

# Computational Intelligence Lab FS2011 <sup>1</sup>

Stefan Heule

August 22, 2011

<sup>1</sup>License: Creative Commons Attribution-Share Alike 3.0 Unported (<http://creativecommons.org/licenses/by-sa/3.0/>)

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Dimension Reduction</b>	<b>2</b>
2.1	Principle Component Analysis (PCA)	2
2.1.1	Computing PCA	2
2.1.2	Successive Variance Maximization	3
2.1.3	Minimum Error Formulation	4
2.1.4	Matrix Viewpoint	4
2.2	Singular Value Decomposition	4
2.2.1	SVD Component Matrices	5
2.2.2	Singular Values	5
2.2.3	SVD Interpretation	5
2.2.4	Range, Nullspace, Rank	6
2.2.5	SVD to Reveal Structure	6
2.2.6	Closest Rank- $k$ Matrix	6
2.2.7	Computing the SVD	7
2.2.8	SVD of Symmetric Matrices	7
2.2.9	Relationship to PCA	7
<b>3</b>	<b>Clustering</b>	<b>8</b>
3.1	The $K$ -means Algorithm	9
3.2	Model Order Selection	10
3.3	Clustering Stability	10
<b>4</b>	<b>Soft Clustering</b>	<b>11</b>
4.1	Probability Basics	12
4.1.1	Categorical Distribution	12
4.1.2	Gaussian Distribution	12
4.1.3	Multivariate Gaussian Distribution	13
4.1.4	Multidimensional Moment Statistics	13
4.1.5	Conditional Probability	13
4.2	Introduction to Mixture Models	13
4.3	Gaussian Mixture Models	14
4.4	Expectation Maximization for GMM	15
<b>5</b>	<b>Sampling</b>	<b>16</b>
5.1	Monte Carlo Methods	16
5.1.1	Expectation from Samples	17
5.2	Importance Sampling	17

5.3	Rejection Sampling . . . . .	18
<b>6</b>	<b>Role Mining and Multi-Assignment Clustering</b>	<b>19</b>
6.1	Role-Based Access Control . . . . .	19
6.2	Problem Definition . . . . .	19
6.3	Evaluation Criteria for Boolean Matrix Decomposition . . . . .	20
6.3.1	Evaluating a Matrix Reconstruction . . . . .	21
6.3.2	Evaluation Inference Quality . . . . .	21
6.4	SVD for Role Mining . . . . .	22
6.5	Cross-Validation . . . . .	23
6.6	Simple Approach Using a Variant of $K$ -means . . . . .	23
6.7	The ROLEMINER Algorithm . . . . .	24
6.8	DBPSOLVER . . . . .	25
<b>7</b>	<b>Sparse Coding</b>	<b>26</b>
7.1	Comparison to Principle Component Analysis . . . . .	26
7.2	Compressive Sensing . . . . .	27
7.3	Summary of Sparse Coding . . . . .	27
7.4	Overcomplete Dictionaries . . . . .	27
7.4.1	Signal Coding . . . . .	28
7.4.2	Noisy Observations . . . . .	28
7.4.3	Matching Pursuit (MP) . . . . .	29
7.4.4	Support Recovery of MP . . . . .	29
7.4.5	Sparse Coding for Inpainting . . . . .	30
7.5	Dictionary Learning . . . . .	31
7.5.1	Coding Step . . . . .	31
7.5.2	Dictionary Update . . . . .	31
7.5.3	Initialization . . . . .	32

# 1 Introduction

The core problem we consider is *matrix factorization* (MF) of a given data matrix into two or three unknown matrices. Many important data analysis techniques can be written as MF problems by specifying:

- **Type of data:** e.g. boolean ( $x_{d,n} \in \{0, 1\}$ ) or non-negative real.
- **Approximation:** e.g. exact  $\mathbf{X} = \mathbf{A} \cdot \mathbf{B}$  or minimal Frobenius norm  $\min_{\mathbf{A}, \mathbf{B}} \|\mathbf{X} - \mathbf{A} \cdot \mathbf{B}\|_F^2$ .
- **Constraints on unknown matrices:** e.g.  $\mathbf{A}$  has to be a basis or  $\mathbf{B}$  must be sparse.

