

NON-NEGATIVE MATRIX FACTORIZATION (PAATERO & TANNER, LEE & SEUNG)

SETTING: GIVEN N OBSERVATIONS OF P NON-NEGATIVE EXAMPLES

IN DATA MATRIX

$$X = \begin{bmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_N^T \end{bmatrix}$$

NOTE: $x_{ij} \geq 0$ IS MOTIVATED BY EXPERIMENTAL DATA

- BIOLOGY / GENETICS: CELL COUNTS, PROTEIN EXPRESSION
- IMAGING: PIXEL INTENSITY $0 \rightarrow 1$ (WHITE \rightarrow BLACK)
SAME FOR RGB IN COLOR IMAGES
- TOPIC MODELING: COUNT OR FREQUENCY OF TERMS APPEARING IN A DOCUMENT

GOAL: FIND A MATRIX $\hat{X} \in \mathbb{R}^{N \times P}$ WITH $\text{RANK}(\hat{X}) \ll \text{RANK}(X)$ THAT MINIMIZES THE LOSS OF INFORMATION CONTAINED IN X .

QUESTION 1: HOW DO WE MEASURE INFORMATION LOSS?

1) FROEBINUS NORM $\|X - \hat{X}\|_F = \sqrt{\sum_{n=1}^N \sum_{j=1}^P |x_{nj} - \hat{x}_{nj}|^2}$

2) DIVERGENCE $D(X \| \hat{X}) = \sum_{n=1}^N \sum_{j=1}^P \left[x_{nj} \log \frac{x_{nj}}{\hat{x}_{nj}} + \hat{x}_{nj} - x_{nj} \right]$

i) IN THE CASE $\sum_n \sum_j x_{nj} = \sum_n \sum_j \hat{x}_{nj} = 1$, $D(X \| \hat{X})$ IS CALLED THE KULLBACK-LIEBLER DIVERGENCE, OR RELATIVE ENTROPY OF X FROM \hat{X} .

ii) $D(X \| \hat{X}) = \underbrace{-\sum_n \sum_j \left(x_{nj} \log(\hat{x}_{nj}) - \hat{x}_{nj} \right)}_{\text{LOG-LIKELIHOOD OF } x_{nj} \sim \text{POISSON}(\hat{x}_{nj})} + \sum_n \sum_j x_{nj} \log(x_{nj})$

THM: LET X BE A DATA MATRIX D/O SVD OF X

$$X = \underbrace{\begin{bmatrix} \vec{u}_1 & \dots & \vec{u}_N \end{bmatrix}}_U \underbrace{\begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_p \\ & & & 0 \end{bmatrix}}_D \underbrace{\begin{bmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_p^T \\ \vdots \end{bmatrix}}_{V^T}$$

Fix $1 \leq k < p$ AND SET

$$\hat{X}_k = U \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ & & \sigma_k & & 0 \\ & & 0 & \ddots & \\ & & & & 0 \end{bmatrix} V^T$$

Then

1) \hat{X}_k IS RANK k AND

2) $\hat{X}_k = \underset{\substack{Y \in \mathbb{R}^{N \times p} \\ \text{RANK}(Y)=k}}{\text{argmin}} \|X - Y\|_F$

NOTE:

1) THE ABOVE THM, SAYS THAT PCA GIVES THE ANSWER WHEN WE USE $\|\cdot\|_F$

2) THIS HAS AN INTERPRETABILITY PROBLEM! RECALL

$$\hat{X}_k = \begin{bmatrix} \vec{x}_1^T \\ \vdots \\ \vec{x}_N^T \end{bmatrix} \begin{bmatrix} \vec{v}_1 & \dots & \vec{v}_k \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_k^T \end{bmatrix}$$

$\vec{v}_1, \dots, \vec{v}_k$ ARE FIRST k PRIN COMP LOADINGS

$$\Rightarrow \vec{x}_i^{(k)} = \sum_{s=1}^k \underbrace{\vec{x}_i^T \vec{v}_s}_{y_{is}} \vec{v}_s^T$$

$$\Rightarrow x_{ij} \approx \sum_{s=1}^k y_{is} v_{sj}$$

RECALL, $x_{ij} \geq 0$ IN THIS SETTING BUT

PCA MAKES NO SUCH CONSTRAINT ON SCORES OR LOADINGS!

$$\vec{x}_i \approx \sum_{s=1}^k y_{is} \vec{v}_s$$

\nwarrow WEIGHTS \swarrow FEATURES/COMPONENTS

FROM THIS VIEW

- 1) EACH \vec{x}_i IS BUILT FROM THE SAME LIBRARY OF FEATURES $\vec{v}_1, \dots, \vec{v}_k$
- 2) IF $y_{is} < 0$ WE ARE 'SUBTRACTING' FEATURE \vec{v}_s

PHYSICAL INTERPRETATION?

