

Matrix Factorization

EE698V - Machine Learning for Signal Processing

Vipul Arora



References (highly recommended)

- Lee, Daniel D., and H. Sebastian Seung, ‘Learning the Parts of Objects by Non-Negative Matrix Factorization’, *Nature*, 401 (1999), 788–91
- Lee, D D, and H S Seung, ‘Algorithms for Non-Negative Matrix Factorization’, in *Proc. Conf. Adv. Neural Info. Process. Syst. (NIPS)* (MIT Press, 2001), XIII, 556–62

Matrix Factorization

1.00	0.0	1.00	0.00	0.00	1.00	0.00	0.0
0.75	0.0	0.75	0.25	0.25	0.75	0.25	0.0
0.50	0.0	0.50	0.50	0.50	0.50	0.50	0.0
0.25	0.0	0.25	0.75	0.75	0.25	0.75	0.0
0.00	0.0	0.00	1.00	1.00	0.00	1.00	0.0

Matrix Factorization

1.00	0.0	1.00	0.00	0.00	1.00	0.00	0.0
0.75	0.0	0.75	0.25	0.25	0.75	0.25	0.0
0.50	0.0	0.50	0.50	0.50	0.50	0.50	0.0
0.25	0.0	0.25	0.75	0.75	0.25	0.75	0.0
0.00	0.0	0.00	1.00	1.00	0.00	1.00	0.0

V

$$V = W \cdot H$$

=

0.00	1.00
0.25	0.75
0.50	0.50
0.75	0.25
1.00	0.00

x

0	0	0	1	1	0	1	0
1	0	1	0	0	1	0	0

H
(activations)

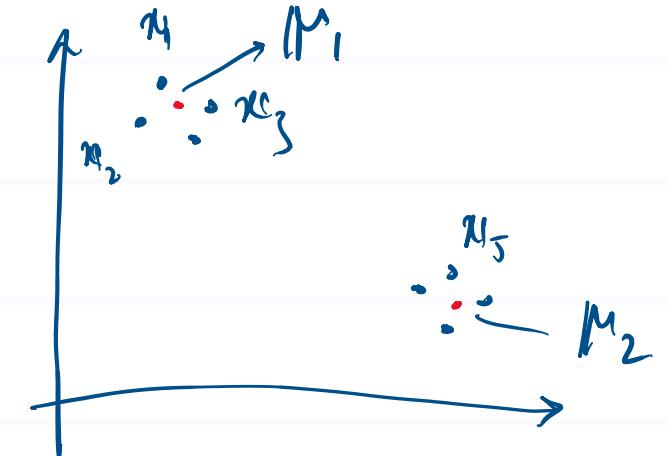
W
(dictionary)

Matrix Factorization: Advantages

- Reduction in size:
 - V : (M, N)
 - W : (M, R)
 - H : (R, N)
 - If R is small, then MR+RN << MN
- Useful for analysis
 - Dimension reduction
 - Classification
- Useful for reconstruction

K-means or Vector Quantization

1.22	0.20	1.00	0.00	0.00	1.05	0.00	0.00
0.83	0.04	0.67	0.35	0.31	0.66	0.34	0.11
0.42	0.04	0.60	0.60	0.60	0.41	0.63	0.00
0.39	0.03	0.45	0.73	0.84	0.21	0.78	0.07
0.20	0.08	0.01	0.69	1.06	0.07	1.11	0.00