Often, MF is computationally NP-hard, therefore we need *efficient* and *low-error* approximation algorithms. We study the following techniques:

1. Dimension reduction
2. Single-assignment clustering
3. Multi-assignment clustering
4. Sparse coding

**Definition 1.1** (Orthogonality). *Two vectors in an inner product space  $V$  are orthogonal if their inner product is zero.*

*An orthogonal matrix is a square matrix with real entries whose column (or rows) are orthogonal unit vectors.*

**Theorem 1.1.** *The inverse of an orthogonal matrix is its transpose.*

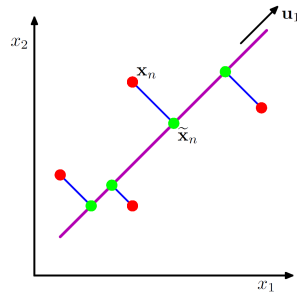


Figure 1: Projection of data points.

## 2 Dimension Reduction

Given  $N$  vectors in  $D$  dimensions, we want to find the  $K$  most important axes and project them. We reduce the dimensions because:

- Some features may be irrelevant
- We want to visualize high dimensional data, e.g., in 2 or 3 dimensions.
- The “intrinsic” dimensionality may be smaller than the number of features recorded in the data.

We want the *most interesting* dimensions:

- data compression
- feature selection
- complexity reduction

### 2.1 Principle Component Analysis (PCA)

Orthogonal linear projection of high dimensional data onto low dimensional subspace, as shown in Figure 1.

We have the following objectives when performing a PCA:

1. Minimize error  $\|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|_2$  of point  $\mathbf{x}_n$  and its approximation  $\tilde{\mathbf{x}}_n$ .
2. Reveal “interesting” information: maximize *variance*.

It is possible to show that both objectives are formally equivalent.

#### 2.1.1 Computing PCA

To perform PCA on a data matrix  $\mathbf{X}$  of size  $D \times N$ :

1. Compute the zero-centered data  $\bar{\mathbf{X}}$  by subtracting the mean of the sample,  $\bar{\mathbf{X}} = \mathbf{X} - \mathbf{M}$ .
2. Calculate the covariance matrix  $\Sigma = \frac{1}{N} \bar{\mathbf{X}} \bar{\mathbf{X}}^T$ .

3. Compute the eigenvectors and eigenvalues of the covariance matrix.
4. Compute the projection of  $\bar{\mathbf{X}}$  on the  $K$  largest principal components  $\mathbf{U}_K = [\mathbf{u}_1, \dots, \mathbf{u}_K]$ .  
By  $\bar{\mathbf{Z}}_K = \mathbf{U}_K^T \bar{\mathbf{X}}$ .

### 2.1.2 Successive Variance Maximization

Consider a set of observations  $\{\mathbf{x}_n\}$ ,  $n = 1, \dots, N$  and  $\mathbf{x}_n \in \mathbb{R}^D$ . We want to project data onto a  $K < D$  dimensional space while maximizing the variance of the projected data.

For  $K = 1$  we define the direction of the projection as  $\mathbf{u}_1$  and set  $\|\mathbf{u}_1\|_2 = 1$ ; only the direction of the projection is important.

For the *original data* with sample mean  $\bar{\mathbf{x}}$ , we have the following covariance:

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

For the *projected data* with new mean  $\mathbf{u}_1^T \bar{\mathbf{x}}$ , the variance is given by

$$\begin{aligned} \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}}\}^2 &= \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^T (\mathbf{x}_n - \bar{\mathbf{x}})\}^2 \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{u}_1^T (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_1 \\ &= \mathbf{u}_1^T \Sigma \mathbf{u}_1 \end{aligned}$$

