

Machine Learning & Pattern Recognition

SONG Xuemeng

songxuемeng@sdu.edu.cn

<http://xuемeng.bitcron.com/>

Unsupervised Feature Extraction

- **Principal Component Analysis (PCA)**
- **Nonnegative Matrix Factorization (NMF)**

Unsupervised Feature Extraction

- **Principal Component Analysis (PCA)**
- **Nonnegative Matrix Factorization (NMF)**

What is feature extraction?

- Feature extraction (dimensionality reduction/feature reduction) refers to the mapping of the original **high-dimensional** data into a **low-dimensional** space.
- Criterion for feature reduction can be different based on different problem setting
 - ✓ Unsupervised setting: minimize the information loss
 - ✓ Supervised setting: maximize the class discrimination

Feature Extraction VS. Feature Selection

- Feature extraction
 - All original features are used.
 - The transformed features are linear combinations of the original features
- Feature selection
 - Only a subset of the original features are used.

Why Feature Extraction?

- Machine learning techniques may not be effective for high-dimensional data
 - **Curse of Dimensionality**
 - Accuracy and efficiency degrade rapidly as the dimension increases
- The **intrinsic** dimension may be small
 - For example, the number of genes responsible for a certain type of disease may be small

Why Feature Extraction?

- **Visualization:** projection of high-dimensional data onto 2D or 3D
- **Data compression:** efficient storage and retrieval
- **Noise removal:** positive effect on query accuracy

Feature Extraction Algorithms

- **Unsupervised**
 - Principal Component Analysis (**PCA**)
 - Nonnegative Matrix Factorization (**NMF**)
 - Independent Component Analysis (ICA) [Reading]
- **Supervised**
 - Linear Discriminant Analysis (**LDA**)
 - General Graph Embedding (GE) [Reading]
 - Canonical Correlation Analysis (CCA) [Reading, encouraged]
- **Semi-supervised**
 - Research topic [Further study, encouraged]

Principal Component Analysis (PCA)

- PCA represents the high-dimensional data in a more tractable, lower-dimensional form, without losing too much information.
 - Reduce the dimensionality of a data set by finding a new set of variables, smaller than the original set of variables.
 - Capture big (principal) variability in the data and ignore small variability.
-
- The new variables, called principal components (PCs), are **uncorrelated**, and are **ordered** by the fraction of the total information each retains.

Algebraic Derivation of PCs

Given a sample set of m observations on a vector of d variables

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\} \in \mathbb{R}^d$$

Define the first PC of the samples by the linear projection $\mathbf{w}_1 \in \mathbb{R}^d$

$$z_{1i} = \mathbf{w}_1^T \mathbf{x}_i = \sum_{k=1}^d w_{1k} x_{ik}$$

where $\mathbf{w}_1 = (w_{11}, w_{12}, \dots, w_{1d})^T$ $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$

$$z_1 = \{z_{11}, z_{12}, \dots, z_{1m}\}$$

\mathbf{w}_1 is chosen such that $\text{var}[z_1]$ is maximum.

Algebraic Derivation of PCs

To find \mathbf{w}_1 , first note that

$$\begin{aligned} \text{var}[z_1] &= E((z_1 - \bar{z}_1)^2) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}_1^T \mathbf{x}_i - \mathbf{w}_1^T \bar{\mathbf{x}})^2 \\ &= \frac{1}{m} \sum_{i=1}^m \mathbf{w}_1^T (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}_1 = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 \end{aligned}$$

where $\mathbf{S} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$ is the covariance matrix.

$\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$ is the mean.

Algebraic Derivation of PCs

To find \mathbf{w}_1 , first note that

$$\begin{aligned} \text{var}[z_1] &= E((z_1 - \bar{z}_1)^2) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}_1^T \mathbf{x}_i - \mathbf{w}_1^T \bar{\mathbf{x}})^2 \\ &= \frac{1}{m} \sum_{i=1}^m \mathbf{w}_1^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{w}_1 = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 \end{aligned}$$

where $\mathbf{S} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$ is the covariance matrix.

$\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$ is the mean.

The covariance matrix \mathbf{S} is **symmetric**.

- The eigenvectors must be **orthogonal** to one another.
- The eigenvalues of \mathbf{S} must all be **≥ 0**

Algebraic Derivation of PCs

To find \mathbf{w}_1 that maximizes $\text{var}[z_1]$ subject to $\mathbf{w}_1^T \mathbf{w}_1 = 1$

Let λ be a Lagrange multiplier

$$L = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 + \lambda(\mathbf{w}_1^T \mathbf{w}_1 - 1)$$

$$\Rightarrow \frac{\partial L}{\partial \mathbf{w}_1} = \mathbf{S} \mathbf{w}_1 - \lambda \mathbf{w}_1 = 0 \quad \Rightarrow \quad (\mathbf{S} - \lambda \mathbf{I}_d) \mathbf{w}_1 = 0$$

Therefore \mathbf{w}_1 is an eigenvector of \mathbf{S} .

The corresponding to the largest eigenvalue $\lambda = \lambda_1$

Algebraic Derivation of PCs

To find \mathbf{w}_1 that maximizes $\text{var}[z_1]$ subject to $\mathbf{w}_1^T \mathbf{w}_1 = 1$

Let λ be a Lagrange multiplier

$$L = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 + \lambda(\mathbf{w}_1^T \mathbf{w}_1 - 1)$$

$$\Rightarrow \frac{\partial L}{\partial \mathbf{w}_1} = \mathbf{S} \mathbf{w}_1 - \lambda \mathbf{w}_1 = 0 \quad \Rightarrow \quad (\mathbf{S} - \lambda \mathbf{I}_d) \mathbf{w}_1 = 0$$

Therefore \mathbf{w}_1 is an eigenvector of \mathbf{S} .