$$x_1 \rightarrow \mu_1$$

$$x_5 \rightarrow \mu_2$$

$$x_3 \rightarrow \mu_1$$

- Winner takes all

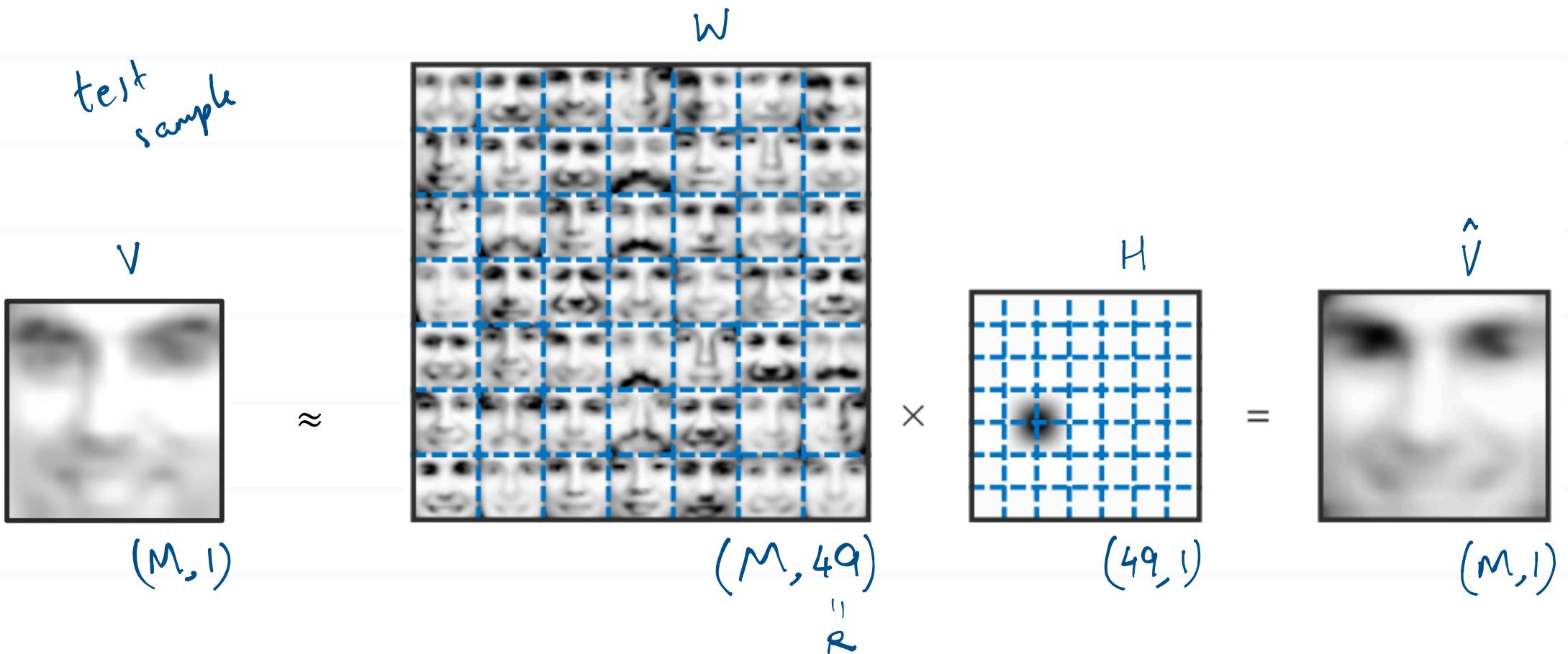
K-means or Vector Quantization

1.22	0.20	1.00	0.00	0.00	1.05	0.00	0.00
0.83	0.04	0.67	0.35	0.31	0.66	0.34	0.11
0.42	0.04	0.60	0.60	0.60	0.41	0.63	0.00
0.39	0.03	0.45	0.73	0.84	0.21	0.78	0.07
0.20	0.08	0.01	0.69	1.06	0.07	1.11	0.00

$$\begin{matrix} & = & \begin{matrix} 0.00 & 1.00 \\ 0.25 & 0.75 \\ 0.50 & 0.50 \\ 0.75 & 0.25 \\ 1.00 & 0.00 \end{matrix} & \times & \begin{matrix} 0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \end{matrix} & H \\ & & \omega & & & & & H \rightarrow S \end{matrix}$$

Vector Quantization: Example

learnt from 1500 samples



PCA as Matrix Factorization

$$y_n = U^T x_n \quad \forall n$$

dict activ.
↓ ↓

$$U = [u_1 \ u_2 \dots \ u_m]$$

(eigenvecs of $S = \frac{1}{N} \sum_n x_n x_n^T$)

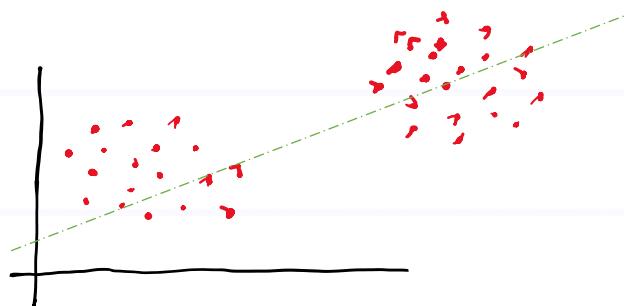
In matrix form: $[y_1 \ y_2 \ \dots \ y_N] = U^T [x_1 \ x_2 \ \dots \ x_N]$