- 3) SUPPOSE $p = 2$

$$x_{ij} = \sum_{s=1}^2 y_{is} v_{sj}$$

$y_{i1} v_{1j} \approx 0 \approx y_{i2} v_{2j} \Rightarrow x_{ij} \approx 0$
 $y_{i1} w_{1j} \gg 0$
 $y_{i2} w_{2j} \ll 0 \Rightarrow x_{ij} \approx 0$

HOW DO WE DIFFERENTIATE BETWEEN

- a) ABSENCE OF FEATURE $w_{1j} = w_{2j} = 0$
- AND b) CANCELLATION OF FEATURE? $w_{1j} = -w_{2j}$

\Rightarrow FEATURES $\vec{w}_1, \dots, \vec{w}_k$ CANNOT BE EXPECTED TO HAVE A MEANINGFUL PHYSICAL INTERPRETATION AS 'BUILDING BLOCKS' OF OUR DATA!

\Rightarrow ENFORCE NON-NEGATIVITY OF ELEMENTS OF FACTORIZATION

~~$$\vec{X} \approx \hat{\vec{X}} = W \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \\ & & & 0 \end{bmatrix} V^T = Y V^T \quad (\text{PCA})$$~~

$$= WH = \begin{bmatrix} w_1^T \\ \vdots \\ w_N^T \end{bmatrix} \begin{bmatrix} h_1^T \\ \vdots \\ h_k^T \end{bmatrix}$$

$w_1, \dots, w_N \in \mathbb{R}^{1 \times k}$ WEIGHTS
 $\vec{h}_1, \dots, \vec{h}_k \in \mathbb{R}^{p \times k}$ FEATURES
 $w_{ij}, h_{ij} \geq 0$

SO THAT

$$\vec{x}_i \approx \sum_{s=1}^k w_{is} \vec{h}_s$$

IS A SUPERIMPOSITION OF FEATURES $\vec{h}_1, \dots, \vec{h}_s$.

GEOMETRIC INTERPRETATION

Ex: $\vec{x}_1, \dots, \vec{x}_N \in \mathbb{R}^3$ AND THERE EXISTS A $W \in \mathbb{R}^{N \times 2}$; $H \in \mathbb{R}^{2 \times 3}$

s.t. $X = WH$

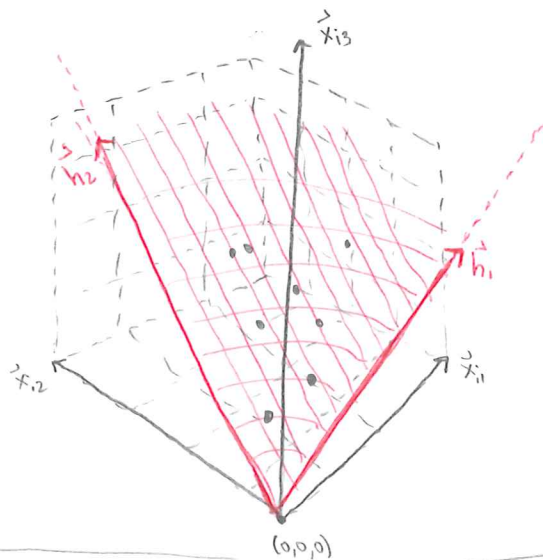
- 1) NON-NEGATIVITY OF x_{ij} IMPLIES $\vec{x}_1, \dots, \vec{x}_N$
LIE IN THE POSITIVE ORTHANT OF \mathbb{R}^3

$$P_3 = \{ \vec{x} \in \mathbb{R}^3, x_j \geq 0, j=1, \dots, 3 \}$$

- 2) EQUALITY $WH = X$ IMPLIES

$$\vec{x}_i = \sum_{s=1}^2 w_{is} \vec{h}_s \in \text{SPAN} \{ \vec{h}_1, \vec{h}_2 \}$$

i.e. DATA RESIDE ON A 2-D SUBSPACE



DEF: THE SIMPLICIAL CONE GENERATED BY THE SET OF VECTORS $\{ \vec{h}_1, \dots, \vec{h}_s \}$ IN \mathbb{R}^p IS THE SET

$$\Gamma(\{ \vec{h}_1, \dots, \vec{h}_s \}) = \left\{ \vec{v} \in \mathbb{R}^p \mid \vec{v} = \sum_{j=1}^s a_j \vec{h}_j, a_1, \dots, a_s \geq 0 \right\}$$