We can formulate this as a maximization problem, where we maximize the variance of the projected data. We seek

$$\max_{\mathbf{u}_1} \mathbf{u}_1^T \Sigma \mathbf{u}_1$$

such that  $\|\mathbf{u}_1\|_2 = 1$ . We use the Lagrangian of this problem

$$\mathcal{L} := \mathbf{u}_1^T \Sigma \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

and set its partial derivative to zero:

$$\frac{\partial}{\partial \mathbf{u}_1} \mathcal{L} \stackrel{!}{=} 0 \quad \implies \quad \Sigma \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

Therefore,  $\mathbf{u}_1$  is an *eigenvector* of  $\Sigma$  and  $\lambda_1$  is its associated *eigenvalue*. Furthermore,  $\lambda_1$  is also the variance of the projected data, since  $\lambda_1 = \mathbf{u}_1^T \Sigma \mathbf{u}_1$ . This is called the *principle direction*.

With a similar approach we will find that the second principle direction is the second largest eigenvector.

The *eigenvalue decomposition* of the covariance matrix

$$\Sigma = \mathbf{U} \Lambda \mathbf{U}^T$$

contains all relevant information.

For a  $K \leq D$  dimensional projection space we can choose the  $K$  eigenvectors with the largest associated eigenvalues.

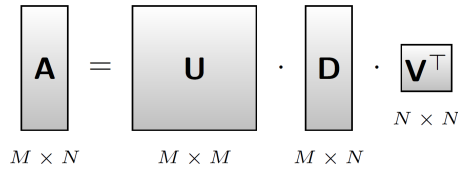


Figure 2: Singular Value Decomposition

### 2.1.3 Minimum Error Formulation

We define an *orthonormal* basis  $\{\mathbf{u}_d\}$ ,  $d = 1, \dots, D$  of  $\mathbb{R}^D$ . The scalar projection of  $\mathbf{x}_n$  onto  $\mathbf{u}_d$  (magnitude) is given by

$$z_{n,d} = \mathbf{x}_n^T \mathbf{u}_d$$

The associated projection onto  $\mathbf{u}_d$  amounts to  $z_{n,d} \mathbf{u}_d$ , and thus we can represent every data point by

$$\mathbf{x}_n = \sum_{d=1}^D z_{n,d} \mathbf{u}_d = \sum_{d=1}^D (\mathbf{x}_n^T \mathbf{u}_d) \mathbf{u}_d$$

### 2.1.4 Matrix Viewpoint

We can also consider PCA from a matrix factorization viewpoint by representing the data as a matrix

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$$

The corresponding zero-centered data is

$$\bar{\mathbf{X}} = \mathbf{X} - \mathbf{M} \quad \text{where} \quad \mathbf{M} = [\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}}]$$

We factor  $\bar{\mathbf{X}}$  by computing the projection of  $\bar{\mathbf{X}}$  on  $\mathbf{U}_K = [\mathbf{u}_1, \dots, \mathbf{u}_K]$  (eigenvectors with  $K$  highest eigenvalues).

$$\bar{\mathbf{Z}}_K = \mathbf{U}_K^T \bar{\mathbf{X}}$$

To approximate  $\bar{\mathbf{X}}$  we return to the original basis:

$$\tilde{\bar{\mathbf{X}}} = \mathbf{U}_K \cdot \bar{\mathbf{Z}}_K$$

Of course, if we use  $K = D$ , we obtain a perfect reconstruction.

## 2.2 Singular Value Decomposition

The *singular value decomposition* (SVD) is a widely used technique to decompose a matrix into several component matrices, exposing many of the useful and interesting properties of the original matrix like rank, nullspace, orthogonal basis of column and row space, etc. Applications which deploy SVD include computing pseudoinverses, least square fitting of data, matrix approximation, and determining the rank, range and null space of a matrix.