The corresponding to the largest eigenvalue $\lambda = \lambda_1$



WHY?

Algebraic Derivation of PCs

Similarly, \mathbf{w}_2 is also an eigenvector of \mathbf{S} , whose eigenvalue $\lambda = \lambda_2$ is the second largest.

In general $\text{var}[z_k] = \mathbf{w}_k^T \mathbf{S} \mathbf{w}_k = \lambda_k$

The k -th largest eigenvalue of \mathbf{S} is the variance of the k -th PC.

The k -th PC z_k retains the k -th greatest fraction of the variation in the sample.

Algebraic Derivation of PCs

- Main Steps for computing PCs:
 - Form the covariance matrix \mathbf{S} .
 - Compute its eigenvectors: $\{\mathbf{w}_i\}_{i=1}^d$
 - The first p eigenvectors $\{\mathbf{w}_i\}_{i=1}^p$ form the p PCs
 - The transformation G consists of the p PCs

$$\mathbf{G} \leftarrow [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_p] \in \mathbb{R}^{d \times p}$$

$$\mathbf{y} = \mathbf{G}^T \mathbf{x} \in \mathbb{R}^p$$

Algebraic Definition of PCs

- In practice, we compute the PCs via singular value decomposition (**SVD**) on the centered data matrix.

- Form the **centered** data matrix:

$$\mathbf{X} = [(\mathbf{x}_1 - \bar{\mathbf{x}}); \dots; (\mathbf{x}_m - \bar{\mathbf{x}})] \in \mathbb{R}^{d \times m}$$

- Compute its SVD:

$$\mathbf{X} = \mathbf{U}_{d \times d} \mathbf{D}_{d \times m} (\mathbf{V}_{m \times m})^T$$

where \mathbf{U} and \mathbf{V} are **orthogonal** matrices, \mathbf{D} is a diagonal matrix.

Algebraic Definition of PCs

- Note that the scatter/covariance matrix can be written as

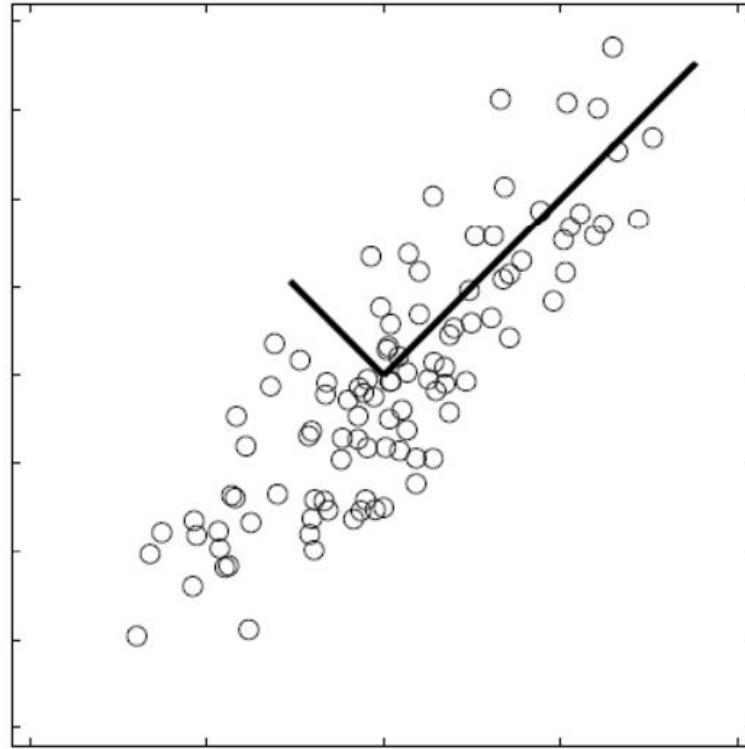
$$\mathbf{S} = \mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{D}^2\mathbf{U}^T \quad \mathbf{X} = \mathbf{U}_{d \times d} \mathbf{D}_{d \times m} (\mathbf{V}_{m \times m})^T$$

- The eigenvectors of \mathbf{S} are the columns of \mathbf{U} and the eigenvalues are the diagonal elements of \mathbf{D}^2 .
- Take **only a few** significant eigenvalue-eigenvector pairs $p \ll d$. The new reconstructed sample from low-dim space is:

$$\hat{\mathbf{x}}_i = \bar{\mathbf{x}} + \mathbf{U}_{d \times p} (\mathbf{U}_{d \times p})^T (\mathbf{x}_i - \bar{\mathbf{x}})$$

