

Machine Learning

PATRICK HALL

DEPARTMENT OF DECISION SCIENCE

Matrix Factorization

- Dimension Reduction:
 - **Principle Component Analysis**
 - **Singular Value Decomposition**
 - **Non-Negative Matrix Factorization**
 - **Generalized Low Rank Models**
 - **Factorization Machines**

Principal Components

Principal Component Analysis

Principal Components are a sequence of projections of the data - mutually correlated and ordered in variance.

- ▶ Linear manifolds approximating a set of N points $x_i = \mathbb{R}^p$
- ▶ Nonlinear generalization of principal curves and surfaces

Principal Components

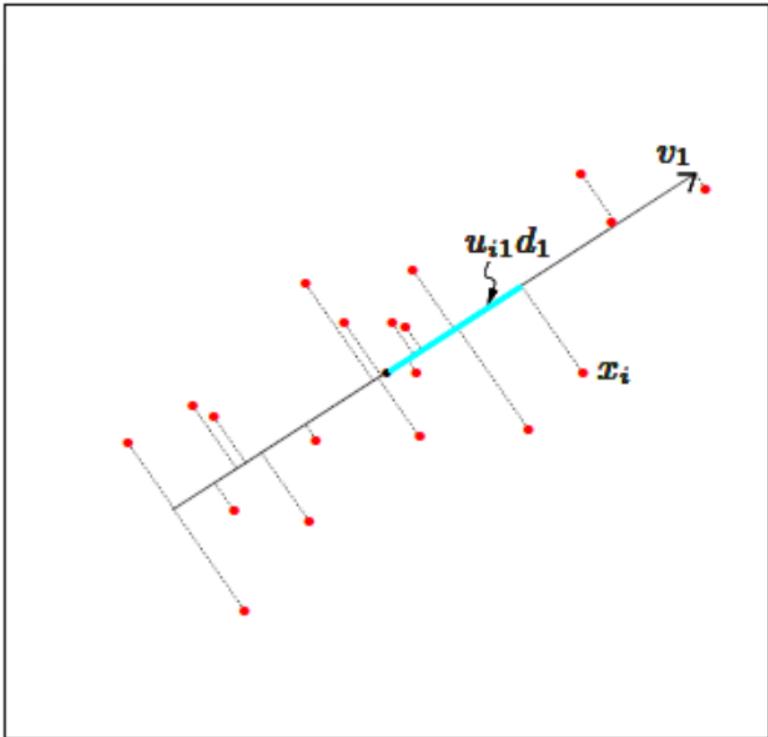


FIGURE 14.20. The first linear principal component of a set of data. The line minimizes the total squared distance from each point to its orthogonal projection onto the line.

Elements of Statistical Learning (pg.534) –

The one-dimensional principal component line in the \mathbb{R}^2 with $q = 1$. For each data point x_i , there is a closest point on the line given by $u_{i1}d_1v_1$, where v_1 is the direction of the line and $\hat{\lambda}_i = u_{i1}d_1$ measures the distance along the line from the origin.

