# Before We Start: Teaching Evaluation!

# Before We Start: Teaching Evaluation!

**Note: Teaching evaluation results will be used to decide whether I can stay professor in Oldenburg.**

**Additional Note: Please use the free-text input for (1) improvement suggestions and for (2) raising issues that are not covered by the form (e.g. technical issues in hybrid lectures).**

**Please try to provide improvement suggestions for all issues raised!**

# Semester Schedule

| CW 42 | 14. Oct | Lecture | 1 | Orga & Intro | 1-26 |
| CW 43 | 21. / 23. Oct | Lecture + Exercises | 2 | Probability, Statistics & Correlation | 27-56 |
| CW 44 | 28. Oct | Lecture | 3 | Data Munging, Cleaning & Bias | 57-94 / "Invisible Women" |
| CW 45 | 04. / 06. Nov | Lecture + Exercises | 4 | Scores & Rankings | 95-120 |
| CW 46 | 11. Nov | Lecture | 5 | Statistical Distributions & Significance | 121-154 |
| CW 47 | 18. / 20. Nov | Lecture + Exercises | 6 | Building & Evaluating Models | 201-236 |
| CW 48 | 25. Nov | Guest Lecture | 7 | Data Visualization | 155-200 |
| CW 49 | 02. / 04. Dec | Lecture + Exercises | 8 | Intro to Machine Learning | 351-390 |
| CW 50 | 09. Dec | Lecture | 9 | Linear Algebra | 237-266 |
| CW 51 | 16. / 18. Dec | Lecture + Exercises | 10 | Linear Regression & Gradient Descent | 267-288 |
| CW 02 | 06. Jan | Lecture | 11 | Logistic Regression & Classification | 289-302 |
| CW 03 | 13. / 15. Jan | Lecture + Exercises | 12 | Nearest Neighbor Methods & Clustering | 303-350 |
| CW 04 | 20. Jan | Lecture | 13 | Data Science in the Wild | 391-426 |
| CW 05 | 27. / 29. Jan | Lecture + Exercises | 14 | Q&A / Feedback | |
| CW 06 | 03. / 04. Feb | Oral Exams (Block 1) | | Preparation in our last session | |
| CW 13 | 24. / 25. Mar | Oral Exams (Block 2) | | ("Oral Exam Briefing") | |

# Get the Matrix!

The most critical part of a data science project often is reducing all the information you can find to one or more data matrices, ideally as large as possible.

Rows are examples.

Columns are distinct features/attributes.

# Linear Algebra

Linear algebra is the mathematics of matrices.

<span style="color:red">This makes it the language of data science.</span>

Many machine learning algorithms are best understood through linear algebra.

You presumably had an undergraduate course in linear algebra, but here I will review what you need to know.

# **What Can n·m Matrices Represent?**

- <span style="color:red">Data</span>: rows are objects, columns features.
- <span style="color:red">Geometric point sets</span>: rows are points, columns are dimensions
- <span style="color:red">Systems of Equations</span>: rows are equations, columns are coefficients for each variable.
- <span style="color:red">Graphs/Networks</span>: M[i,j] denotes the number of edges from vertex i to vertex j.
- <span style="color:red">Vectors</span>: any row, column or d·1 matrix

# Linear Algebra Formulae

- Concise formulas written as products of matrices provides great power.
- Algebraic substitution coupled with a rich set of identities yields elegant, mechanical ways to manipulate such formulae.
- But such strings of operations can, be difficult to interpret and understand.

$$w = (A^T A)^{-1} A^T b \leftarrow \text{Normal Form Equation}$$

# Algebraic Proof: 2=1

$$
\begin{aligned}
a &= b \\
a^2 &= ab \\
a^2 - b^2 &= ab - b^2 \\
(a+b)(a-b) &= b(a-b) \\
a+b &= b \\
2b &= b \\
2 &= 1
\end{aligned}
$$

Division by zero not allowed!

# Lessons from the Proof

- Algebraic proofs often do not transport intuition about *why* things work very well
- They are easier to verify than to create.
- Even so, there are special cases / singularities to watch for, like division by 0.
- In linear algebra, such cases include singular / non-invertible matrices.