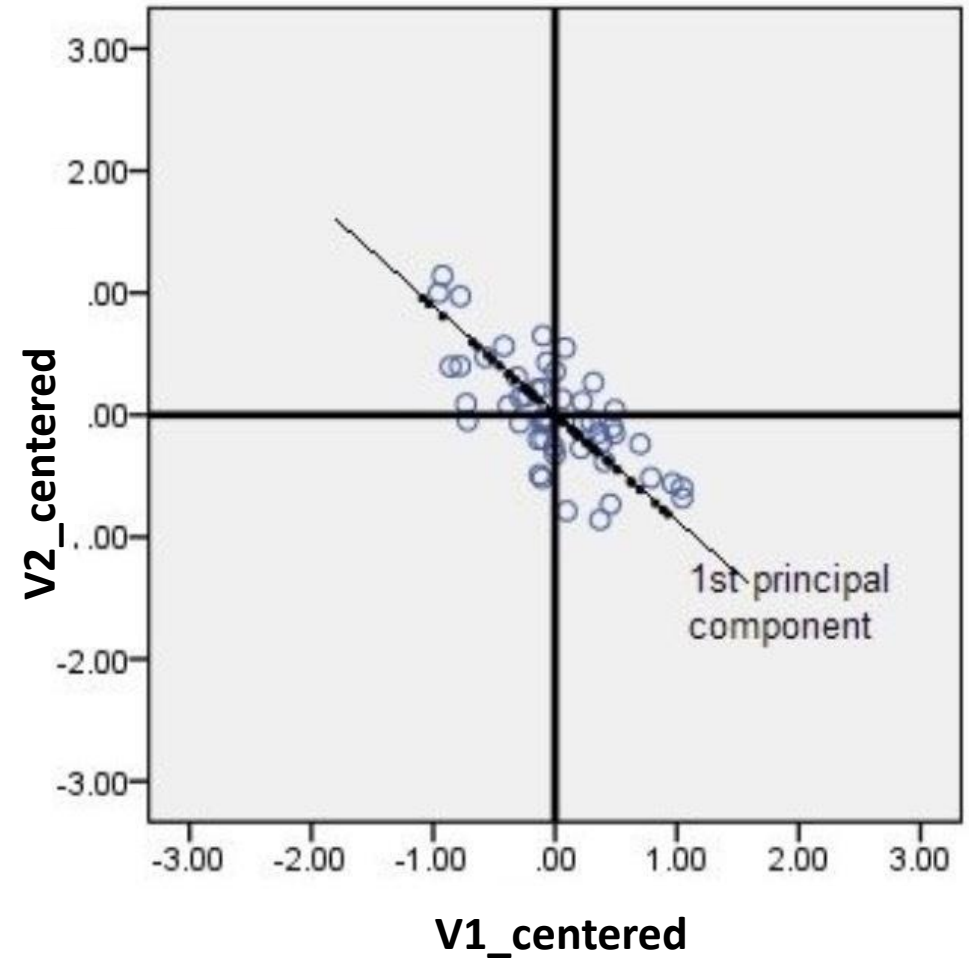
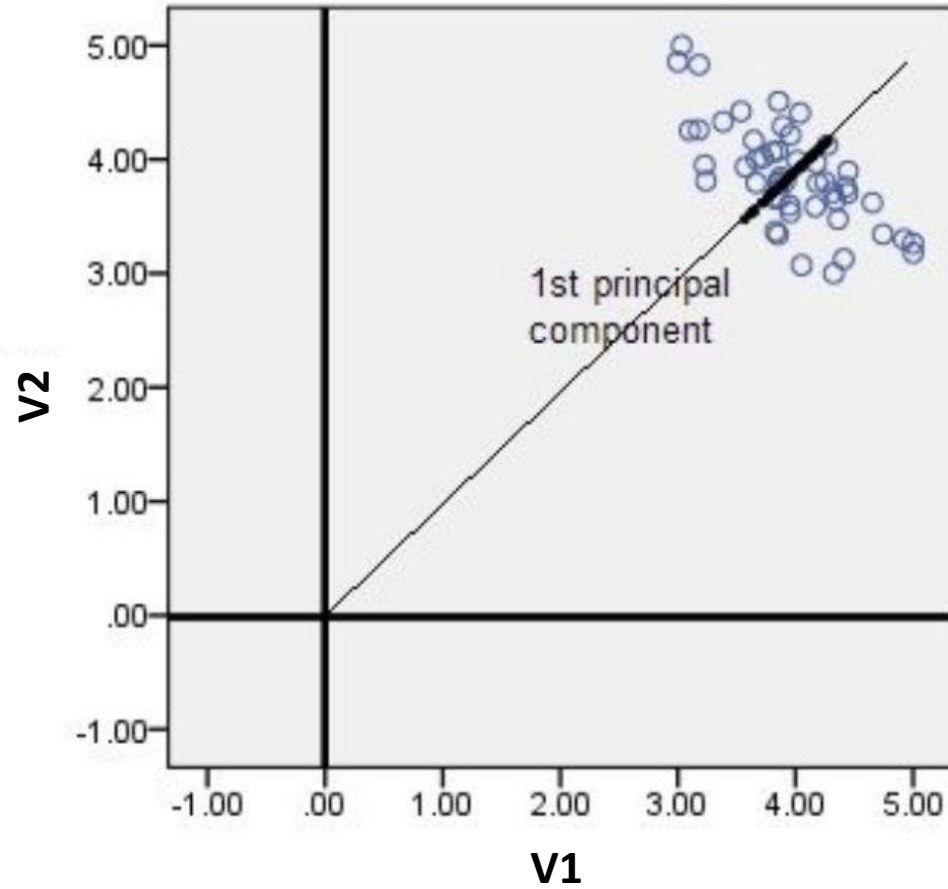
WHY?

Visualize PCs



Data points are represented in a rotated orthogonal coordinate system: the origin is the mean of the data points and the axes are provided by the eigenvectors.

The Necessity of Centralization



How Many PCs to Keep?

How Many PCs to Keep?