Principal Components

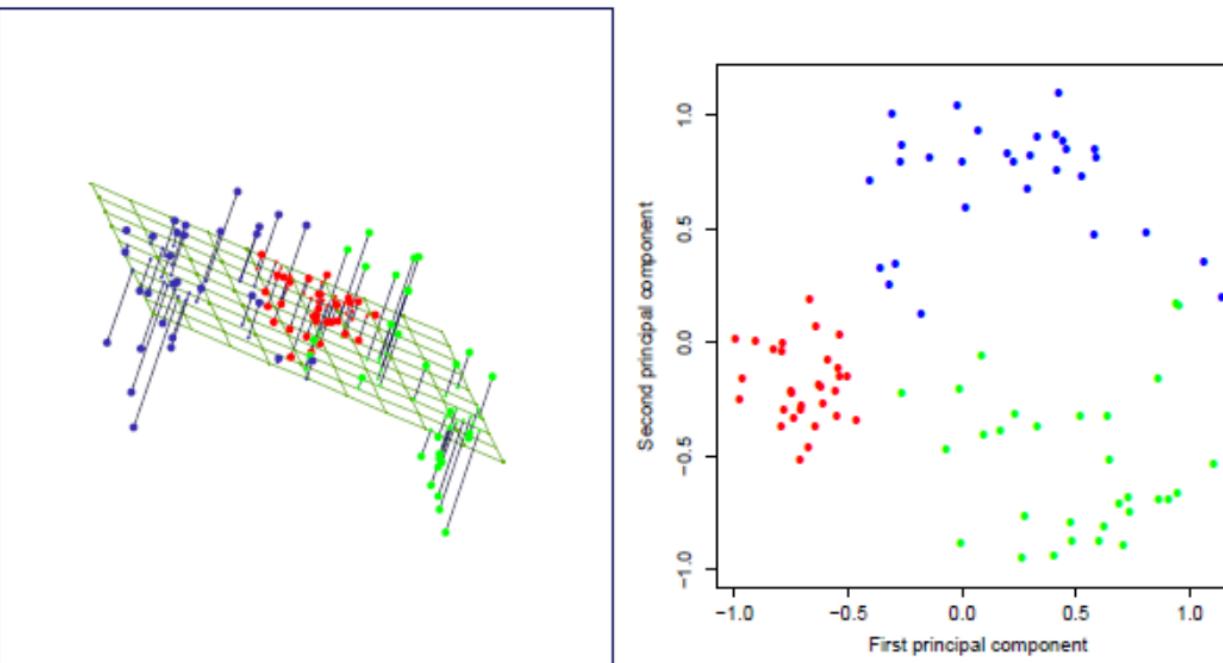


FIGURE 14.21. The best rank-two linear approximation to the half-sphere data. The right panel shows the projected points with coordinates given by $\mathbf{U}_2\mathbf{D}_2$, the first two principal components of the data.

Elements of Statistical Learning (pg.536)

The two-dimensional principal component surface with $q = 2$ to fit the half-sphere data.

Principal Components

Handwritten Digits

Principal components are useful tools for dimension reduction and compression.

Consider sample of 130 handwritten 3's - each a digitized 16x16 grayscale image as shown in Figure 14.22.

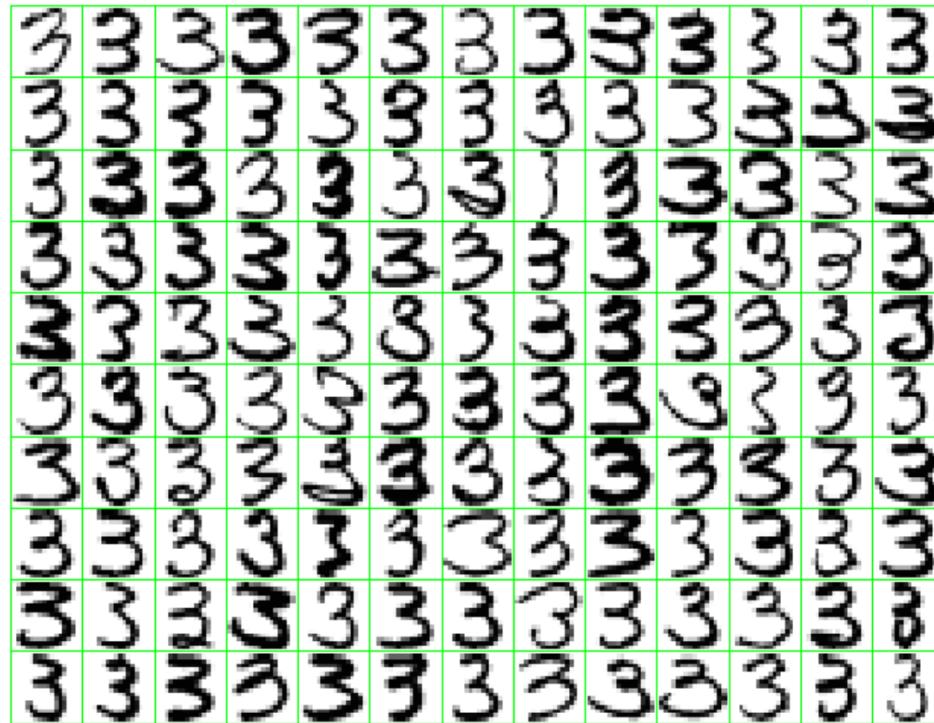


FIGURE 14.22. A sample of 130 handwritten 3's shows a variety of writing styles.