**Theorem 2.1** (SVD Theorem). *Let  $\mathbf{A}$  be any real  $M$  by  $N$  matrix,  $\mathbf{A} \in \mathbb{R}^{M \times N}$ , then  $\mathbf{A}$  can be decomposed as  $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T$ .*

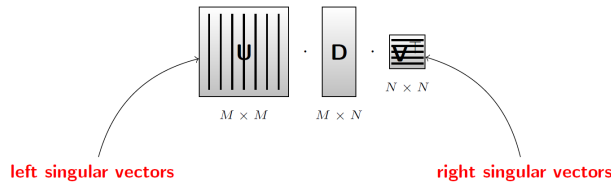


Figure 3: Left and right singular vectors

- $\mathbf{U}$  is an  $M \times M$  orthogonal matrix such that  $\mathbf{U}^T \mathbf{U} = I_{(M)}$ .
- $\mathbf{D}$  is an  $M \times N$  diagonal matrix.
- $\mathbf{V}^T$  is an  $N \times N$  orthogonal matrix such that  $\mathbf{V}^T \mathbf{V} = I_{(N)}$ .

### 2.2.1 SVD Component Matrices

The elements of  $\mathbf{D}$  are only non-zero on the diagonal, and are called the *singular values*. An important relation to the rank  $r$  of a matrix is

$$\begin{cases} d_{i,i} > 0 & 1 \leq i \leq r \\ d_{i,i} = 0 & r + 1 \leq i \leq n \\ d_{i,j} = 0 & i \neq j \end{cases}$$

By convention, the ordering of the singular vectors is determined by high-to-low sorting of the singular values, with the highest singular value in the upper left corner of  $\mathbf{D}$ .

$$d_1 \geq d_2 \geq \dots \geq d_r > d_{r+1} = \dots = d_n = 0$$

The first  $r$  columns of  $\mathbf{U}$  are called the *left singular vectors*, they form an orthogonal basis for the space spanned by the columns of the original matrix  $\mathbf{A}$ . Similarly the first  $r$  rows of  $\mathbf{V}^T$  are the right singular vectors, they form an orthonormal basis for the row space of  $\mathbf{A}$ . This is illustrated in Figure 3.

We can also interpret the columns of  $\mathbf{U}$  as the eigenvectors of  $\mathbf{A}\mathbf{A}^T$ . This can be verified by the SVD:

$$\mathbf{A}\mathbf{A}^T = \mathbf{U}\mathbf{D}\mathbf{V}^T\mathbf{V}\mathbf{D}\mathbf{U}^T = \mathbf{U}\mathbf{D}^2\mathbf{U}^T$$

Similarly for the rows of  $\mathbf{V}^T$  (or the columns of  $\mathbf{V}$ ), which are the eigenvectors of  $\mathbf{A}^T\mathbf{A}$ , since

$$\mathbf{A}^T\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{U}^T\mathbf{U}\mathbf{D}\mathbf{V}^T = \mathbf{V}\mathbf{D}^2\mathbf{V}^T$$

### 2.2.2 Singular Values

The singular values of a matrix  $\mathbf{A}$  represent the *lengths* of the semi-axes of the hyper-ellipsoid  $E$  defined as  $E = \{\mathbf{A}\mathbf{x} \mid \|\mathbf{x}\|_2 = 1\}$ . Hence, they are non-negative.

### 2.2.3 SVD Interpretation

Suppose we are given a matrix  $\mathbf{A}$  that records the ratings of movies by users (rows correspond to users, columns to movies). Then,



- $U$ : user-to-concepts affinity matrix.
- $V$ : movie-to-concepts affinity matrix.
- $D$ : the diagonal elements of  $D$  represent the “expressiveness” of each concept in the data.

#### 2.2.4 Range, Nullspace, Rank

SVD provides an explicit representation of the range and nullspace of a matrix.

- The right singular vectors corresponding to vanishing singular values of  $A$  span the nullspace of  $A$ .
- The left singular vectors correspond to the non-zero singular values of  $A$  and span the range of  $A$ .

#### 2.2.5 SVD to Reveal Structure

SVD can be viewed from three mutually compatible points of view:

1. SVD is a method for transforming correlated variables into a set of uncorrelated ones that better expose the various relationships among the original data items.
2. SVD is a method for identifying and ordering the dimensions along which data points exhibit the most variation.
3. SVD as data reduction method, once we have identified the directions of most variation we can obtain the best approximation of the original data points using fewer dimensions.