To choose p based on percentage of energy to retain, we can use the following criterion (smallest p):

$$\frac{\sum_{i=1}^p \lambda_i}{\sum_{i=1}^d \lambda_i} \geq \textit{Threshold} \quad (e.g., 0.95)$$

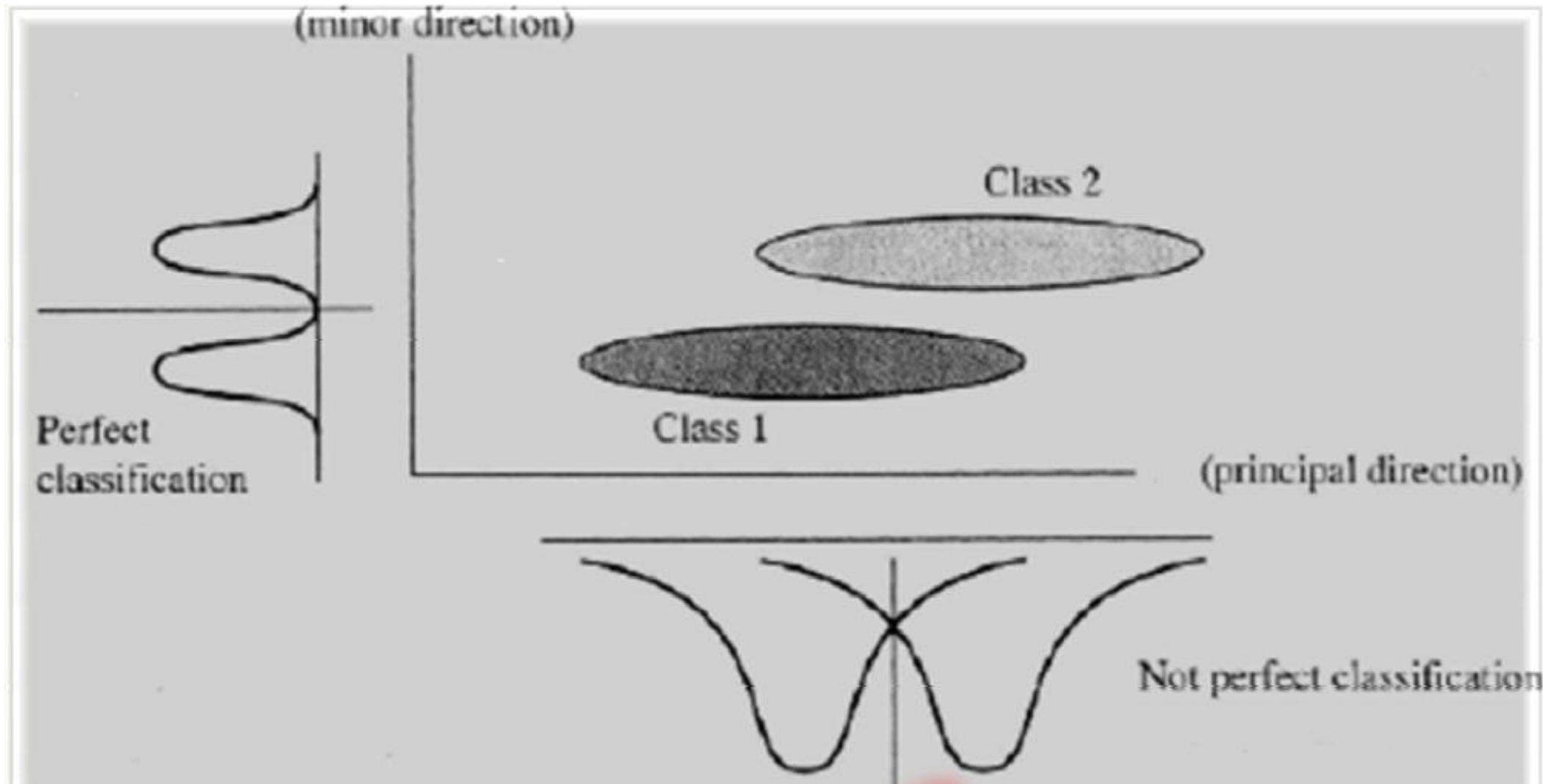
PCA and Classification

- Classification with PCA
 - Project both training and testing data into the PCs space
 - For each testing sample, use NN for classification
 - **Issue:** accuracy is sensitive to the number of PCs

PCA and Classification

- Classification with PCA
 - Project both training and testing data into the PCs space
 - For each testing sample, use NN for classification
 - **Issue:** accuracy is sensitive to the number of PCs
- PCA may not be always an optimal feature extraction technique for classification.
 - Suppose there are C classes in the training data
 - PCA is based on the sample covariance which characterizes the scatter of the **entire** data set, **irrespective of class-membership**.
 - The projection axes chosen by PCA might not provide good discrimination power.

PCA and Classification



Summary of PCA

Algorithm 1 Algorithm for PCA

Input: Samples $\{x_1, x_2, \dots, x_N\}$.

1. Compute the covariance matrix:

$$S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T;$$

2. Perform Eigenvalue Decomposition: $[U] = \text{eig}(S)$;

3. Output PCs matrix $U(:, 1:p)$.

Unsupervised Feature Extraction

- **Principal Component Analysis (PCA)**
- **Nonnegative Matrix Factorization (NMF)**

Non-negative Matrix Factorization

- **PCA review**
 - Find a set of orthogonal principal components (basis)
 - The reconstructed image is a linear combination of PCs plus mean
 - PCA involves adding up some basis vectors and subtracting others.
 - Basis vectors are not physically intuitive for many applications.
 - Sometimes subtracting does not make sense.
 - How to subtract a face? Negative pixel?
- **NMF**: Like PCA, except that the coefficients in the linear combination cannot be negative.