# Points vs. Vectors

Points in *d* dimensions can be represented as unit vectors (points on the sphere), plus their magnitudes.

Distances between points become angles between vectors, for purposes of comparison.
Ignoring magnitudes is a form of scaling, making all points directly comparable.

# Angles between Vectors

To compute angle AB: $\cos(\theta) = \dfrac{A \cdot B}{\|A\|\|B\|}.$

Cos(0)=1, Cos(Pi/2)=0, Cos(Pi)=-1

Scores like correlation coefficients:

Cos = correlation of mean zero variables!

For unit vectors, ||A||=||B||=1, so the angle between A and B is defined by the dot product.

**Recall:**

# Covariance & Pearson Correlation

$$r = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^{n}(X_i - \bar{X})^2}\sqrt{\sum_{i=1}^{n}(Y_i - \bar{Y})^2}}$$

*Covariance*

*Std. Dev. of X*  *Std. Dev. of Y*

A point *(x,y)* makes a positive contribution to *r* when both are above or below their means.

# Visualizing Matrix Operations

We use images to represent matrices

# A Matrix and its Transpose

The transpose of a matrix M interchanges rows and columns, turning an a·b matrix to a b·a matrix.



$$M_{ij}^T = M_{ji}$$

Note that colors get rescaled when magnitudes change.

# **Addition and Transposition**

A mix of scalar multiplication and addition.

# Linear Combination: B=(A+C)/2

Mix of scalar multiplication/ addition. $\alpha \cdot A + (1 - \alpha) \cdot B$

# Matrix Multiplication / Dot Products

The product A·B is defined by:  $C_{i,j} = \sum\limits_{i=1}^{k} A_{i,k} \cdot B_{k,j}$

$A \cdot B$ must share inner dimensions to multiply.

Each element of the product matrix is a dot product of row/column vectors.

Dot products measure how "in sync" the two vectors are, as in computing covariance or correlation.

# **Properties of Matrix Multiplication**

It is associative but not commutative:

$$\left(\begin{bmatrix} 3 & 2 \\ -1 & 0 \end{bmatrix}\begin{bmatrix} -2 & 3 \\ 4 & 2 \end{bmatrix}\right)\begin{bmatrix} -1 & 5 \\ 1 & 2 \end{bmatrix}$$
$$=\begin{bmatrix} 2 & 13 \\ 2 & -3 \end{bmatrix}\begin{bmatrix} -1 & 5 \\ 1 & 2 \end{bmatrix}$$
$$=\begin{bmatrix} 11 & 36 \\ -5 & 4 \end{bmatrix}$$

$$\begin{bmatrix} 3 & 2 \\ -1 & 0 \end{bmatrix}\left(\begin{bmatrix} -2 & 3 \\ 4 & 2 \end{bmatrix}\begin{bmatrix} -1 & 5 \\ 1 & 2 \end{bmatrix}\right)$$
$$=\begin{bmatrix} 3 & 2 \\ -1 & 0 \end{bmatrix}\begin{bmatrix} 5 & -4 \\ -2 & 24 \end{bmatrix}$$
$$=\begin{bmatrix} 11 & 36 \\ -5 & 4 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix}$$

Multiplication by the identity commutes: $IA = AI = A$

Although matrix algorithms can be simple to program, faster and more numerically stable algorithms often exist in highly optimized libraries.

# Multiplying Feature Matrices

Suppose A is an n·d data matrix. What is A times its transpose?

- $A \cdot A^T$ is an n·n matrix of dot products, measuring "in sync-ness" among *points*.
- $A^T \cdot A$ is a d·d matrix of dot products, measuring "in sync-ness" among *features*.

These are called covariance matrices.