$$\underbrace{Y}_{M \times N} = \underbrace{U^T}_{M \times 784} \underbrace{X}_{784 \times N}$$

all u_i are
orthonormal

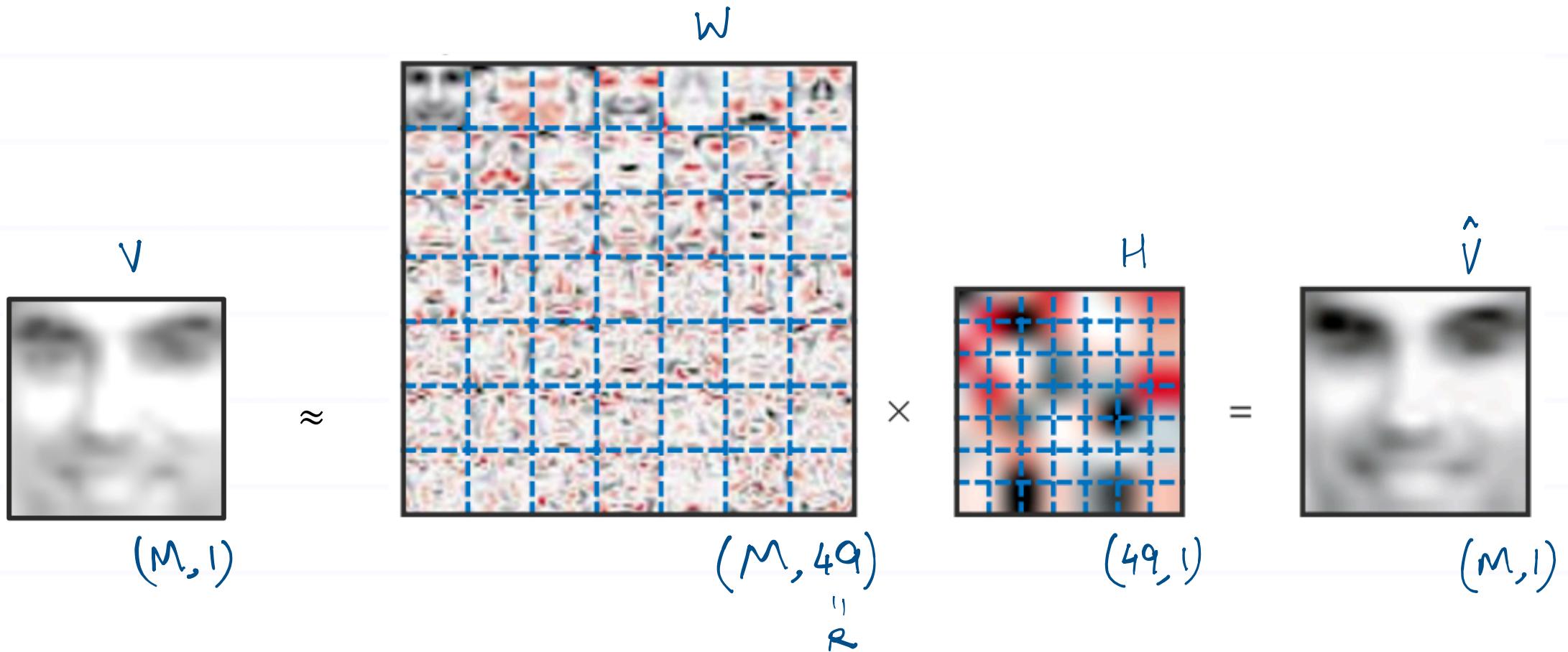
- This is a way of matrix factorization
- Here, Y, U, X can have any real values
- But for images, negative values do not make sense.
- Can we constrain U to have only non-negative values? Then Y will also be non-negative

PCA as Matrix Factorization



- Dimensionality reduction: R eigen vectors

PCA: Example



Non Negative Matrix Factorization (NMF)

- Many times matrices are non-negative
 - Images
 - Spectral magnitudes
 - Energies
 - Counts (NLP)
- Non-negative dictionary and activations give interpretability and dexterity; dictionary learns different parts
- $V \geq 0, W \geq 0, H \geq 0$