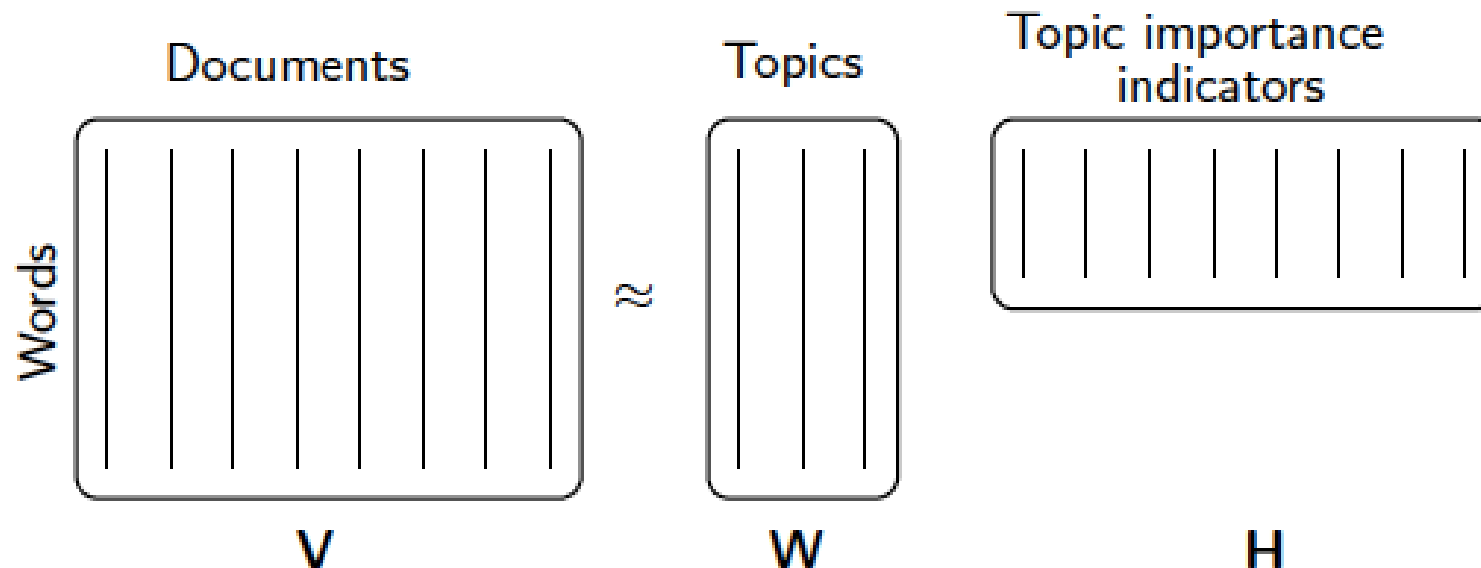
NMF Basis Vectors

- Only allowing **adding** of basis vectors makes intuitive sense
- Forcing the reconstruction coefficients to be nonnegative leads to nice basis vectors
 - To reconstruct vector (image), all you can do is to add in more basis vectors
 - This leads to basis vectors that represent parts

Objective Function

Assume $V \in \mathbb{R}^{m \times n}$ is the sample matrix, the task is to approximate the original data matrix with two nonnegative data matrices $V \approx WH$, $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$:

$$\min_{W, H} \|V - WH\|_F^2, \text{ s.t. } W \geq \mathbf{0} \text{ and } H \geq \mathbf{0}$$



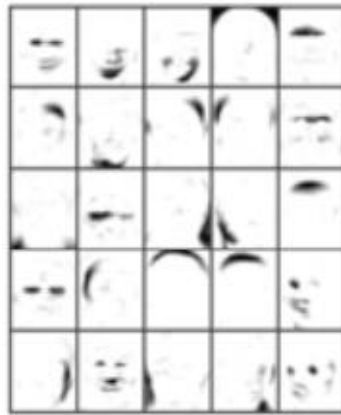
Objective Function

Assume $V \in \mathbb{R}^{m \times n}$ is the sample matrix, the task is to approximate the original data matrix with two nonnegative data matrices $V \approx WH$, $W \in \mathbb{R}^{m \times k}$, $H \in \mathbb{R}^{k \times n}$:

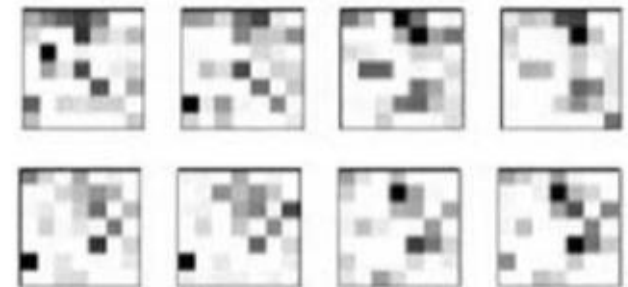
$$\min_{W, H} \|V - WH\|_F^2, \text{ s.t. } W \geq 0 \text{ and } H \geq 0$$



V : an image is a **column** vector.
Vector shown as a matrix here.



W : a **basis** vector is a **column** vector



H : a **coefficient** vector is a **column** vector

Optimization

- Use gradient descent to find a local minimum

$$\begin{aligned} J &= \|V - WH\|_F^2 = \text{tr}((V - WH)(V - WH)^T) \\ &= \text{tr}(VV^T - WHV^T - VH^TW^T + WHH^TW^T) \\ &= \text{tr}(VV^T) - 2\text{tr}(WHV^T) + \text{tr}(WHH^TW^T) \end{aligned}$$

- The gradient descent update rule is (H only, W similar):

$$\frac{\partial J}{\partial H} = ?$$

Matrix calculus

- Numerator layout: lay out according to \mathbf{y} and \mathbf{x}^T . (Jacobian formulation)
- Denominator layout: lay out according to \mathbf{y}^T and \mathbf{x} . (Hessian formulation)

Numerator layout:

分子布局

$$\frac{\partial y}{\partial \mathbf{x}} = \left[\frac{\partial y}{\partial x_1} \frac{\partial y}{\partial x_2} \dots \frac{\partial y}{\partial x_n} \right]$$

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_n}{\partial x} \end{bmatrix}$$

Denominator layout:

分母布局

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial x} = \left[\frac{\partial y_1}{\partial x} \frac{\partial y_2}{\partial x} \dots \frac{\partial y_n}{\partial x} \right]$$

