal of Statistical Software*, vol. 87, pp. 1–31, Nov. 2018, doi: 10.18637/jss.v087.i10.
- [6] Jutho, “Jutho/TensorKit.jl,” Aug. 2024. <https://github.com/Jutho/TensorKit.jl> (accessed Aug. 15, 2024).
- [7] M. Abbott *et al.*, “mcabbott/Tullio.jl: v0.3.7,” Oct. 2023. <https://doi.org/10.5281/zenodo.10035615>
- [8] A. Peter, “under-Peter/OMEinsum.jl,” Aug. 2024. <https://github.com/under-Peter/OMEinsum.jl> (accessed Aug. 16, 2024).
- [9] Y.-J. Wu, “yunjhongwu/TensorDecompositions.jl,” Feb. 2024. <https://github.com/yunjhongwu/TensorDecompositions.jl> (accessed Aug. 16, 2024).
- [10] Y. Xu and W. Yin, “A Block Coordinate Descent Method for Regularized Multiconvex Optimization with Applications to Nonnegative Tensor Factorization and Completion,” *SIAM Journal on Imaging Sciences*, vol. 6, no. 3, pp. 1758–1789, Jan. 2013, doi: 10.1137/120887795.
- [11] J. Kim, Y. He, and H. Park, “Algorithms for nonnegative matrix and tensor factorizations: a unified view based on block coordinate descent framework,” *Journal of Global Optimization*, vol. 58, no. 2, pp. 285–319, Feb. 2014, doi: 10.1007/s10898-013-0035-4.

- [12] Z. Yang and E. Oja, “Unified Development of Multiplicative Algorithms for Linear and Quadratic Nonnegative Matrix Factorization,” *IEEE Transactions on Neural Networks*, vol. 22, no. 12, pp. 1878–1891, Dec. 2011, doi: 10.1109/TNN.2011.2170094.
- [13] T. G. Kolda and B. W. Bader, “Tensor Decompositions and Applications,” *SIAM Review*, vol. 51, no. 3, pp. 455–500, Aug. 2009, doi: 10.1137/07070111X.
- [14] L. Qi, Y. Chen, M. Bakshi, and X. Zhang, “Triple Decomposition and Tensor Recovery of Third Order Tensors,” Mar. 01, 2020. <http://arxiv.org/abs/2002.02259> (accessed Aug. 01, 2023).
- [15] F. Wu, C. Li, and Y. Li, “Manifold Regularization Nonnegative Triple Decomposition of Tensor Sets for Image Compression and Representation,” *Journal of Optimization Theory and Applications*, vol. 192, no. 3, pp. 979–1000, Mar. 2022, doi: 10.1007/s10957-022-02001-6.
- [16] I. V. Oseledets, “Tensor-Train Decomposition,” *SIAM Journal on Scientific Computing*, vol. 33, no. 5, pp. 2295–2317, Jan. 2011, doi: 10.1137/090752286.
- [17] J. B. Kruskal, “Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics,” *Linear Algebra and its Applications*, vol. 18, no. 2, pp. 95–138, Jan. 1977, doi: 10.1016/0024-3795(77)90069-6.
- [18] A. Bhaskara, M. Charikar, and A. Vijayaraghavan, “Uniqueness of Tensor Decompositions with Applications to Polynomial Identifiability,” in *Proceedings of The 27th Conference on Learning Theory*, May 2014, pp. 742–778. Accessed: Jan. 08, 2025. [Online]. Available: <https://proceedings.mlr.press/v35/bhaskara14a.html>
- [19] S. A. Vavasis, “On the Complexity of Nonnegative Matrix Factorization,” *SIAM Journal on Optimization*, vol. 20, no. 3, pp. 1364–1377, Jan. 2010, doi: 10.1137/070709967.
- [20] Y. Nesterov, “Nonlinear Optimization,” *Lectures on Convex Optimization*. Springer International Publishing, Cham, pp. 3–58, 2018. doi: 10.1007/978-3-319-91578-4_1.
- [21] A. Virmaux and K. Scaman, “Lipschitz regularity of deep neural networks: analysis and efficient estimation,” in *Advances in Neural Information Processing Systems*, 2018, vol. 31. Accessed: Jan. 11, 2025. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/hash/d54e99a6c03704e95e6965532dec148b-Abstract.html>
- [22] P. Tseng and S. Yun, “A coordinate gradient descent method for nonsmooth separable minimization,” *Mathematical Programming*, vol. 117, no. 1, pp. 387–423, Mar. 2009, doi: 10.1007/s10107-007-0170-0.