Principal Components

The principal components are computed via SVD and Figure 14.23 shows the first two principal components of these data. In the parametrized model of the form, equation(4), the two-component model is defined as:

$$\begin{aligned}\hat{f}(\lambda) &= \bar{x} + \lambda_1 v_1 + \lambda_2 v_2 \\ &= \boxed{3} + \lambda_1 \cdot \boxed{3} + \lambda_2 \cdot \boxed{3}.\end{aligned}$$

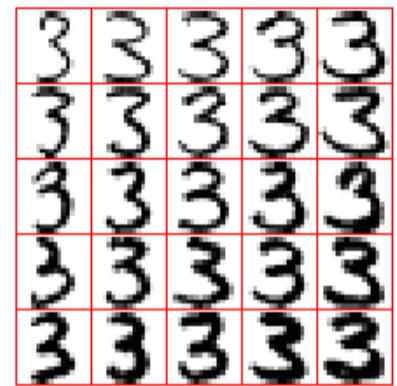
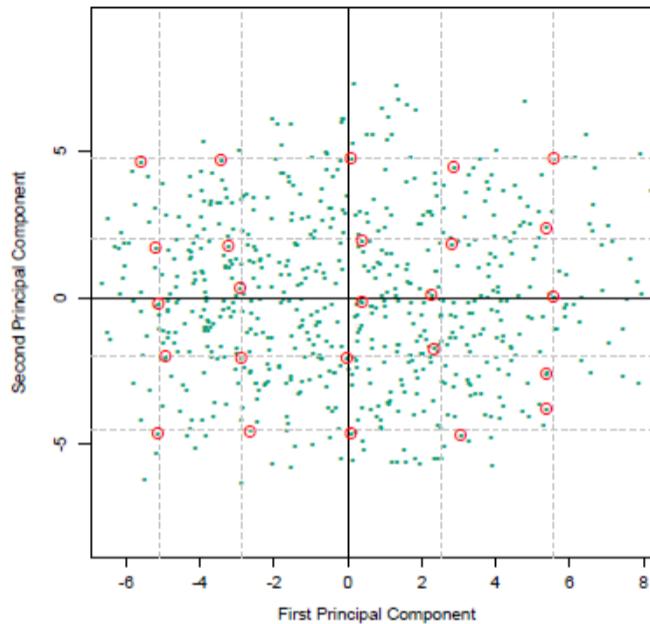


FIGURE 14.23. (Left panel:) the first two principal components of the handwritten threes. The circled points are the closest projected images to the vertices of a grid, defined by the marginal quantiles of the principal components. (Right panel:) The images corresponding to the circled points. These show the nature of the first two principal components.

Principal Components

Compare the singular values to those obtained for equivalent uncorrelated data by randomly scrambling each column of \mathbf{X} . The pixels in a digitized image are inherently correlated, and since these are all the same digit the correlations are even stronger.

This example well illustrates where a relatively small subset of the principal components serve as excellent lower-dimensional features for representing the high-dimensional data.