NMF: Example

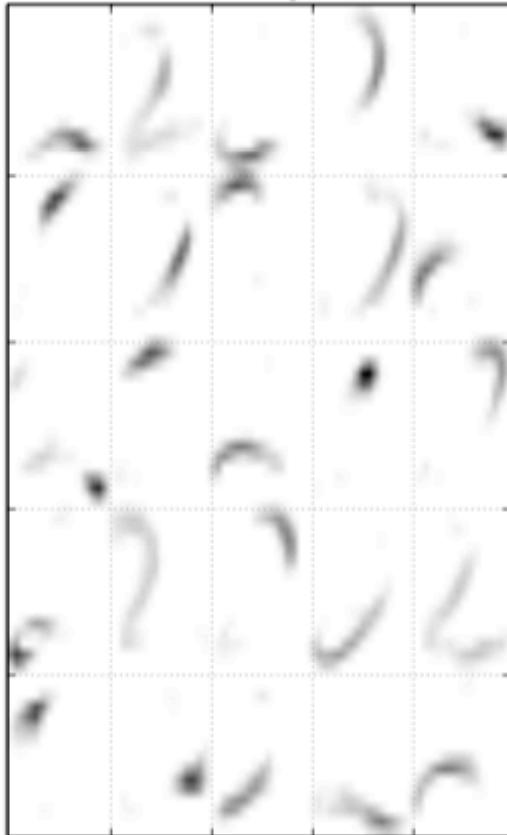
$$v \approx w \times h = \hat{v}$$

Diagram illustrating Non-negative Matrix Factorization (NMF) decomposition:

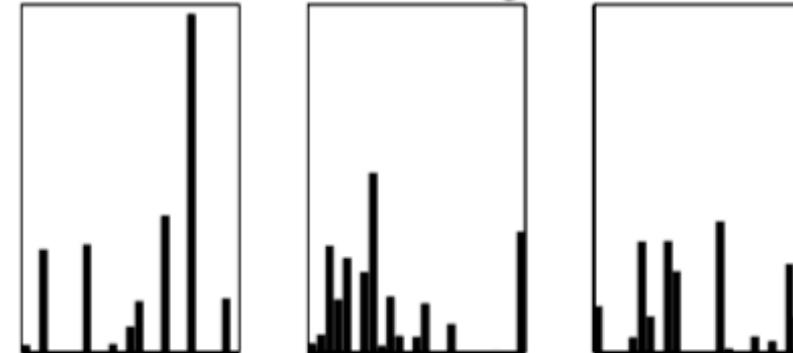
- Input:** A grayscale image v of size $(M, 1)$.
- Approximation:** The input v is approximated by the product of two matrices: w and h .
- Matrix w :** A grayscale image of size $(M, 49)$. It is shown with a blue dashed grid overlay, indicating it has 7 columns and 7 rows.
- Matrix h :** A grayscale image of size $(49, 1)$. It is shown with a blue dashed grid overlay, indicating it has 49 columns and 1 row.
- Output:** The product $w \times h$ results in a reconstructed image \hat{v} of size $(M, 1)$, which visually appears very similar to the original input v .

NMF: Another Example

Basis Components



Mixture Weights



Data Vectors



NMF

- V is known
- W, H are unknowns
- Estimate them in tandem
 - Given V , initialize W
 - Estimate H , given W
 - Estimate W , given H
 - Iterate until convergence

NMF: parameter estimation

$$\begin{aligned} E &= \|V - \hat{V}\|^2 \\ &= \|V - WH\|^2 \quad \text{constraint: } W \geq 0, H \geq 0 \end{aligned}$$

or $E = \sum_{ij} \left| V_{ij} - \sum_r w_{ir} h_{rj} \right|^2 ; \quad w_{ir} \geq 0, h_{rj} \geq 0 \quad \forall i, r, j$

To optimize:

- $\frac{\partial E}{\partial W \text{ or } H} = 0$

- grad. desc.

NMF: parameter estimation

- $W_{ir} \leftarrow W_{ir} \frac{(VH^T)_{ir}}{(WHH^T)_{ir}}$
- $H_{rj} \leftarrow H_{rj} \frac{(W^TV)_{rj}}{(W^TWH)_{rj}}$
- See derivation in Reference: Lee, D D, and H S Seung, ‘Algorithms for Non-Negative Matrix Factorization’, in *Proc. Conf. Adv. Neural Info. Process. Syst. (NIPS)* (MIT Press, 2001), xiii, 556–62

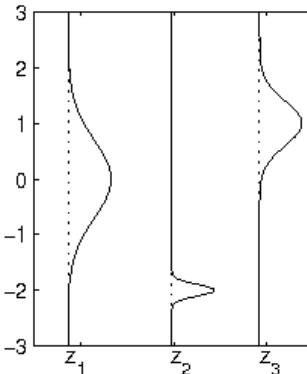
Probabilistic Latent Component Analysis (PLCA): one kind of NMF