Optimization

- Use gradient descent to find a local minimum

$$\begin{aligned}
 J &= \|V - WH\|_F^2 = \text{tr}((V - WH)(V - WH)^T) \\
 &= \text{tr}(VV^T - WHV^T - VH^TW^T + WHH^TW^T) \\
 &= \text{tr}(VV^T) - 2\text{tr}(WHV^T) + \text{tr}(WHH^TW^T)
 \end{aligned}$$

Condition	Expression	Numerator layout	Denominator layout
A is not a function of X	^[5] $\frac{\partial \text{tr}(\mathbf{X}^\top \mathbf{A} \mathbf{X})}{\partial \mathbf{X}} =$	$\mathbf{X}^\top (\mathbf{A} + \mathbf{A}^\top)$	$(\mathbf{A} + \mathbf{A}^\top) \mathbf{X}$
A is not a function of X	^[5] $\frac{\partial \text{tr}(\mathbf{X}^{-1} \mathbf{A})}{\partial \mathbf{X}} =$	$-\mathbf{X}^{-1} \mathbf{A} \mathbf{X}^{-1}$	$-(\mathbf{X}^{-1})^\top \mathbf{A}^\top (\mathbf{X}^{-1})^\top$
A , B are not functions of X	$\frac{\partial \text{tr}(\mathbf{A} \mathbf{X} \mathbf{B})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{B} \mathbf{A} \mathbf{X})}{\partial \mathbf{X}} =$	$\mathbf{B} \mathbf{A}$	$\mathbf{A}^\top \mathbf{B}^\top$
A , B , C are not functions of X	$\frac{\partial \text{tr}(\mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^\top \mathbf{C})}{\partial \mathbf{X}} =$	$\mathbf{B} \mathbf{X}^\top \mathbf{C} \mathbf{A} + \mathbf{B}^\top \mathbf{X}^\top \mathbf{A}^\top \mathbf{C}^\top$	$\mathbf{A}^\top \mathbf{C}^\top \mathbf{X} \mathbf{B}^\top + \mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B}$

Optimization

- Use gradient descent to find a local minimum

$$\begin{aligned} J &= \|V - WH\|_F^2 = \text{tr}((V - WH)(V - WH)^T) \\ &= \text{tr}(VV^T - WHV^T - VH^TW^T + WHH^TW^T) \\ &= \text{tr}(VV^T) - 2\text{tr}(WHV^T) + \text{tr}(WHH^TW^T) \end{aligned}$$

- The gradient descent update rule is (***H*** only, ***W*** similar):

$$H_{a\mu} \leftarrow H_{a\mu} + \eta_{a\mu} [(W^TV)_{a\mu} - (W^TWH)_{a\mu}]$$

WHY?

Deriving Update Rules

- Gradient Descent Rule:

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

Justify later!

- Set $\eta_{a\mu} = \frac{H_{a\mu}}{(W^T W H)_{a\mu}}$, the update rule becomes

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

What is Significant about This?

- This is a **multiplicative** update instead of an additive update.
 - If the **initial** values of W and H are all **non-negative**, then the W and H can never become negative.
- This lets us produce a **non-negative** factorization

How do we know that this will converge?

WARNING

**Math
Ahead**

REMAIN CALM

How do we know that this will converge?

Definition 1: $G(h, h')$ is an auxiliary function for $F(h)$ if the conditions

$$G(h, h') \geq F(h) \quad G(h, h) = F(h)$$

are satisfied.

How do we know that this will converge?

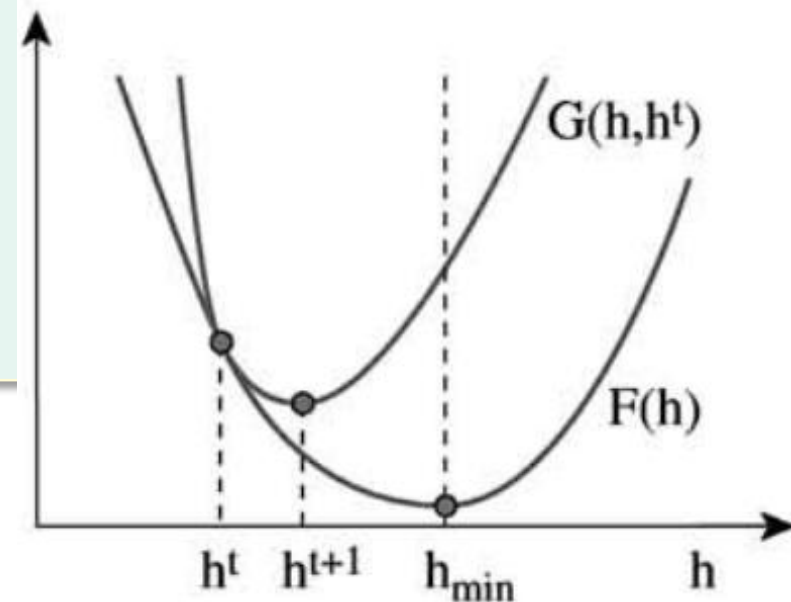
Definition 1: $G(h, h')$ is an auxiliary function for $F(h)$ if the conditions

$$G(h, h') \geq F(h) \quad G(h, h) = F(h)$$

are satisfied.

Lemma 1: If G is an auxiliary function of F , then F is nonincreasing under the update

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



How do we know that this will converge?

