# Non-negative matrix factorization

Daniel Hsu

COMS 4772

# Singular value decomposition

- $\boldsymbol{A} = \boldsymbol{U}\boldsymbol{S}\boldsymbol{V}^\top$
  - $\boldsymbol{U}^\top\boldsymbol{U} = \boldsymbol{V}^\top\boldsymbol{V} = \boldsymbol{I}$
  - $\boldsymbol{S} \succ 0$ diagonal
  - Truncations at rank $k$ are optimal for spectral/Frobenius error
- What if we want to add constraints to factors?

# Non-negative matrix factorization (NMF)

- **Given**: $X \in \mathbb{R}^{m \times n}$ non-negative
  - Columns are, e.g., word frequencies of documents, pixel intensities of images.
- **Goal**: factor $X = VY$ where $V \in \mathbb{R}^{m \times r}$ and $Y \in \mathbb{R}^{r \times n}$ have only non-negative entries
  - NP-hard to decide if this is possible (Vavasis, 2007)

# Heuristic (Lee & Seung, 1999)

- Write approximation objective $f(V, Y) := \|X - VY\|_F^2$ as

$$
\begin{aligned}
f(V, Y) &= \sum_{i,j} X_{i,j}^2 - 2 X_{i,j}(UV)_{i,j} + (VY)_{i,j}^2 \\
&= \underbrace{\|X\|_F^2 + \|VY\|_F^2}_{\geq 0} - \underbrace{2\,\mathrm{tr}(X^\top VY)}_{\geq 0} \\
&= f_+(V, Y) - f_-(V, Y) .
\end{aligned}
$$

- **Multiplicative updates** (preserves non-negativity):

$$
V_{i,k} \leftarrow V_{i,k} \cdot \frac{\frac{\partial}{\partial V_{i,k}} f_-(V, Y)}{\frac{\partial}{\partial V_{i,k}} f_+(V, Y)}, \qquad Y_{k,j} \leftarrow Y_{k,j} \cdot \frac{\frac{\partial}{\partial Y_{k,j}} f_-(V, Y)}{\frac{\partial}{\partial Y_{k,j}} f_+(V, Y)}
$$

- Update factor $\geq 1$ iff $f'(V, Y) \leq 0$.
- **Fixed points**: $V = 0$, $Y = 0$, or stationary point of $f$.

# Example



Figure 1: NMF for face images

# Recovery problem

- ▶ Suppose $\boldsymbol{X} = \boldsymbol{V}\boldsymbol{Y}$ for some non-negative $\boldsymbol{V}$ and $\boldsymbol{Y}$ of rank $r$.
    - ▶ Assume (WLOG) rows of $\boldsymbol{X}$, $\boldsymbol{V}$, and $\boldsymbol{Y}$ sum to 1.
    - ▶ Each row of $\boldsymbol{X}$ is a convex combination of rows of $\boldsymbol{Y}$.

- ▶ **Given**: $\boldsymbol{X}$.
- ▶ **Goal**: recover factors $\boldsymbol{V}$ and $\boldsymbol{Y}$.
- ▶ **Separability assumption**: $\boldsymbol{V}$ has positive definite diagonal submatrix.
    - ▶ Ensures uniqueness (Donoho & Stodden, 2003; Arora, Ge, Kannan, & Moitra, 2012)
    - ▶ Each row of $\boldsymbol{Y}$ appears as a row of $\boldsymbol{X}$ (possibly scaled).
    - ▶ (Scaling factor is 1 under assumption that rows of $\boldsymbol{V}$ sum to 1.)

# Recovery algorithm (Arora, Ge, Kannan, & Moitra, 2012)

- **Main idea**: identify the rows of $X$ that are exactly rows of $Y$.
- For each $i = 1, 2, \ldots, m$:
  - If $i$-th row of $X$ is in convex hull of all other rows of $X$, then delete the $i$-th row of $X$
- What remains is exactly $r$ rows of $X$, each being a row of $Y$.

# Application: topic models

- $X_{w,d}$ = number of times word $w$ appears in document $d$
- $V_{w,t} = \Pr(\text{word } w \mid \text{topic } t)$
- $Y_{t,d} \propto \Pr(\text{topic } t \mid \text{document } d)$
- $\mathbb{E}(X) = VY$
- **Separability assumption**: for every topic $t$, there is a word $w_t$ that has non-zero probability in $V$ only under topic $t$.
  - E.g., word "backprop" only appears in documents about topic "machine learning"

- **Goal**: estimate $V$ from documents
  - When model is well-specified,

$$X = VY + \text{zero-mean noise}.$$

# Using co-occurrences (Arora, Ge, & Moitra, 2012)

- ▶ Assume each document has two tokens (i.e., length $\geq 2$)
- ▶ **Bag-of-words assumption with $(V, Y)$ model**: for document $d$,
  - ▶ First token is word $w$ with probability $\sum_t V_{w,t} Y_{t,d}$
  - ▶ Second token is word $w$ with probability $\sum_t V_{w,t} Y_{t,d}$ (independent of first token)

- ▶ **Co-occurrence matrix**: $M_{w,w'} =$ number of documents where first token is $w$ and second token is $w'$.

$$\mathbb{E}(M) = V Y Y^\top V^\top.$$

- ▶ Separability of $V$ can be used with $\mathbb{E}(M)$.
- ▶ If documents are independent, then $M$ is sum of independent random matrices; can exploit matrix concentration to bound $\|M - \mathbb{E}(M)\|_2$.