Reference:

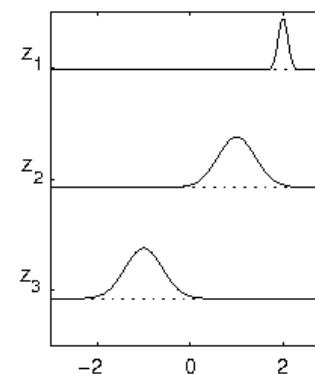
Shashanka, M, 'Latent Variable Framework for Modeling and Separating Single Channel Acoustic Sources' (Boston University, 2007) **Section 3.3.2**

PLCA

W



H



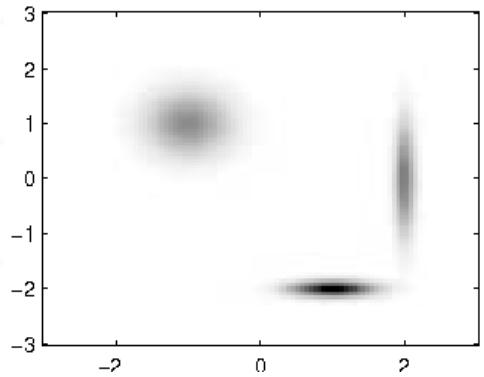
(original)

$$V \approx WH$$

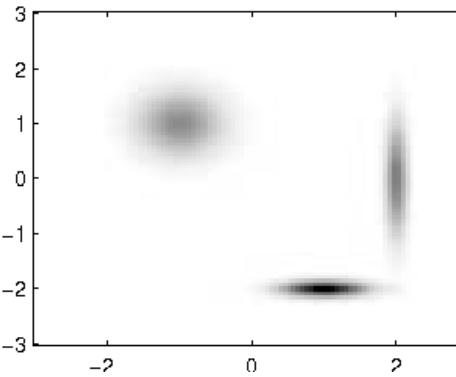
$$= \hat{V}$$

(approx.)

V

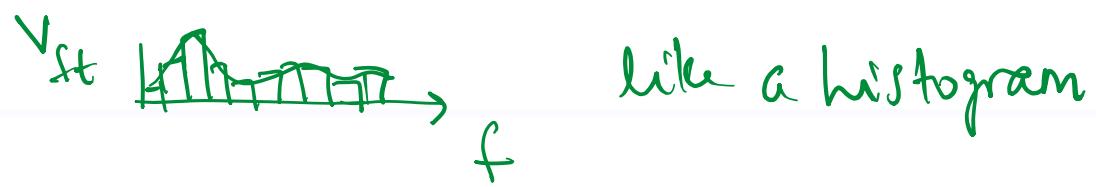


\hat{V}



(original)

PLCA



Given V_{ft}

$$f \begin{array}{|c|} \hline t \\ \hline \end{array} = f \begin{array}{|c|} \hline z \\ \hline \end{array} \times \begin{array}{|c|} \hline z \\ \hline t \\ \hline \end{array}$$

Let $V_{ft} \sim P(f|t)$

$$P(f|z) \quad P(z|t)$$

We want to approximate $P(f|t)$ with PLCA
using Latent variable z (discrete r.v.)

$$\begin{aligned} P(f|t) &= \sum_z P(z|t) \underbrace{P(f|z,t)}_{= P(f|z)} \quad \text{independent of } t \end{aligned}$$

PLCA

(t is not a r.v. here)

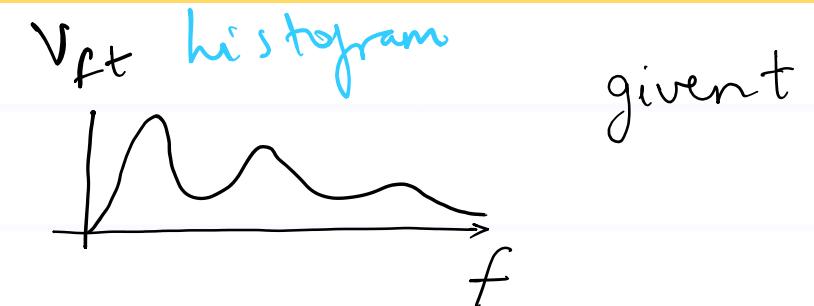
$$v_{ft} \sim p_t(f) = \sum_z p(f|z) p_t(z)$$

