

Algorithms for Non-negative Matrix Factorization

- By D. D. Lee and H. S. Seung, 2001

Contents

- NMF
- Cost Functions
- Multiplicative update rules
- Interpretation as gradient descent
- Proving the monotonic convergence

NMF

- Nonnegative Matrix: $V \geq 0 \Leftrightarrow \forall i,j V_{ij} \geq 0$
- Non-negative matrix factorization(NMF):Using two non-negative matrices to approximate another non-negative matrix.

$$\begin{cases} V \approx WH (V \in R^{n \times m}, W \in R^{n \times r}, H \in R^{r \times m}) \\ W \geq 0, H \geq 0 \end{cases}$$

- Significance: Relatively few basis vectors are used to represent many data vectors.

Cost Function

- Quantify the quality of approximation.
- Square of the Euclidean distance between A and B

$$\|\mathbf{A} - \mathbf{B}\|^2 = \sum_{ij} (\mathbf{A}_{ij} - \mathbf{B}_{ij})^2 \geq 0$$

- Generalized Kullback-Leibler divergence of A from B

$$D(\mathbf{A} \parallel \mathbf{B}) = \sum_{ij} \left(\mathbf{A}_{ij} \log \frac{\mathbf{A}_{ij}}{\mathbf{B}_{ij}} - \mathbf{A}_{ij} + \mathbf{B}_{ij} \right) \geq 0$$

Optimization Problems

- Problem 1:
$$\begin{cases} \min_{\mathbf{W}, \mathbf{H}} \|\mathbf{V} - \mathbf{WH}\| \\ \mathbf{W}, \mathbf{H} \geq 0 \end{cases}$$
- Problem 2:
$$\begin{cases} \min_{\mathbf{W}, \mathbf{H}} D(\mathbf{V} \parallel \mathbf{WH}) \\ \mathbf{W}, \mathbf{H} \geq 0 \end{cases}$$

How to tackle the problems above?

- They are convex in W only or H only, but not convex in both variables together.
- Goal: Finding local minima

Gradient decent?

- Convergence can be slow
- Sensitive to the step size

Conjugate gradient?

- More complicated to complement

Multiplicative update rules

Euclidean distance

$$H_{a\mu} \leftarrow H_{a\mu} \frac{(W^T V)_{a\mu}}{(W^T W H)_{a\mu}}$$

$$W_{ia} \leftarrow W_{ia} \frac{(V H^T)_{ia}}{(W H H^T)_{ia}}$$

Divergence

$$H_{a\mu} \leftarrow H_{a\mu} \frac{\sum_i W_{ia} V_{i\mu} / (W H)_{i\mu}}{\sum_k W_{ka}}$$

$$W_{ia} \leftarrow W_{ia} \frac{\sum_\mu H_{a\mu} V_{i\mu} / (W H)_{i\mu}}{\sum_v H_{av}}$$

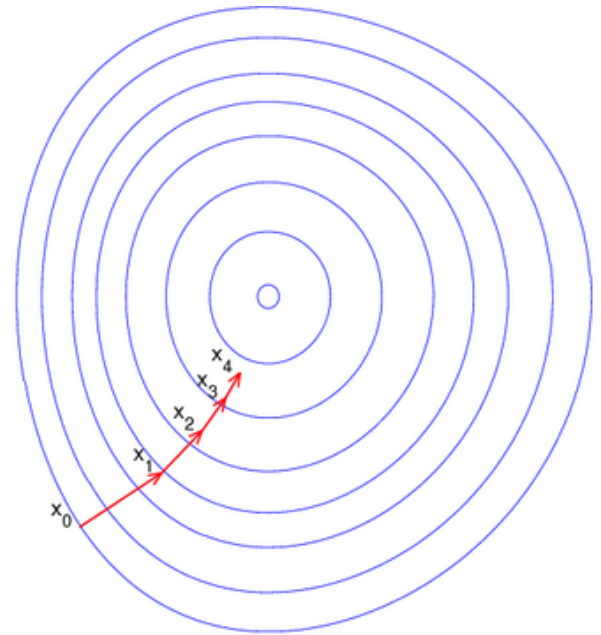
The Euclidean distance and divergence are not increasing under the update rules.

Gradient decent

- Gradient decent for $f(\mathbf{x})$

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \lim_{\eta \rightarrow 0^+} \eta \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}, n \geq 0 \Rightarrow f(\mathbf{x}_0) \geq f(\mathbf{x}_1) \geq f(\mathbf{x}_2) \geq \dots$$

- The sequence $\{\mathbf{x}_n\}$ converges to the desired local minimum.