ALSO CALLED THE POSITIVE SPAN OF $\{ \vec{h}_1, \dots, \vec{h}_s \}$

CONICAL HULL OF $\{ \vec{h}_1, \dots, \vec{h}_s \}$

NOTE: WE CAN CHOOSE DIFFERENT STEP SIZES FOR EACH COORDINATE

for $(j=1, \dots, p)$ {

$$x_j^{(i+1)} = x_j^{(i)} - \Delta_{x_j} \nabla f \begin{pmatrix} x_1^{(i+1)} \\ \vdots \\ x_{j-1}^{(i+1)} \\ x_j^{(i)} \\ \vdots \\ x_p^{(i)} \end{pmatrix}$$

}

AND ITERATE $i=1, \dots$ UNTIL CONVERGENCE

Ex: NMF : DIVERGENCE

$$\text{ERROR} = D(X \parallel WH) = \sum_{s=1}^n \sum_{t=1}^p \left[X_{st} \log \frac{X_{st}}{(WH)_{st}} + (WH)_{st} - X_{st} \right]$$

ASSUME WE HAVE STEP SIZES $\Delta_{s\ell}^W, \Delta_{\ell t}^H$ $1 \leq s \leq n$
 $1 \leq \ell \leq k$
 $1 \leq t \leq p$

1) FIND $\frac{\partial}{\partial w_{ij}} D(X \parallel WH)$

$$\frac{\partial}{\partial w_{ij}} D(X \parallel WH) = \frac{\partial}{\partial w_{ij}} \left(\sum_{st} [-X_{st} \log (WH)_{st} + (WH)_{st}] + \sum_{st} [X_{st} \log X_{st} + X_{st}] \right)$$

DOES NOT DEPEND ON W OR H!

$$(WH)_{st} = \sum_{\ell=1}^k w_{s\ell} h_{\ell t}$$

$$= \sum_{st} \left[-X_{st} \frac{\partial}{\partial w_{ij}} \log (WH)_{st} + \frac{\partial}{\partial w_{ij}} (WH)_{st} \right]$$

$$= \sum_t \left[\left(-X_{it} \frac{1}{(WH)_{it}} + 1 \right) \cdot \frac{\partial}{\partial w_{ij}} (WH)_{it} \right]$$

$$= - \sum_t \left[\left(\frac{X_{it}}{(WH)_{it}} - 1 \right) H_{jt} \right]$$

2) Find $\frac{\partial}{\partial H_{ij}} D(X \| WH)$

$$\begin{aligned} \frac{\partial}{\partial H_{ij}} D(X \| WH) &= \sum_{st} \left[-X_{st} \frac{\partial}{\partial H_{ij}} \log (WH)_{st} + \frac{\partial}{\partial H_{ij}} (WH)_{st} \right] \\ &= \sum_s \left[\left(-X_{sj} \frac{1}{(WH)_{sj}} + 1 \right) \frac{\partial}{\partial H_{ij}} (WH)_{sj} \right] \\ &= - \sum_s \left[\left(\frac{X_{sj}}{(WH)_{sj}} - 1 \right) W_{sj} \right] \end{aligned}$$

3) ITERATIVE GRADIENT DESCENT

$$\begin{aligned} W_{ij} &\leftarrow W_{ij} + \Delta_{ij}^W \left[\sum_t \frac{H_{jt} X_{it}}{(WH)_{it}} - \sum_t H_{jt} \right] & i=1, \dots, N \\ & & j=1, \dots, k \\ H_{j\ell} &\leftarrow H_{j\ell} + \Delta_{j\ell}^H \left[\sum_s \frac{W_{sj} X_{s\ell}}{(WH)_{sj}} - \sum_s W_{sj} \right] & \ell=1, \dots, p \end{aligned}$$