PLCA

- AIM: Given V_{ft} , we want to estimate the underlying distribution with the help of latent variables
- Analogy:
 - underlying clusters
 - underlying classes
 - parts

Maximum likelihood Estimation

$$\text{likelihood} = \prod_i P_t(f_i)$$



$$\log \text{likelihood} = \sum_i \log P_t(f_i)$$

$$= \sum_f V_{ft} \log P_t(f)$$

\because counts $\propto V_{ft}$

$f_i \leftarrow f$

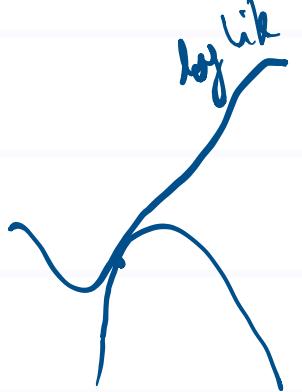
MLE

Over all t ,

$$\begin{aligned}\text{log likelihood} &= \sum_{f,t} V_{ft} \log P_t(f) \\ &= \sum_{f,t} V_{ft} \log \left[\sum_z P_t(z) P(f|z) \right]\end{aligned}$$

MLE

Over all t ,



$$\begin{aligned}
 \text{log likelihood} &= \sum_{f,t} V_{ft} \log P_t(f) \\
 &= \sum_{f,t} V_{ft} \log \left[\underbrace{\sum_z \frac{P_t(z) P(f|z)}{q_t(z|f)}}_{\text{pdf}} \right] \times q_t(z|f) \\
 &\geq \sum_{f,t} V_{ft} \sum_z q_t(z|f) \log \left[\frac{P_t(z) P(f|z)}{q_t(z|f)} \right]
 \end{aligned}$$

$$\therefore \text{Aux. fn., } \mathcal{L}(q, \theta) = \sum_{f,t} V_{ft} \sum_z q_t(z|f) \log \left[\frac{P_t(z) P(f|z)}{q_t(z|f)} \right]$$

Recall: EM Algorithm

- E-step

- $q \leftarrow \operatorname{argmin}_q \text{KL} [q(Z) || p(Z|X, \theta)]$

$$q(z) = p(z|x, \theta)$$

- M-step

- $\theta \leftarrow \operatorname{argmax}_{\theta} \mathbb{E}_q [\log p(X, Z | \theta)]$

analytically solvable

EM Algorithm

$$E\text{-Step: } q_t(z|f) = \frac{P_t(z) P(f|z)}{\sum_z P_t(z) P(f|z)}$$

$$M\text{-step: } \theta = \arg \max_{\theta} \mathcal{L}(q, \theta)$$

M-step

$$\mathcal{L}(q, \theta) = \sum_{f,t} V_{ft} \sum_z q_t(z|f) \log \left[\frac{P_t(z) P(f|z)}{q_t(z|f)} \right]$$

θ includes $\underbrace{P_t(z)}_H$ and $\underbrace{P(f|z)}_W$

Constraints:

$$\sum_z P_t(z) = 1 \quad \& \quad \sum_f P(f|z) = 1$$

M-step

$$\mathcal{L}(q, \theta) = \sum_{f,t} V_{ft} \sum_z q_t(z|f) \log \left(P_t(z) P(f|z) \right) + \sum_t \lambda_t \left(1 - \sum_z P_t(z) \right)$$

$$+ \sum_z \lambda'_z \left(1 - \sum_f P(f|z) \right)$$

$$\frac{\partial \mathcal{L}}{\partial P_{t'}(z')} = \sum_f V_{ft'} q_{t'}(z'|f) \cdot \frac{1}{P_{t'}(z')} + \lambda_{t'} (-1) = 0$$

$$\Rightarrow \sum_f V_{ft'} q_{t'}(z'|f) = \lambda_{t'} P_{t'}(z')$$

Summing over z' gives $\lambda_{t'}$

1

$$P_{t'}(z') = \frac{\sum_f V_{ft'} q_{t'}(z'|f)}{\lambda_{t'}}$$

M-step

∴

$$P_t(z) = \frac{\sum_f V_{ft} q_t(z|f)}{\sum_{z'} \sum_f V_{ft} q_t(z'|f)}$$

Similarly,

$$P(f|z) = \frac{\sum_t V_{ft} q_t(z|f)}{\sum_f \sum_t V_{ft} q_t(z|f)}$$