Multiplicative vs additive update rules

- Euclidean distance

$$\mathbf{H}_{a\mu} \leftarrow \mathbf{H}_{a\mu} + \eta_{a\mu} \left[(\mathbf{W}^T \mathbf{V})_{a\mu} - (\mathbf{W}^T \mathbf{W} \mathbf{H})_{a\mu} \right]$$

$$\eta_{a\mu} = \frac{\mathbf{H}_{a\mu}}{(\mathbf{W}^T \mathbf{W} \mathbf{H})_{a\mu}}$$

$$\mathbf{H}_{a\mu} \leftarrow \mathbf{H}_{a\mu} \frac{(\mathbf{W}^T \mathbf{V})_{a\mu}}{(\mathbf{W}^T \mathbf{W} \mathbf{H})_{a\mu}}$$

- Divergence

$$\mathbf{H}_{a\mu} \leftarrow \mathbf{H}_{a\mu} + \eta_{a\mu} \left[\sum_i \mathbf{W}_{ia} \frac{\mathbf{V}_{i\mu}}{(\mathbf{W} \mathbf{H})_{i\mu}} - \sum_i \mathbf{W}_{ia} \right]$$

$$\eta_{a\mu} = \frac{\mathbf{H}_{a\mu}}{\sum_i \mathbf{W}_{ia}}$$

$$\mathbf{H}_{a\mu} \leftarrow \mathbf{H}_{a\mu} \frac{\sum_i \mathbf{W}_{ia} \mathbf{V}_{i\mu} / (\mathbf{W} \mathbf{H})_{i\mu}}{\sum_k \mathbf{W}_{ka}}$$

- Is it convergent, even if η is not necessarily small enough?

Proofs of converge

- Define an auxiliary function $G(h, h^t)$ for $F(h)$
 $G(\mathbf{h}, \mathbf{h}^t) \geq F(\mathbf{h}), G(\mathbf{h}, \mathbf{h}) = F(\mathbf{h})$

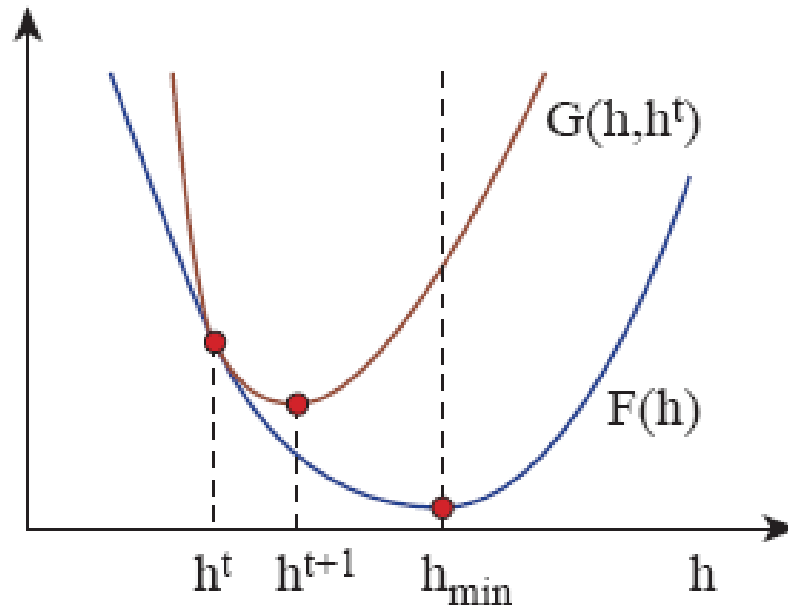
- Find a local minimum of G by iterating the update

$$\mathbf{h}^{t+1} = \arg \min_h G(\mathbf{h}, \mathbf{h}^t)$$

- The sequence converges to a local minimum of $F(h)$

$$F(\mathbf{h}^{t+1}) \leq G(\mathbf{h}^{t+1}, \mathbf{h}^t) \leq G(\mathbf{h}^t, \mathbf{h}^t) \leq F(\mathbf{h}^t)$$

Auxiliary Function



$$\mathbf{h}^{t+1} = \arg \min_h G(\mathbf{h}, \mathbf{h}^t)$$

Updates for H of Euclidean distance

$$F(\mathbf{h}) = \frac{1}{2} \sum_i \left(\mathbf{v}_i - \sum_a \mathbf{W}_{ia} \mathbf{h}_a \right)^2$$

$$G(\mathbf{h}, \mathbf{h}^t) = F(\mathbf{h}^t) + (\mathbf{h} - \mathbf{h}^t)^T \nabla F(\mathbf{h}^t) + \frac{1}{2} (\mathbf{h} - \mathbf{h}^t)^T \mathbf{K}(\mathbf{h}^t) (\mathbf{h} - \mathbf{h}^t)$$

$$K_{ab}(\mathbf{h}^t) = \delta_{ab} (\mathbf{W}^T \mathbf{W} \mathbf{h}^t)_a / \mathbf{h}_a^t$$