* NOTHING PROHIBITS NEGATIVE VALUES OF W_{ij} OR $H_{j\ell}$

4) LEE & SEUNG (2001)

i) SET $\Delta_{ij}^W = \frac{W_{ij}}{\sum_t H_{jt}}$, $\Delta_{j\ell}^H = \frac{H_{j\ell}}{\sum_s W_{sj}}$

$$\begin{aligned} W_{ij} &\leftarrow W_{ij} \left[\frac{\sum_t (H_{jt} X_{it} / (WH)_{it})}{\sum_t H_{jt}} \right] \\ H_{j\ell} &\leftarrow H_{j\ell} \left[\frac{\sum_s (W_{sj} X_{s\ell} / (WH)_{sj})}{\sum_s W_{sj}} \right] \end{aligned}$$

NOTE:

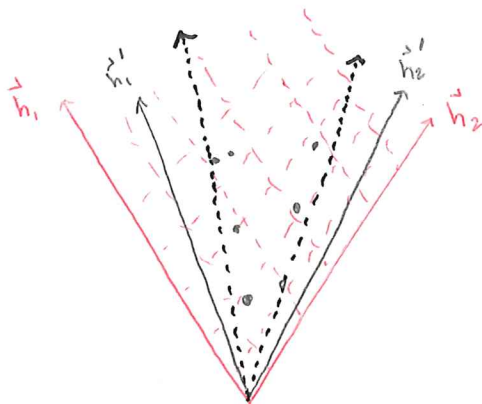
- 1) MULTIPLICATIVE CONSTRAINTS ≥ 0 IMPLIES ITERATES REMAIN ≥ 0
- 2) IF $X = WH$ THEN MULTIPLICATIVE CONSTRAINTS = 1
- 3) IMPLEMENTED IN PACKAGE ALMT

IN THIS EXAMPLE $\vec{x}_1, \dots, \vec{x}_N$ LIVE IN THE SIMPLICIAL CONE OF $\{\vec{h}_1, \vec{h}_2\}$

NOTE: $\vec{x}_1, \dots, \vec{x}_N$ ALWAYS RESIDE IN $\Gamma(\{\vec{x}_1, \dots, \vec{x}_N\})$,
THE CONICAL HULL $\{\vec{x}_1, \dots, \vec{x}_N\}$, TRIVIALY

QUESTION: CAN WE FIND TWO VECTORS $\vec{h}_1' \neq \vec{h}_1, \vec{h}_2' \neq \vec{h}_2$ SUCH THAT

$$\{\vec{x}_1, \dots, \vec{x}_N\} \subseteq \Gamma(\{\vec{h}_1', \vec{h}_2'\})?$$



TRIVIAL CHOICES

- $\vec{h}_1 \rightleftharpoons \vec{h}_2, w_{11} \rightleftharpoons w_{12}$
SWAP FEATURES & WEIGHTS
- $\vec{h}_1 \mapsto c\vec{h}_1, w_{11} \mapsto \frac{1}{c} w_{11}$
 $c > 0$
RESCALE FEATURES/WEIGHTS
- EXCLUDE THESE CASES

YES! INFINITE NUMBER OF CHOICES

- NEW FEATURES $\vec{h}_1', \vec{h}_2' \Rightarrow$ NEW WEIGHTS w_{11}', w_{12}'
 \Rightarrow NEW MATRICES W', H'

$$X = WH = W'H'$$

- EVEN WHEN THERE IS A LOWER-DIMENSIONAL REPRESENTATION OF X , NO GUARANTEE THAT THE SOLUTION $X = WH$ IS UNIQUE!

QUESTION: IGNORING TRIVIAL CHOICES, IF $X = WH$, WHEN ARE W, H UNIQUE?

ANSWER: 1) CAN PLACE DATA ON THE BOUNDARY OF P !

IN THE PICTURE ABOVE, NO GAP BETWEEN \vec{h}_1, \vec{h}_2 AND DATA!

2) OTHER CASES, (SEE DONOHU & STODEN) BUT DIFFICULT TO VERIFY W/ DATA
LAURBERG ET AL.