Definition 1: $G(h, h')$ is an auxiliary function for $F(h)$ if the conditions

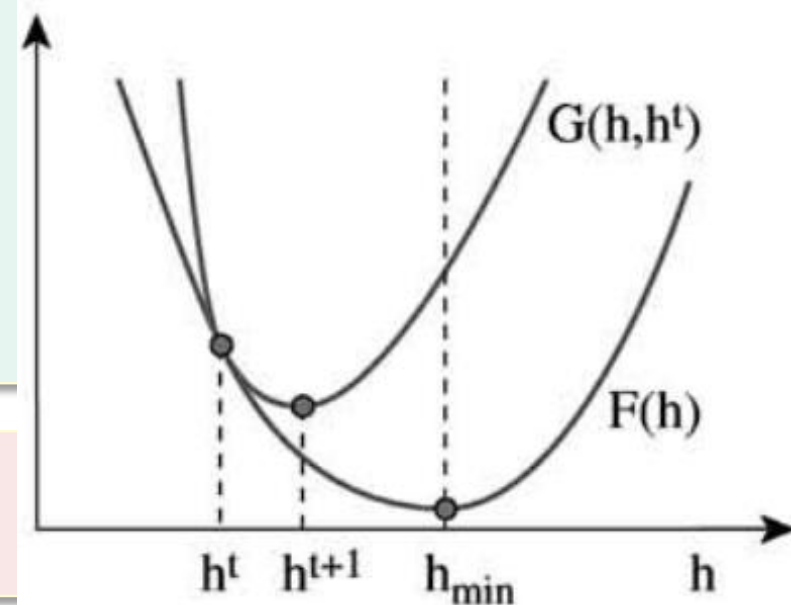
$$G(h, h') \geq F(h) \quad G(h, h) = F(h)$$

are satisfied.

Lemma 1: If G is an auxiliary function of F , then F is nonincreasing under the update

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

Proof: $F(h^{t+1}) \leq G(h^{t+1}, h^t) \leq G(h^t, h^t) = F(h^t)$



How do we know that this will converge?

Lemma 2: If $K(\mathbf{h}^t)$ is the diagonal matrix

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

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

is an auxiliary function for $F(\mathbf{h}) = \frac{1}{2} \sum_i (v_i - \sum_a W_{ia} h_a)^2$

❖ **Kronecker delta:** $\delta_{ab} = I_{ab}$, \mathbf{I} is the identity matrix.

How do we know that this will converge?

Lemma 2: If $K(\mathbf{h}^t)$ is the diagonal matrix

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

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

is an auxiliary function for $F(\mathbf{h}) = \frac{1}{2} \sum_i (v_i - \sum_a W_{ia} h_a)^2$

❖ **Kronecker delta:** $\delta_{ab} = I_{ab}$, I is the identity matrix.

$$\begin{aligned} F(\mathbf{H}) &= \frac{1}{2} \|\mathbf{V} - \mathbf{W}\mathbf{H}\|_F^2 = \frac{1}{2} \sum_{j=1}^n \|\mathbf{V}_{\cdot j} - \mathbf{W}\mathbf{H}_{\cdot j}\|_2^2 = \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^d (V_{ij} - \mathbf{W}_{i\cdot} \mathbf{H}_{\cdot j})^2 \\ &= \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^d (V_{ij} - \sum_{a=1}^c \mathbf{W}_{ia} \mathbf{H}_{aj})^2 \end{aligned}$$

How do we know that this will converge?

Lemma 2: If $\mathbf{K}(\mathbf{h}^t)$ is the diagonal matrix

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


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

is an auxiliary function for $F(\mathbf{h}) = \frac{1}{2} \sum_i (v_i - \sum_a W_{ia} h_a)^2$

❖ **Kronecker delta:** $\delta_{ab} = I_{ab}$, \mathbf{I} is the identity matrix.

$$F(\mathbf{H}) = \frac{1}{2} \|\mathbf{V} - \mathbf{W}\mathbf{H}\|_F^2 = \frac{1}{2} \sum_{j=1}^n \|\mathbf{V}_{\cdot j} - \mathbf{W}\mathbf{H}_{\cdot j}\|_2^2 = \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^d (V_{ij} - \mathbf{W}_{i\cdot} \mathbf{H}_{\cdot j})^2$$

$$= \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^d (V_{ij} - \sum_{a=1}^c \mathbf{W}_{ia} \mathbf{H}_{aj})^2$$


$$F(\mathbf{H}_{\cdot j}) = \frac{1}{2} \sum_{i=1}^d (V_{ij} - \sum_{a=1}^c \mathbf{W}_{ia} \mathbf{H}_{aj})^2$$

How do we know that this will converge?

Lemma 2: If $K(\mathbf{h}^t)$ is the diagonal matrix

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


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

is an auxiliary function for $F(\mathbf{h}) = \frac{1}{2} \sum_i (v_i - \sum_a W_{ia} h_a)^2$

❖ **Kronecker delta:** $\delta_{ab} = I_{ab}$, I is the identity matrix.

$$F(\mathbf{H}) = \frac{1}{2} \|\mathbf{V} - \mathbf{W}\mathbf{H}\|_F^2 = \frac{1}{2} \sum_{j=1}^n \|\mathbf{V}_{\cdot j} - \mathbf{W}\mathbf{H}_{\cdot j}\|_2^2 = \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^d (V_{ij} - \mathbf{W}_{i\cdot} \mathbf{H}_{\cdot j})^2$$

$$= \frac{1}{2} \sum_{j=1}^n \sum_{i=1}^d (V_{ij} - \sum_{a=1}^c \mathbf{W}_{ia} \mathbf{H}_{aj})^2$$


$$F(\mathbf{H}_{\cdot j}) = \frac{1}{2} \sum_{i=1}^d (V_{ij} - \sum_{a=1}^c \mathbf{W}_{ia} \mathbf{H}_{aj})^2 \xrightarrow[\mathbf{V}_{\cdot j} \text{ as } \mathbf{v}]{\mathbf{H}_{\cdot j} \text{ as } \mathbf{h}} F(\mathbf{h}) = \frac{1}{2} \sum_i (v_i - \sum_a W_{ia} h_a)^2$$