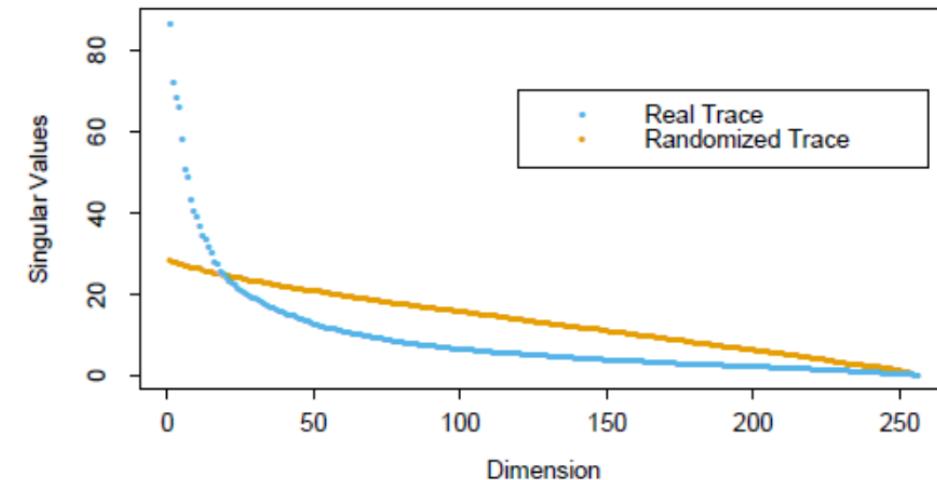


FIGURE 14.24. The 256 singular values for the digitized threes, compared to those for a randomized version of the data (each column of \mathbf{X} was scrambled).

Principal Components

Principal Components

The Principal components of a set of data in $x_i = \mathbb{R}^p$ provide a sequence of best linear approximation to that data with all ranks $q \leq p$. If we denote the set of observations by x_1, x_2, \dots, x_N and consider the *rank – q* linear model:

$$f(x) = \mu + \mathbf{V}_q \lambda \quad (4)$$

where λ is a location vector in $x_i = \mathbb{R}^p$, \mathbf{V}_q is $p \times q$ matrix with q orthogonal unit vectors as columns, and λ is a q vector of parameters.

Principal Components

Fitting a model to the data by least squares amounts to minimizing the *reconstruction error*

$$\min_{\mu, \{\lambda_i\}, \mathbf{V}_q} \sum_{i=1}^N \|x_i - \mu - \mathbf{V}_q \lambda_i\|^2 \quad (5)$$

We can partially optimize for μ and λ_i to obtain:

$$\hat{\mu} = \bar{x}, \quad (6)$$

$$\hat{\lambda}_i = \mathbf{V}_q^T (x_i - \bar{x}) \quad (7)$$

Principal Components

Which leads to find the orthogonal matrix \mathbf{V}_q :

$$\min_{\mathbf{V}_q} = \sum_{i=1}^N \|(\mathbf{x}_i - \bar{\mathbf{x}}) - \mathbf{V}_q \mathbf{V}_q^T (\mathbf{x}_i - \bar{\mathbf{x}})\|^2 \quad (8)$$

Here, $\mathbf{H}_q = \mathbf{V}_q \mathbf{V}_q^T$, (8), is a *pxp projection matrix* that maps each point \mathbf{x}_i onto its rank-q reconstruction $\mathbf{H}_q \mathbf{x}_i$, a orthogonal projection of \mathbf{x}_i onto the subspace spanned by the columns of \mathbf{V}_q .

Principal Components

The solution can be expressed as the stacked (centered) observations into the rows of an $N \times p$ matrix \mathbf{X} and we can construct the *singular value decomposition* of \mathbf{X} :

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \tag{9}$$

where \mathbf{U} is $N \times p$ orthogonal matrix whose columns \mathbf{u}_j are called the *left singular vectors*; and \mathbf{V} is $p \times p$ orthogonal matrix whose columns \mathbf{v}_j are called the *right singular vectors*; and \mathbf{D} is $p \times p$ diagonal matrix with diagonal elements $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ known as the *singular values*. The columns of \mathbf{UD} are called the principal components of \mathbf{X} .

Singular Value Decomposition

Matrix Factorization

Diagonalization: In many cases, the eigenvalues-eigenvetor of matrix \mathbf{A} can be descriptively portrayed in a useful factorization of the form $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ where \mathbf{D} is a diagonal matrix. Moreover, diagonal entries of \mathbf{D} are the eigenvalues of \mathbf{A} that corresponds, respectively, to the eigenvectors of \mathbf{P} . Unfortunately, not all matrices can be factored as $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ with \mathbf{D} diagonal matrix. However, a factorization of a form $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{P}^{-1}$ is possible for any $m \times n$ matrix \mathbf{A} . A special factorization of this type, $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{P}^{-1}$, is called the *singular value decomposition*, and is one of the most useful matrix factorization in applied linear algebra and its related applications.