ESTIMATING A NMF

GIVEN $X \in \mathbb{R}^{n \times p}$ AND $k < \min\{n, p\}$ WE WANT TO

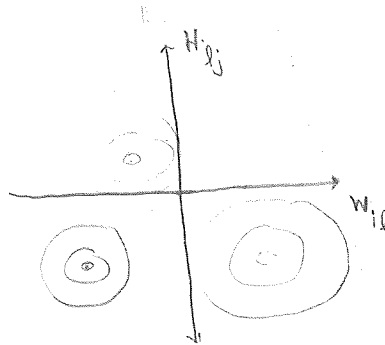
MINIMIZE

$$\|X - WH\|_F \quad \text{OR} \quad D(X \| WH) \quad \text{OR GENERAL ERROR } E(X, WH)$$

FOR $W \in \mathbb{R}^{n \times k}$ $H \in \mathbb{R}^{k \times p}$

1) NO CLOSED FORM EXPRESSION FOR W OR H IN GENERAL

2) NON-CONVEX \rightarrow MANY LOCAL MINIMA



* DIFFICULT TO FIND THE
GLOBAL MINIMIZER!

\Rightarrow USE GRADIENT DESCENT (OR VARIANTS) BUT
TRY MANY RANDOM INITIAL CONDITIONS THEN
PICK BEST

GRADIENT DESCENT

GOAL: MINIMIZE $f(x)$ $\hat{x} \in \mathbb{R}^p$

i) PICK SMALL STEP SIZE Δx

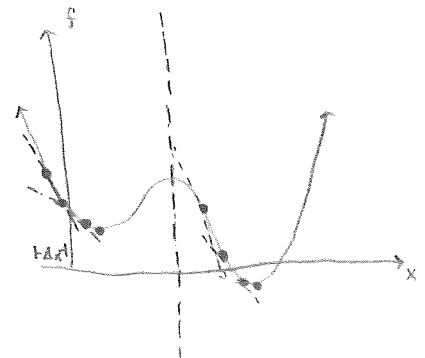
ii) FIND $\nabla f = \begin{bmatrix} \partial_{x_1} f \\ \vdots \\ \partial_{x_p} f \end{bmatrix}$

iii) ITERATE

$$\hat{x}_{i+1} = \hat{x}_i - \Delta x \nabla f(\hat{x}_i)$$

"TAKE A STEP DOWNHILL"

UNTIL SEQUENCE CONVERGES!



Ex: HANDWRITTEN 4s

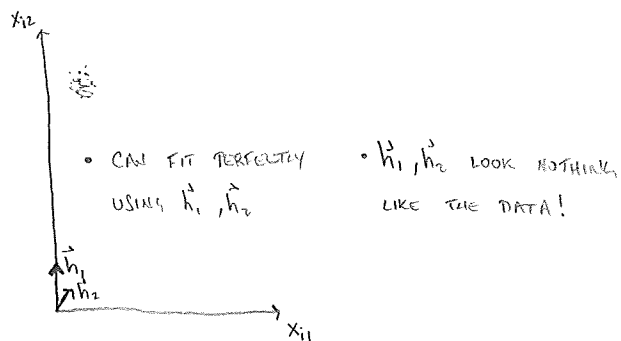
- * SEE MARKDOWN DOCUMENT ON CANVAS > FILES > NUMERICAL EXAMPLES
- * NOTE DIFFERENCES FROM PCA, OUTPUT AND # OF DATA POINTS

NMF IN PRACTICE

- 0) (OPTIONAL) RESCALE DATA (DON'T STANDARDIZE WHICH GIVES < 0 ENTRIES)
- 1) IF PRIOR KNOWLEDGE SUGGESTS A SPECIFIC RANK k , USE IT. OTHERWISE, CHOOSE A SUITABLE RANGE OF RANKS.
- 2) FOR EACH RANK, GENERATE MULTIPLE FITS W/ DIFFERENT INITIALIZATION OF THE ALGORITHM
 - i) CHOOSE BEST AT EACH RANK
 - ii) COMPARE RESIDUALS AT DIFFERENT RANKS AND MAKE CHOICE OF k BASED ON RESULTS
 - * BALANCE OVERFITTING W/ LOW RANK
- 3) CAN ALSO CONSIDER
 - DIFFERENT CHOICES OF ERROR $\| \cdot \|_F$, $D(\| \cdot \|_{WH})$ OR PENALTIES TO IMPOSE SPARSITY IN WH (HOVER 2004)
 - FIXING/INITIALIZING W OR H BASED ON PRIOR KNOWLEDGE