Proof of Lemma 2

Since $G(\mathbf{h}, \mathbf{h}) = F(\mathbf{h})$ is obvious, we need to show that $G(\mathbf{h}, \mathbf{h}^t) \geq F(\mathbf{h})$. To do this, we compare the 2nd order Taylor polynomial

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

$$\nabla F(\mathbf{h}) = \frac{\partial F(\mathbf{h})}{\partial \mathbf{h}} = -\mathbf{W}^T \mathbf{v} + \mathbf{W}^T \mathbf{W} \mathbf{h} \qquad \frac{\partial F(\mathbf{h})}{\partial \mathbf{h} \partial \mathbf{h}} = \mathbf{W}^T \mathbf{W}$$

Since we know that

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

We find that prove

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

Proof of Lemma 2

$$0 \leq (\mathbf{h} - \mathbf{h}^t)^T [\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}] (\mathbf{h} - \mathbf{h}^t) \iff \mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W} \text{ is psd.}$$

We consider matrix $\mathbf{M} = \mathbf{H}^t [\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}] \mathbf{H}^t$ $\mathbf{H}^t = \mathbf{H}^{tT} = \begin{bmatrix} h_1^t & & \\ & \ddots & \\ & & h_n^t \end{bmatrix}$

$$M(\mathbf{h}^t)_{ab} = h_a^t (\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W})_{ab} h_b^t$$

Then $\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}$ is psd if and only if \mathbf{M} is.  Why?

If \mathbf{M} is positive semidefinite, then for any vector \mathbf{v} , we have

$$\mathbf{v}^T \mathbf{M} \mathbf{v} = \mathbf{v}^T \mathbf{H}^t [\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}] \mathbf{H}^t \mathbf{v} = (\mathbf{H}^t \mathbf{v})^T [\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}] (\mathbf{H}^t \mathbf{v}) \geq 0$$

if and only if $\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}$ is also positive semidefinite.

Proof of Lemma 2

$$\begin{aligned} \mathbf{v}^T \mathbf{M} \mathbf{v} &= \sum_{ab} v_a M_{ab} v_b = \sum_{ab} v_a h_a^t K_{ab}(h^t) h_b^t v_b - \sum_{ab} v_a h_a^t (W^T W)_{ab} h_b^t v_b \\ &= \sum_a v_a h_a^t \frac{(W^T W h^t)_a}{h_a^t} h_a^t v_a - \sum_{ab} v_a h_a^t (W^T W)_{ab} h_b^t v_b \\ &= \sum_{ab} h_a^t (W^T W)_{ab} h_b^t v_a^2 - \sum_{ab} v_a h_a^t (W^T W)_{ab} h_b^t v_b \\ &\quad \text{symmetric} \\ &= \sum_{ab} (W^T W)_{ab} h_a^t h_b^t \left[\frac{1}{2} v_a^2 + \frac{1}{2} v_b^2 - v_a v_b \right] = \frac{1}{2} \sum_{ab} (W^T W)_{ab} h_a^t h_b^t (v_a - v_b)^2 \\ &\geq 0 \end{aligned}$$

$$\mathbf{M} = \mathbf{H}^t [\mathbf{K}(\mathbf{h}^t) - \mathbf{W}^T \mathbf{W}] \mathbf{H}^t$$

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

$$(W^T W h^t)_a = (W^T W)_a h^t = \sum_b (W^T W)_{ab} h_b^t$$

How do we know that this will converge?

- Given the auxiliary function

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

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

- According to lemma 1, F is nonincreasing under the update

$$\mathbf{h}^{t+1} = \arg \min_{\mathbf{h}} G(\mathbf{h}, \mathbf{h}^t) \quad \Rightarrow \quad \mathbf{h}^{t+1} = \mathbf{h}^t - \underbrace{\mathbf{K}(\mathbf{h}^t)^{-1}}_{\eta} \nabla F(\mathbf{h}^t)$$

- The rule can be explicitly written as

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

Deriving Update Rules

- Gradient Descent Rule:

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

Justify later!

- Set $\eta_{a\mu} = \frac{H_{a\mu}}{(W^T W H)_{a\mu}}$, the update rule becomes

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

Deriving Update Rules

Note: The whole procedure iterates between the optimizations of \mathbf{H} and \mathbf{W} until converged, given that \mathbf{W} and \mathbf{H} are initialized as nonnegative values.

- The update rule:

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

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

Deriving Update Rules

Note: The whole procedure iterates between the optimizations of \mathbf{H} and \mathbf{W} until converged, given that \mathbf{W} and \mathbf{H} are initialized as nonnegative values.

- The update rule:

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

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

Your
homework😊

Summary of NMF

Algorithm 2 Algorithm for NMF

Input: Sample matrix $V = [v_1, v_2, \dots, v_N]$.

Initialize W^0 and H^0 as arbitrary positive matrices.

for $t = 0 : 1 : T_{max}$ **do**

$$H_{a\mu}^{t+1} = H_{a\mu}^t \frac{(W^{tT} V)_{a\mu}}{(W^{tT} W^t H^t)_{a\mu}};$$

$$W_{a\mu}^{t+1} = W_{a\mu}^t \frac{(V H^{t+1T})_{a\mu}}{(W^t H^{t+1} H^{t+1T})_{a\mu}};$$

If $\|W^t - W^{t+1}\| < \epsilon$ and $\|H^t - H^{t+1}\| < \epsilon$

 return;

end for

3. Output matrices W and H .

Discussion

- For new image, how to obtain the reconstruction coefficients?