Proving steps:

Step one

- Show $G(\mathbf{h}, \mathbf{h}^t)$ is an auxiliary function for $F(\mathbf{h})$

Step two

- Obtain the update rule by setting the gradient of $G(\mathbf{h}, \mathbf{h}^t)$ to zero

Step three

- Check the equivalence between the update rule and $H_{a\mu} \leftarrow H_{a\mu} \frac{(\mathbf{W}^T \mathbf{V})_{a\mu}}{(\mathbf{W}^T \mathbf{W} \mathbf{H})_{a\mu}}$

Auxiliary function $G(\mathbf{h}, \mathbf{h}^t)$ for $F(\mathbf{h})$

- $F(\mathbf{h}) = \frac{1}{2} \sum_i (\mathbf{v}_i - \sum_a \mathbf{W}_{ia} \mathbf{h}_a)^2$
- $F(\mathbf{h}) = F(\mathbf{h}^t) + (\mathbf{h} - \mathbf{h}^t)^T \nabla F(\mathbf{h}^t) + \frac{1}{2} (\mathbf{h} - \mathbf{h}^t)^T (\mathbf{W}^T \mathbf{W}) (\mathbf{h} - \mathbf{h}^t)$
- $G(\mathbf{h}, \mathbf{h}^t) = F(\mathbf{h}^t) + (\mathbf{h} - \mathbf{h}^t)^T \nabla F(\mathbf{h}^t) + \frac{1}{2} (\mathbf{h} - \mathbf{h}^t)^T \mathbf{K}(\mathbf{h}^t) (\mathbf{h} - \mathbf{h}^t)$

$$G(\mathbf{h}, \mathbf{h}^t) \geq F(\mathbf{h}) \Leftrightarrow (\mathbf{h} - \mathbf{h}^t)^T [\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}] (\mathbf{h} - \mathbf{h}^t) \geq 0$$

$$\mathbf{M}_{ab}(\mathbf{h}^t) = \mathbf{h}_a^t [\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}]_{ab} \mathbf{h}_b^t$$

$$\begin{aligned} \mathbf{v}^T \mathbf{M} \mathbf{v} &= \sum_{ab} \mathbf{v}_a \mathbf{M}_{ab} \mathbf{v}_b = \sum_{ab} \mathbf{h}_a^t (\mathbf{W}^T \mathbf{W})_{ab} \mathbf{h}_b^t \mathbf{v}_a^2 - \mathbf{v}_a \mathbf{h}_a^t (\mathbf{W}^T \mathbf{W})_{ab} \mathbf{h}_b^t \mathbf{v}_b \\ &= \sum_{ab} (\mathbf{W}^T \mathbf{W})_{ab} \mathbf{h}_a^t \mathbf{h}_b^t \left[\frac{1}{2} \mathbf{v}_a^2 + \frac{1}{2} \mathbf{v}_b^2 - \mathbf{v}_a \mathbf{v}_b \right] \\ &= \frac{1}{2} \sum_{ab} (\mathbf{W}^T \mathbf{W})_{ab} \mathbf{h}_a^t \mathbf{h}_b^t (\mathbf{v}_a + \mathbf{v}_b)^2 \geq 0 \end{aligned}$$

Minimum of $G(\mathbf{h}, \mathbf{h}^t)$ and update rules

$$\frac{\partial G(\mathbf{h}, \mathbf{h}^t)}{\partial \mathbf{h}} = \nabla F(\mathbf{h}^t) + (\mathbf{h} - \mathbf{h}^t)K(\mathbf{h}^t) = 0 \Rightarrow \mathbf{h}^{t+1} = \mathbf{h}^t - K(\mathbf{h}^t)^{-1} \nabla F(\mathbf{h}^t)$$

$$\mathbf{h}_a^{t+1} = \mathbf{h}_a^t \frac{(\mathbf{W}^T \mathbf{v})_a}{(\mathbf{W}^T \mathbf{W} \mathbf{h}^t)_a}$$

$$\mathbf{H}_{a\mu} \leftarrow \mathbf{H}_{a\mu} \frac{(\mathbf{W}^T \mathbf{v})_{a\mu}}{(\mathbf{W}^T \mathbf{W} \mathbf{H})_{a\mu}}$$

By reversing the role of \mathbf{H} and \mathbf{W} , F can be shown to be non-increasing under the update rules of for \mathbf{W} .