Singular Value Decomposition

Singular Values of an $m \times n$ Matrix

Let \mathbf{A} be an $m \times n$ matrix. Then $\mathbf{A}^T \mathbf{A}$ is symmetric and can be orthogonally diagonalized. Let $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ be an orthonormal basis for \mathbb{R}^n consisting of eigenvectors of $\mathbf{A}^T \mathbf{A}$, and $\lambda_1, \dots, \lambda_n$ be associated eigenvalues of $\mathbf{A}^T \mathbf{A}$. Then, for $1 \leq i \leq n$,

$$\begin{aligned}\|\mathbf{A}\mathbf{v}_i\|^2 &= (\mathbf{A}\mathbf{v}_i)^T \mathbf{A}\mathbf{v}_i = \mathbf{v}_i^T \mathbf{A}^T \mathbf{A}\mathbf{v}_i \\ &= \mathbf{v}_i^T (\lambda_i \mathbf{v}_i) \\ &= \lambda_i\end{aligned}\tag{1}$$

with, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$. The **singular values** of \mathbf{A} are the square roots of the eigenvalues of $\mathbf{A}^T \mathbf{A}$, $\sigma_1, \dots, \sigma_n$.

Singular Value Decomposition

Now, suppose that $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ is an orthonormal basis of \mathbb{R}^n consisting of eigenvectors of $\mathbf{A}^T \mathbf{A}$, arranged so that the corresponding eigenvalues of $\mathbf{A}^T \mathbf{A}$ satisfy $\lambda_1 \geq \dots \geq \lambda_n$, and suppose that \mathbf{A} has r nonzero singular values. Then, $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_n\}$ is an orthogonal basis for $Col\mathbf{A}$ and $rank\mathbf{A} = r$

Singular Value Decomposition

The Singular Value Decomposition

The "Diagonal" Matrix Σ : The decomposition of A involves an $m \times n$ Σ of the form

$$\Sigma = \begin{bmatrix} \mathbf{D} & 0 \\ 0 & 0 \end{bmatrix} \quad (2)$$

where \mathbf{D} is an $r \times r$ diagonal matrix for some $r \leq \min(m, n)$ with $m - r$ rows and $n - r$ columns of 0 entries.

Singular Value Decomposition

Now, we can define:

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T \quad (3)$$

where \mathbf{A} is an $m \times n$ matrix with $\text{rank } \mathbf{A} = r$. Then there exists an $m \times n$ matrix Σ as defined in (2), where the diagonal entries in \mathbf{D} are the first r singular values of \mathbf{A} , with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$, and there exists an $m \times m$ orthogonal matrix \mathbf{U} and an $n \times n$ orthogonal matrix \mathbf{V} . ANY factorization of the form $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$ with positive diagonal entries in \mathbf{D} is called a **singular value decomposition** of \mathbf{A} .

Singular Value Decomposition

Example:

Find a singular value decomposition of $A = \begin{bmatrix} 1 & -1 \\ -2 & 2 \\ 2 & -2 \end{bmatrix}$

$$\begin{bmatrix} 1 & -1 \\ -2 & 2 \\ 2 & -2 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} & \frac{2}{\sqrt{5}} & \frac{-2}{\sqrt{45}} \\ \frac{-2}{3} & \frac{1}{\sqrt{5}} & \frac{4}{\sqrt{45}} \\ \frac{2}{3} & 0 & \frac{5}{\sqrt{45}} \end{bmatrix} \cdot \begin{bmatrix} 3\sqrt{2} & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

Non-negative Matrix Factorization

Non-negative Matrix factorization - a recent alternative approach to PCA in which the data and components are assumed to be non-negative. It is useful for modeling non-negative dataset such as images.

The $N \times p$ data matrix \mathbf{X} is approximated by

$$\mathbf{X} \approx \mathbf{W}\mathbf{H} \tag{10}$$

where \mathbf{W} is $N \times p$ and \mathbf{H} is $r \times p$, $r \leq \max(N, p)$. We assume that $x_{ij}, w_{ik}, H_{kj} \geq 0$.