# Row or Column Covariance Matrix?

$$A \cdot A^T \qquad A^T \cdot A$$

# Interpreting Matrix Multiplication

- Multiplying 0/1 adjacency matrices yield paths of length two:   $a[i,k]=a[i,j] \cdot a[j,k]$
- Multiplication by permutation matrices rearrange rows/columns:

$$P = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{pmatrix} \quad M = \begin{pmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \\ m_{41} & m_{42} & m_{43} & m_{44} \end{pmatrix}$$

$$PM = \begin{pmatrix} m_{31} & m_{32} & m_{33} & m_{34} \\ m_{11} & m_{12} & m_{13} & m_{14} \\ m_{41} & m_{42} & m_{43} & m_{44} \\ m_{21} & m_{22} & m_{23} & m_{24} \end{pmatrix}$$

# **Interpreting Matrix Multiplication**

● Rotating points in space:



Multiplying something by the right matrix can have magic properties, in arbitrary dimensions.

# Dividing Matrices

The inverse operation to multiplication is division.

An important special case of division is inversion: $A \cdot A^{-1} = I$ implies $A^{-1} = \dfrac{I}{A}$

In fact it is equivalent, because $\dfrac{A}{B} = A \cdot B^{-1}$

# Matrix Inversion

$A^{-1}$ is the multiplicative inverse of A if $A \cdot A^{-1} = I$ , where I is the identity matrix.

$$\mathbf{A}^{-1} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}.$$

If matrix A has an inverse, it can be computed by solving a linear system using Gaussian elimination.

$$[A|I] = \begin{bmatrix} 6 & 4 & 1 & 1 & 0 & 0 \\ 10 & 7 & 2 & 0 & 1 & 0 \\ 5 & 3 & 1 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & -1 \\ 0 & 1 & 0 & 0 & 1 & -2 \\ 5 & 3 & 1 & 0 & 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 & 1 & -1 & 1 \\ 0 & 1 & 0 & 0 & 1 & -2 \\ 5 & 3 & 1 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 1 & -1 & 1 \\ 0 & 1 & 0 & 0 & 1 & -2 \\ 0 & 0 & 1 & -5 & 2 & 2 \end{bmatrix}$$

$$\rightarrow A^{-1} = \begin{bmatrix} 1 & -1 & 1 \\ 0 & 1 & -2 \\ -5 & 2 & 2 \end{bmatrix}$$

# Inverse of Lincoln

$$L \qquad\qquad L^{-1} \qquad\qquad LL^{-1}$$

# Matrix Inversion and Linear Systems

Multiplying both sides of $Aw = b$ by the inverse of A yields: $(A^{-1}A)w = A^{-1}b$ or $w = A^{-1}b$

Thus solving linear equations is equivalent to matrix inversion.

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 5 & 1 \\ 2 & 3 & 8 \end{bmatrix}^{-1} \begin{bmatrix} 10 \\ 8 \\ 3 \end{bmatrix} = \frac{1}{25} \begin{bmatrix} -232 \\ 129 \\ 19 \end{bmatrix} = \begin{bmatrix} -9.28 \\ 5.16 \\ 0.76 \end{bmatrix}.$$

The inverse makes it cheap to evaluate many b vectors. However, Gaussian elimination is more numerically stable than inversion.

# Matrix Rank

Systems of equations are underdetermined if rows can be expressed as linear combinations of other rows.

The <span style="color:red">rank</span> of a matrix is a measure of the number of linearly independent rows.

An n·n matrix should be rank n for all operations to be properly defined on it.

# **Increasing Lincoln Memorial's Rank**

Some rows of the Lincoln Memorial are not linearly independent, so it is not full rank.

Adding small amounts of random noise increases rank without serious image distortion.

```
In[37]:= MatrixRank[m]

        508

In[44]:= noise = Table[ Table[ RandomReal[{-0.5, 0.5}], {512}], {512}];

In[45]:= MatrixRank[m + noise]

Out[45]= 512
```

# Factoring Matrices

Many important machine learning algorithms can be viewed as factoring a matrix.

Suppose $n \cdot m$ matrix *A* can be expressed as the product $B \cdot C$, i.e an $n \cdot k$ matrix times a $k \cdot m$ matrix.

If k<min(n,m), B and C compress matrix A.

Further, B is a small feature matrix replacing A.