Updates for H of divergence

$$F(\mathbf{h}) = \sum_i v_i \log \left(\frac{v_i}{\sum_a W_{ia} h_a} \right) - v_i + \sum_a W_{ia} h_a$$

$$G(\mathbf{h}, \mathbf{h}^t) = \sum_i (v_i \log v_i - v_i) + \sum_{ia} W_{ia} h_a - \sum_{ia} v_i \frac{W_{ia} h_a^t}{\sum_b W_{ib} h_b^t} \left(\log W_{ia} h_a - \log \frac{W_{ia} h_a^t}{\sum_b W_{ib} h_b^t} \right)$$

Step one

- Show $G(\mathbf{h}, \mathbf{h}^t)$ is an auxiliary function for $F(\mathbf{h})$

Step two

- Obtain the update rule by setting the gradient of $G(\mathbf{h}, \mathbf{h}^t)$ to zero

Step three

- Check the equivalence between the update rule and $H_{a\mu} \leftarrow H_{a\mu} \frac{\sum_i W_{ia} v_i / (WH)_{i\mu}}{\sum_k W_{ka}}$

Proof of $G(\mathbf{h}, \mathbf{h}^t)$ by convexity

$$F(\mathbf{h}) = \sum_i \mathbf{v}_i \log \left(\frac{\mathbf{v}_i}{\sum_a \mathbf{W}_{ia} \mathbf{h}_a} \right) - \mathbf{v}_i + \sum_a \mathbf{W}_{ia} \mathbf{h}_a$$

$$G(\mathbf{h}, \mathbf{h}^t) = \sum_i (\mathbf{v}_i \log \mathbf{v}_i - \mathbf{v}_i) + \sum_{ia} \mathbf{W}_{ia} \mathbf{h}_a - \sum_{ia} \mathbf{v}_i \frac{\mathbf{W}_{ia} \mathbf{h}_a^t}{\sum_b \mathbf{W}_{ib} \mathbf{h}_b^t} \left(\log \mathbf{W}_{ia} \mathbf{h}_a - \log \frac{\mathbf{W}_{ia} \mathbf{h}_a^t}{\sum_b \mathbf{W}_{ib} \mathbf{h}_b^t} \right)$$

- The convexity of log function

$$-\log \sum_a \mathbf{W}_{ia} \mathbf{h}_a = -\log \sum_a \alpha_a \frac{\mathbf{W}_{ia} \mathbf{h}_a}{\alpha_a} \leq -\sum_a \alpha_a \log \frac{\mathbf{W}_{ia} \mathbf{h}_a}{\alpha_a}$$

$$\alpha_a = \frac{\mathbf{W}_{ia} \mathbf{h}_a^t}{\sum_b \mathbf{W}_{ib} \mathbf{h}_b^t}$$

$$-\log \sum_a \mathbf{W}_{ia} \mathbf{h}_a \leq -\sum_a \frac{\mathbf{W}_{ia} \mathbf{h}_a^t}{\sum_b \mathbf{W}_{ib} \mathbf{h}_b^t} \left(\log \mathbf{W}_{ia} \mathbf{h}_a - \log \frac{\mathbf{W}_{ia} \mathbf{h}_a^t}{\sum_b \mathbf{W}_{ib} \mathbf{h}_b^t} \right) \Rightarrow G(\mathbf{h}, \mathbf{h}^t) \geq F(\mathbf{h})$$

Minimum of $G(\mathbf{h}, \mathbf{h}^t)$ and update rules

$$\frac{\partial G(\mathbf{h}, \mathbf{h}^t)}{\partial \mathbf{h}_a} = \sum_i \mathbf{W}_{ia} - \frac{1}{\mathbf{h}_a} \sum_i \mathbf{v}_i \frac{\mathbf{W}_{ia} \mathbf{h}_a^t}{\sum_b \mathbf{W}_{ib} \mathbf{h}_b^t} = 0$$

$$\mathbf{h}^{t+1} = \frac{\mathbf{h}_a^t}{\sum_i \mathbf{W}_{ia}} \sum_i \frac{\mathbf{v}_i}{\sum_b \mathbf{W}_{ib} \mathbf{h}_b^t} \mathbf{W}_{ia}$$

$$\mathbf{H}_{a\mu} \leftarrow \mathbf{H}_{a\mu} \frac{\sum_i \mathbf{W}_{ia} \mathbf{v}_{i\mu} / (\mathbf{WH})_{i\mu}}{\sum_k \mathbf{W}_{ka}}$$

- By reversing the role of \mathbf{H} and \mathbf{W} , the update rule for \mathbf{W} can similarly be shown to be non-increasing.