ARCHETYPE ANALYSIS

DEFAULT NMF MAY NOT GIVE FEATURES H WHICH ARE REPRESENTATIVE OF DATA

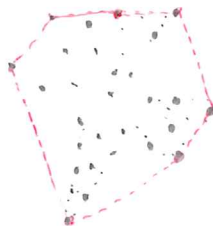


* REQUIRE h_1, \dots, h_k TO BE CONVEX COMBINATIONS OF THE DATA

$$\vec{h}_k^T = \sum_{n=1}^N \pi_{kn} \vec{x}_n^T \quad \pi_{kn} \geq 0 \quad \sum_n \pi_{kn} = 1$$

$$\Rightarrow H = \begin{bmatrix} h_1^T \\ \vdots \\ h_k^T \end{bmatrix} = \begin{bmatrix} \pi_{11} & \dots & \pi_{1n} \\ \vdots & & \vdots \\ \pi_{k1} & \dots & \pi_{kn} \end{bmatrix} X$$

Ex: CONSIDER THE DATA BELOW. IDENTIFY THE ^{SMALLEST} REGION CONTAINING ALL POSSIBLE CONVEX COMBOS OF THE DATA POINTS, THE CONVEX HULL OF THE DATA!



i.e. THROW A LOOP AROUND POINTS AND PULL IT TIGHT!

ARCHETYPAL ANALYSIS

• MINIMIZE $\|X - WBX\|_F$ SUBJECT TO CONSTRAINTS

$$X \in \mathbb{R}^{N \times p}$$

$$i) W_{ij} \geq 0$$

$$W \in \mathbb{R}^{N \times k}$$

$$ii) B_{ij} \geq 0$$

$$B \in \mathbb{R}^{k \times N}$$

$$iii) B\mathbf{1} = \mathbf{1}, \mathbf{1} = \text{VECTOR OF ALL ONES}$$

• THE ROWS OF BX ARE CALLED ARCHETYPES!

NOTE: IN NMF, WE REQUIRE $k \leq \min\{N, p\}$

• IN ARCHETYPAL ANALYSIS, WE CAN LET $k > p$ TO SEARCH FOR k CLUSTERS/ARCHETYPES/PROTOTYPES IN OUR DATA.

* RELATED TO k -MEANS CLUSTERING!

THE NUMERICAL NMF EXAMPLE CONTAINS AN ANALYSIS USING ARCH. ANALYSIS!

DIMENSION REDUCTION: NMF V. ARCH. ANALYSIS

GOALS OF TWO METHODS ARE SIMILAR BUT DIFFERENT FOCUS

1) IN NMF, LOW DIMENSIONAL REPRESENTATION USUALLY
FROM WEIGHTS IN W

$$\text{i.e. } X \in \mathbb{R}^{N \times p}, W \in \mathbb{R}^{N \times k}, H \in \mathbb{R}^{k \times p} \quad \boxed{k < \min\{N, p\}}$$

REDUCE FROM p DIMENSIONS TO k DIMENSIONS $\mathbb{R}^{N \times p} \mapsto \mathbb{R}^{N \times k}$

2) IN ARCH. ANALYSIS, ARCHETYPE ARE OF GREATER

EMPHASIS, CAN HAVE MORE THAN p !

$$\text{i.e. } X \in \mathbb{R}^{N \times p}, W \in \mathbb{R}^{N \times k}, B \in \mathbb{R}^{k \times N} \quad k < N$$

REDUCE FROM p DIMENSIONS TO k LABELS OR ARCHETYPES

• \vec{x}_i IS LIKE ARCHETYPE j !
 $\mathbb{R}^{N \times p} \rightarrow \{1, \dots, k\}$ - CLUSTER POINTS BASED ON MOST
SIMILAR ($w_{ij} \approx 1$) ARCHETYPE!

$$\begin{array}{l} \vec{x}_1, \dots, \vec{x}_N \\ \downarrow \\ \text{ARCHETYPES} \\ \begin{bmatrix} \vec{y}_1^T \\ \vdots \\ \vec{y}_k^T \end{bmatrix} = B X \end{array} \quad \mathbb{R}^{N \times p} \mapsto \mathbb{R}^{k \times p}$$

* IN PSET 4 TAKE NOTE OF LOCATION OF ARCHETYPES

FOR INCREASING k . CONSIDER CONNECTIONS W/

CONVEX HULLS AND POSITIVE SPANS!