# **Factoring Word-Document Matrices**

If A is a document/word co-occurrence matrix, and A=BC, where B is d·k and C is k·w:

- ● B is a compressed feature vector for docs
- ● C is a compressed feature vector for words

# Factoring Word-Document Matrices

# LU Decomposition

Factoring a matrix M representing lower and upper triangular matrices L and U prove useful in solving linear systems.

The determinant of M is the product of the main diagonal elements of U.

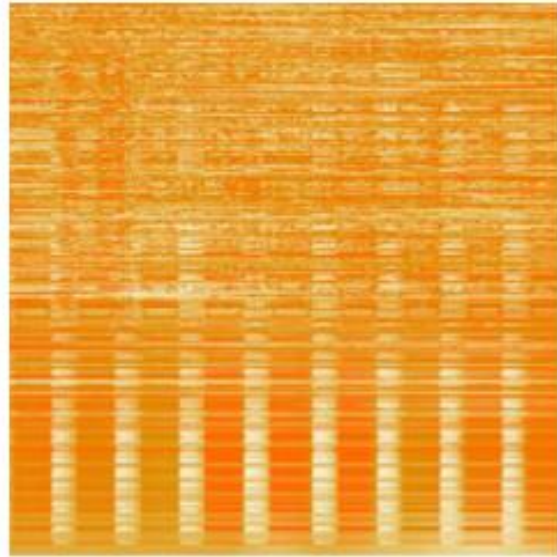A determinant of 0 means the matrix is not full rank.

# LU Decomposition: Example

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \cdot \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix}$$

# LU Decomposition of Memorial

Rows were permuted by this solver.

# Lessons from Lincoln

- Multiplying the factors of the matrix did not reconstruct it exactly, due to numerical instability.
- A high matrix condition number can tip you off that you are in trouble.
- Still, the gross features of the data are largely preserved.

# **Eigenvalues and Eigenvectors**

Multiplying a vector U by a matrix A can have the same effect as multiplying it by a scalar I.

$$
\overset{A}{\begin{bmatrix} -5 & 2 \\ 2 & -2 \end{bmatrix}} \cdot \overset{U_1}{\begin{bmatrix} 1 \\ 2 \end{bmatrix}} = \overset{I_1}{-1} \overset{U_1}{\begin{bmatrix} 1 \\ 2 \end{bmatrix}} \qquad \overset{A}{\begin{bmatrix} -5 & 2 \\ 2 & -2 \end{bmatrix}} \cdot \overset{U_2}{\begin{bmatrix} 2 \\ -1 \end{bmatrix}} = \overset{I_2}{-6} \overset{U_2}{\begin{bmatrix} 2 \\ -1 \end{bmatrix}}
$$

Thus the eigenvalue-eigenvector pairs ($U_i$,$I_i$) must encode a lot of information about matrix A!

# Computing Eigenvalues

The n distinct eigenvalues of a rank n matrix can be found by factoring its characteristic equation:

$$\det(A - \lambda I) = \begin{vmatrix} -\lambda & 1 & 1 \\ 1 & -\lambda & 1 \\ 1 & 1 & -\lambda \end{vmatrix}$$

$$= -\lambda^3 + 3\lambda + 2$$

$$= (\lambda - 2)(\lambda + 1)^2 \qquad \lambda_1 = 2, \; \lambda_{2,3} = -1$$

Faster algorithms exist to find the largest eigenvalues, which are the most important.

# Computing Eigenvectors

The vector associated with a given eigenvalue can be computed by solving a linear system:

$$\mathbf{A} \cdot \mathbf{v}_1 = \lambda_1 \cdot \mathbf{v}_1$$

$$\left( \mathbf{A} - \lambda_1 \right) \cdot \mathbf{v}_1 = 0$$

$$\begin{bmatrix} -\lambda_1 & 1 \\ -2 & -3-\lambda_1 \end{bmatrix} \cdot \mathbf{v}_1 = 0$$

$$\begin{bmatrix} 1 & 1 \\ -2 & -2 \end{bmatrix} \cdot \mathbf{v}_1 = \begin{bmatrix} 1 & 1 \\ -2 & -2 \end{bmatrix} \cdot \begin{bmatrix} v_{1,1} \\ v_{1,2} \end{bmatrix} = 0$$

Another approach uses $v' = \dfrac{A \cdot v}{I}$ to compute approximations to v until it converges.