Non-negative Matrix Factorization

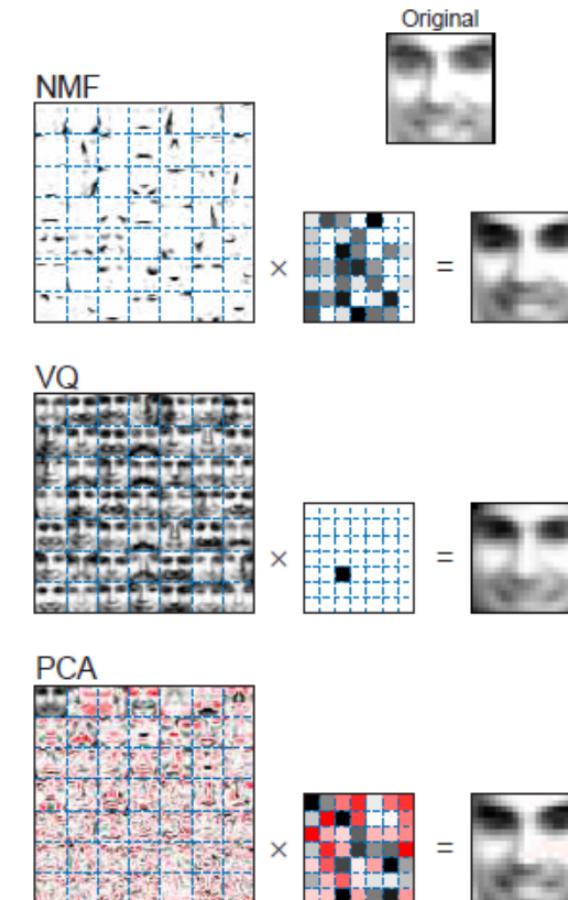
Elements of Statistical Learning (pg.534) – Figure 14.33

Compare the NMF, VQ, and PCA.

All three learning methods were applied to a database of $N=2429$ facial images, each with 19×19 pixels.

Note that, unlike VQ and PCA, NMF learns to represent faces with set of basis images resembling parts of faces

FIGURE 14.33. Non-negative matrix factorization (NMF), vector quantization (VQ, equivalent to k-means clustering) and principal components analysis (PCA) applied to a database of facial images. Details are given in the text. Unlike VQ and PCA, NMF learns to represent faces with a set of basis images resembling parts of faces.



Non-negative Matrix Factorization

The matrices \mathbf{W} and \mathbf{H} are found by maximizing

$$L(\mathbf{L}, \mathbf{H}) = \sum_{i=1}^N \sum_{j=1}^p [x_{ij} \log (\mathbf{WH})_{ij} - (\mathbf{WH})_{ij}] \quad (11)$$

where above (11) is the log-likelihood from a model in which x_{ij} has Poisson distribution with mean $(\mathbf{WH})_{ij}$ - a reasonable for a positive data.

Non-negative Matrix Factorization

The following algorithm converges to a local maximum of $L(\mathbf{W}, \mathbf{H})$:

$$w_{ik} \leftarrow w_{ik} \frac{\sum_{j=1}^p h_{kj} x_{ij} / (\mathbf{WH})_{ij}}{\sum_{j=1}^p h_{kj}} \quad (12)$$

$$h_{kj} \leftarrow h_{kj} \frac{\sum_{i=1}^N w_{ik} x_{ij} / (\mathbf{WH})_{ij}}{\sum_{i=1}^N w_{ik}} \quad (13)$$

Non-negative Matrix Factorization

Archetypal Analysis approximates each data point by a convex combination of a collection of prototypes that lie on the "convex hull of the data cloud." In this sense, the prototypes are "pure" or "archetypal".

Recall, the $N \times p$ matrix \mathbf{X} is modeled as

$$\mathbf{X} \approx \mathbf{W}\mathbf{H} \tag{14}$$

where \mathbf{W} is $N \times r$ and \mathbf{H} is $r \times p$. Here we assume that $w_{ik} \geq 0$ and $\sum_{k=1}^r w_{ik} = 1 \forall i$. Hence, the N data points (row of \mathbf{X}) in $p - \text{dimensional}$ space are represented by convex combination of the r archetypes (rows of \mathbf{H}).

