# INTRO TO DATA SCIENCE
## LECTURE 15: ADVANCED UNSUPERVISED LEARNING

## I. LDA
## II. LDA EXERCISE WITH PYTHON AND GENSIM
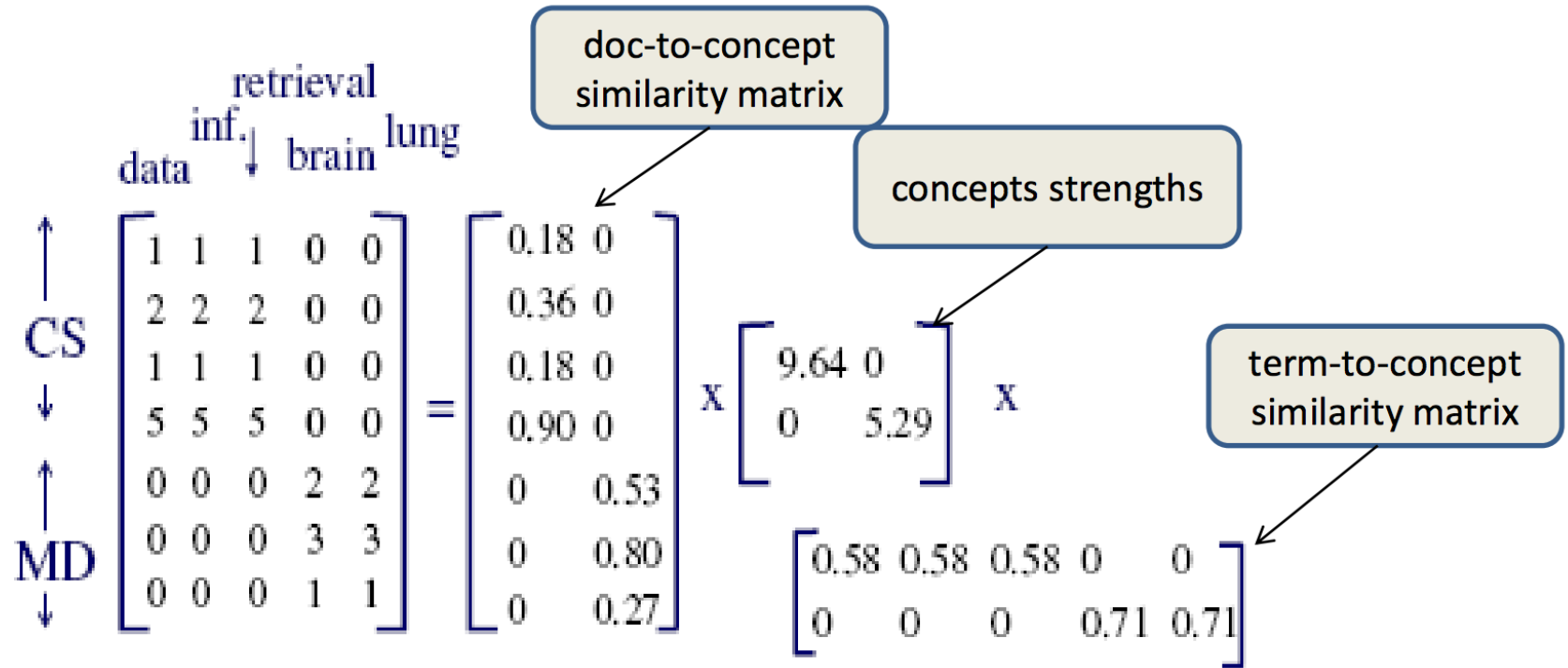
# REVIEW: DIMENSIONALITY REDUCTION

*Q: What is dimensionality reduction?*

*A: A set of techniques for reducing the size (in terms of features, records, and/or bytes) of the dataset under examination.*

*In general, the idea is to regard the dataset as a matrix and to decompose the matrix into simpler, meaningful pieces.*

*Dimensionality reduction is frequently performed as a pre-processing step before another learning algorithm is applied.*

# REVIEW: SINGULAR VALUE DECOMPOSITION

$$
\begin{array}{c}
\text{CS} \\
\\
\text{MD}
\end{array}
\begin{bmatrix}
1 & 1 & 1 & 0 & 0 \\
2 & 2 & 2 & 0 & 0 \\
1 & 1 & 1 & 0 & 0 \\
5 & 5 & 5 & 0 & 0 \\
0 & 0 & 0 & 2 & 2 \\
0 & 0 & 0 & 3 & 3 \\
0 & 0 & 0 & 1 & 1
\end{bmatrix}
=
\begin{bmatrix}
0.18 & 0 \\
0.36 & 0 \\
0.18 & 0 \\
0.90 & 0 \\
0 & 0.53 \\
0 & 0.80 \\
0 & 0.27
\end{bmatrix}
\times
\begin{bmatrix}
9.64 & 0 \\
0 & 5.29
\end{bmatrix}
\times
\begin{bmatrix}
0.58 & 0.58 & 0.58 & 0 & 0 \\
0 & 0 & 0 & 0.71 & 0.71
\end{bmatrix}
$$

Column labels: data, inf., retrieval, brain, lung

doc-to-concept similarity matrix

concepts strengths

term-to-concept similarity matrix

*Consider a matrix* A *with* n *rows and* d *features.*

*The* **singular value decomposition** *of* A *is given by:*

$$A = U \Sigma V^\mathsf{T}$$

(n x d)　　　　(n x k)　(k x k)　(k x d)

*The* **singular value decomposition** *of* A *is given by:*

$$A = U \Sigma V^{\mathsf{T}}$$

(n x d)        (n x k)   (k x k)   (k x d)

*The nonzero entries of* $\Sigma$ *are the* **singular values** *of* A. *These are real, nonnegative, and rank-ordered (decreasing from left to right).*

# I. LATENT DIRICHLET ALLOCATION