# **Properties of Eigenvalues/vectors**

- A full-rank matrix has n vector-value pairs.
- Each pair of vectors from a symmetric matrix are mutually orthogonal, like x-y axes. E.g. the dot product is of (2,-1) and (1,2) is zero.
- Thus eigenvectors can play the role of dimensions or basis in n-dimensional space.

Vectors/values are found by solving linear systems.
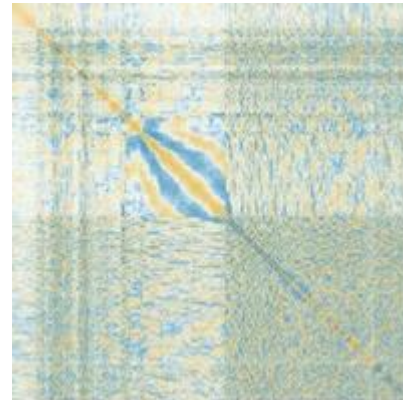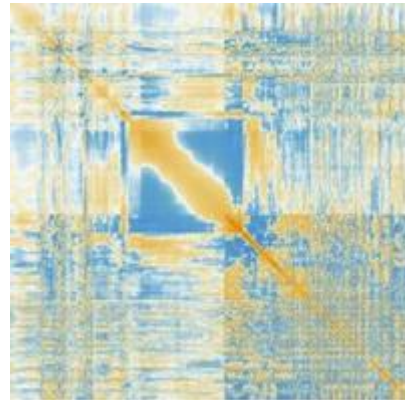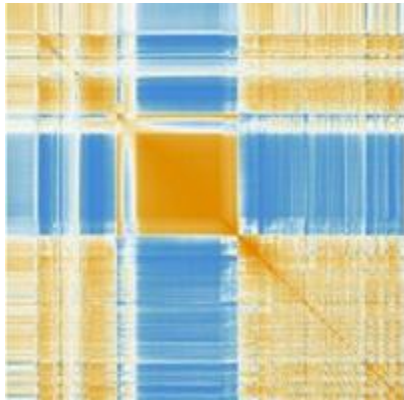
# Eigenvalue Decomposition

Any *n·n* symmetric matrix *M* can be decomposed into its n eigenvector products:

$$M = \sum_{i=1}^{n} \lambda_i U_i U_i^T$$

Larger eigenvalues correspond to more important vector products.

# Reconstructing a Covariance Matrix

Summing only the largest vectors performs dimension reduction, identifying the most important features of the matrix.



Covariance matrix error for the Lincoln memorial reduces when summing the 1, 5, and 50 largest eigenvectors for n=512

# Singular Value Decomposition (SVD)

The SVD of an n·m matrix M factors it $M = UDV^T$ where D is diagonal (weighted identity matrix)

Thus UD weights each column of U by D, as does $DV^T$.

Retaining only the rows/column with large weights permits us to compress m features with relatively little loss.