Non-negative Matrix Factorization

We also assume that

$$\mathbf{H} = \mathbf{B}\mathbf{X} \tag{15}$$

where \mathbf{B} is rxN with $b_{ki} \geq 0$ and $\sum_i b_{ki} = 1 \forall k$. Thus the archetypes are convex combinations of the data points. Using both equation(14) and equation(15), we can minimize

$$\begin{aligned} J(\mathbf{W}, \mathbf{B}) &= \|\mathbf{X} - \mathbf{WH}\|^2 \\ &= \|\mathbf{X} - \mathbf{WBX}\|^2 \end{aligned} \tag{16}$$

over the weights \mathbf{W} and \mathbf{B} , and the algorithm converges to a local minimum of the criterion.

Generalized Low Rank Model

GLRM - extension of PCA where:

- ▶ Appropriate loss function instead of least-square approach as in PCA
- ▶ Can add regularization on the low dimensional factors to improve generalization error and avoid overfitting
- ▶ Can impose some structure in the low dimensional factors - such as sparsity

In essence, the term *Generalized Low Rank Model (GLRM)* refer to any low rank approximation of a dataset obtained by minimizing a loss function on the approximation error together with regularization of the low dimensional factors.

Generalized Low Rank Model

Note, in general, low rank approximations problems cannot be solved globally and efficiently. However, can be solved locally by methods that alternate between updating the two factors in the low rank approximation.

That is, each step involves solving a convex or non-convex problem that can be solved exactly. While these alternating methods do not find the globally best low rank approximation, they are very useful and effective for the original data analysis.

Generalized Low Rank Model

Suppose now that our data \mathbf{A} is a database consisting of m examples and n features, with entries \mathbf{A}_{ij} drawn from a feature set \mathcal{F}_j . We observe only entries \mathbf{A}_{ij} for

$(i, j) \in \Omega \subseteq \{1, \dots, m\} \times \{1, \dots, n\}$ from the matrix \mathbf{A} , so the other entries are unknown. Now the user provides a loss function

$L_{ij} : \mathbb{R} \times \mathcal{F}_j \rightarrow \mathbb{R}$. The loss $L_{ij}(u, a)$ describes the approximation error incurred when we represent a feature value $a \in \mathcal{F}_i$ by the number $u \in \mathbb{R}$. The user also provides regularizers

$r_i : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{\infty\}$ for $i = 1, \dots, n$ and $j = 1, \dots, m$. These regularizers prevent overfitting.

Generalized Low Rank Model

Now, we can define a GLRM on the database \mathbf{A} as:

$$\text{minimize} \quad \sum_{(i,j) \in \Omega} L_{ij}(x_i y_j, \mathbf{A}_{ij}) + \sum_{i=1}^m r_i(x_i) + \sum_{j=1}^n \tilde{r}_j(y_j) \quad (23)$$

with variable $X \in \mathbb{R}^{n \times k}$ and $Y \in \mathbb{R}^{k \times m}$, and with loss L_{ij} and regularizers $r_i(x_i) : \mathbb{R}^{1 \times k} \rightarrow \mathbb{R}$ and $\tilde{r}_j(y_j) : \mathbb{R}^{k \times 1} \rightarrow \mathbb{R}$

Generalized Low Rank Model

Few examples of GLRM and $L_{ij}(u, a)$ are illustrated.

Model	$L_{ij}(u, a)$	$r(x)$	$\tilde{r}(y)$
PCA	$(u - a)^2$	0	0
Quadratically regularized PCA	$(u - a)^2$	$\ x\ _2^2$	$\ y\ _2^2$
nonnegative matrix factorization	$(u - a)^2$	$I_+(x)$	$I_+(y)$
Sparse PCA	$(u - a)^2$	$\ x\ _1$	$\ y\ _1$
k -means	$(u - a)^2$	$I_1(x)$	0
Robust PCA	$ u - a $	$\ x\ _2^2$	$\ y\ _2^2$
Boolean PCA (MMMF)	$(1 - au)_+$	$\ x\ _2^2$	$\ y\ _2^2$
Logistic PCA	$\log(1 + \exp(-au))$	$\ x\ _2^2$	$\ y\ _2^2$
Poisson PCA	$\exp(u) - au + a \log a - a$	$\ x\ _2^2$	$\ y\ _2^2$
Ordinal PCA	$\sum_{a' < a} (1 - u + a')_+ + \sum_{a' > a} (1 + u - a')_+$	$\ x\ _2^2$	$\ y\ _2^2$

Table 1: A few examples of GLRMs. Here I_+ is the indicator of the nonnegative orthant, and I_1 is the indicator of the 1-sparse unit vectors.

Factorization Machines

Factorization Machines(FM) can estimate reliable parameters under very high sparsity. The factorization machine models all nested variable interactions, but uses a factorized parametrization, and can be computed in linear time and that it depends only on a linear number of parameters. This allows direct optimization and storage of model parameters without the need of storing any training data for prediction.

The advantages of **FM** are:

- ▶ FM allows parameter estimation under very spare data
- ▶ FM have linear complexity
- ▶ FM are a general predictor that can work with any real valued feature vector

Factorization Machines

FM Model Equation: FM of degree 2 is defined as

$$\hat{y}(x) = w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle v_i, v_j \rangle x_i x_j \quad (20)$$

where the model parameters that have to be estimated are:

$$w_0 \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^n, \text{ and } \mathbf{V} \in \mathbb{R}^{n \times k}$$

Factorization Machines

And $\langle \cdot, \cdot \rangle$ is the dot product of two vectors of size k :

$$\langle v_i, v_j \rangle = \sum_{f=1}^k v_{i,f} \cdot v_{j,f} \quad (21)$$

A row v_i within \mathbf{V} is the i -th variable with k factors and $k \in \mathbb{N}_0^+$ is the hyperparameter that defines the dimensionality of the factorization.

Factorization Machines

FM can be applied in variety of prediction tasks such as:

- ▶ Regression: $\hat{y}(x)$ can be used directly as the predictor and the optimization criterion is the minimum least square error in D .
- ▶ Binary Classification: the sign of $\hat{y}(x)$ is used and the parameters are optimized for hinge loss or digit loss
- ▶ Ranking: the vectors x are ordered by the scores of $\hat{y}(x)W$ and optimization is done over pairs of instance vectors $\langle x^{(a)}, x^{(b)} \rangle \in D$ with pairwise classification loss

Multidimensional Scaling

Multidimensional Scaling has a similar goal of PCA and NMF (lower-dimensional manifold), but approaches the problem in different way.

Lets suppose we have observations $x_1, x_2, \dots, x_N \in \mathbb{R}^p$ and let d_{ij} be the distance between i and j . We often choose the Euclidean distance, $d_{ij} = \|x_i - x_j\|$, but other distances may be used such as *dissimilarity* measure d_{ij} . Multidimensional scaling seeks values $z_1, z_2, \dots, z_N \in \mathbb{R}^k$ to minimize the stress function

$$S_M(z_1, z_2, \dots, z_N) = \sum_{i \neq j} (d_{ij} - \|z_i - z_j\|)^2 \quad (17)$$

This is known as *least squares* or *Kruskal-Shephard* scaling that preserves the pairwise distance as well as possible.

Multidimensional Scaling

A variation on *least squares* scaling is - *Sammon mapping* that seeks to minimizes

$$S_{Sm}(z_1, z_2, \dots, z_N) = \sum_{i \neq i'} \frac{(d_{ii'} - \|z_i - z_{i'}\|)^2}{d_{ii'}} \quad (18)$$

with an emphasis on preserving smaller pairwise distance.

In *classical scaling*, we use the centered inner product

$s_{ii'} = \langle x_i - \bar{x}, x_{i'} - \bar{x} \rangle$ and minimize

$$S_C(z_1, z_2, \dots, z_N) = \sum_{i, i'} (s_{ii'} - \langle z_i - \bar{z}, z_{i'} - \bar{z} \rangle)^2 \quad (19)$$

over the $z_1, z_2, \dots, z_N \in \mathbb{R}^k$