# Reconstruction from SVD

The outer product of vectors yields a matrix

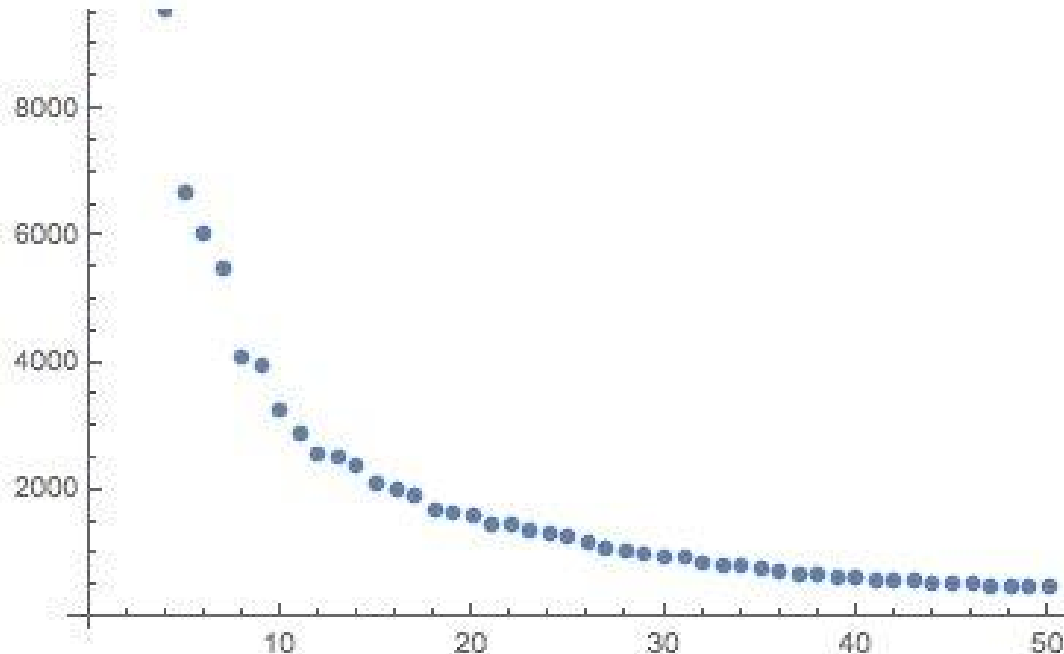$$P = X \bigotimes Y \qquad\qquad P[j, k] = X[j]Y[k]$$

Matrix M can be expressed a sum of outer products from SVD:

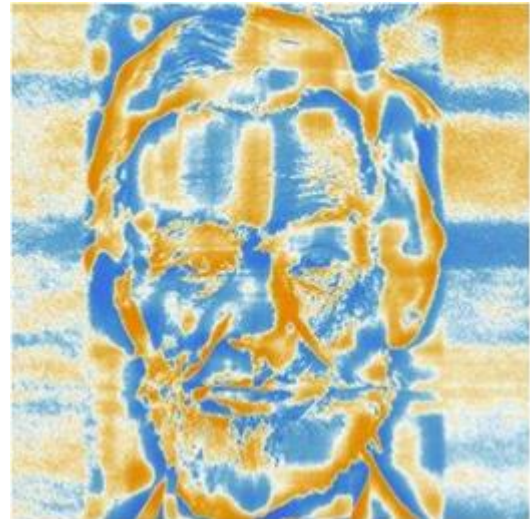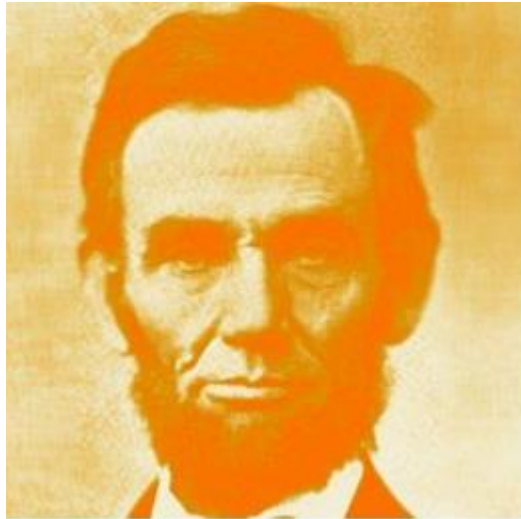$$C = A \cdot B = \sum_k A_k \bigotimes B_k^T$$

Summing only the largest matrix products produces an approximation of M
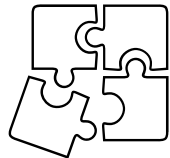
# Error Declines with Dimensionality

# Reconstructing Lincoln

Lincoln's face from 5 and 50 singular values, a substantial compression of the original matrix.

# Wrapup: Linear Algebra

- Linear algebra is a powerful tool for solving complex data science problems
- Matrix operations are particularly useful to work with data that can be represented in matrix form
- There are implementations available that are significantly more efficient than naive ones