#### 2.2.6 Closest Rank- $k$ Matrix

We can use this decomposition to find the  $k$ -closest approximation. Let  $A$  be a matrix of rank  $R$ . If we wish to approximate  $A$  using a matrix of lower rank  $K$ , then

$$A_k = \sum_{k=1}^K d_k \mathbf{u}_k \mathbf{v}_k^T$$

is the closest matrix in the  $L_2$  norm. More formally, we have

$$\min_{B, \text{rank}(B)=k} \|A - B\|_2 = \|A - A_k\|_2$$

Moreover, the magnitudes of the nonzero singular values provide a measure of how close  $A$  is to a matrix of lower rank in the Euclidean norm:

$$\|A - A_k\|_2 = d_{k+1}$$

The *spectral norm* of a matrix  $A$  is the largest singular value of  $A$ :

$$\|A\|_2 = \sigma_{\max}(A)$$

### 2.2.7 Computing the SVD

Given a matrix  $\mathbf{A}$

1. Compute its transpose  $\mathbf{A}^T$  and then  $\mathbf{A}^T \mathbf{A}$ .
2. Determine the eigenvalues of  $\mathbf{A}^T \mathbf{A}$  and sort them in descending order. Square the roots to obtain the singular values of  $\mathbf{A}$ .
3. Construct the diagonal matrix  $\mathbf{D}$  by placing the singular values in descending order along its diagonal.
4. Use the eigenvalues from the second step to compute the eigenvectors of  $\mathbf{A}^T \mathbf{A}$ . Place them (in the right order) along the columns of  $\mathbf{V}$  and compute  $\mathbf{V}^T$ .
5. Compute the inverse of the diagonal matrix,  $\mathbf{D}^{-1}$ , using  $(\mathbf{D}^{-1})_{ii} = 1/D_{ii}$ . Compute the matrix  $\mathbf{U} = \mathbf{A} \mathbf{V} \mathbf{D}^{-1}$ .

### 2.2.8 SVD of Symmetric Matrices

**Theorem 2.2.** *If  $\mathbf{S}$  is a real and symmetric matrix (i.e.,  $\mathbf{S} = \mathbf{S}^T$ ), then  $\mathbf{S} = \mathbf{U} \mathbf{D} \mathbf{U}^T$  where the columns of  $\mathbf{U}$  are the eigenvectors of  $\mathbf{S}$  and  $\mathbf{D}$  is a diagonal matrix with the eigenvalues of  $\mathbf{S}$ .*

### 2.2.9 Relationship to PCA

Given any  $m \times n$  matrix  $\mathbf{A}$ ,

- The matrix  $\mathbf{A} \mathbf{A}^T$  is a symmetric matrix describing the similarity between the rows of the matrix  $\mathbf{A}$ . The eigenvectors  $[\mathbf{u}_1, \dots, \mathbf{u}_n]$  of  $\mathbf{A} \mathbf{A}^T$  are the principle components in a PCA of  $\mathbf{A}$  (equal up to a constant, assuming  $\mathbf{A}$  is centered) and are the left singular vectors in the SVD decomposition of  $\mathbf{A}$ .
- The squared singular values of  $\mathbf{A}$  are the eigenvalues of  $\mathbf{A} \mathbf{A}^T$  and  $\mathbf{A}^T \mathbf{A}$ .

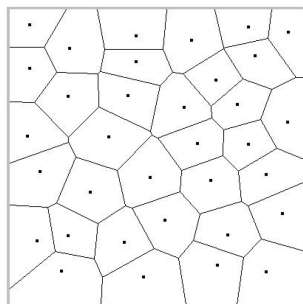


Figure 4: Quantization of the space  $\mathbb{R}^D$

### 3 Clustering

Given a set of data points in a  $d$ -dimensional Euclidean space, we want to find a “meaningful” partition of the data. The partition should group together similar data points (put them in one group/cluster), while the different clusters should be as dissimilar as possible. Partitioning the data in this way can uncover similarities between data points and give rise to data compression schemes.

We can quantize the space  $\mathbb{R}^D$  (as shown in Figure 4). A data point  $\mathbf{x} \in \mathbb{R}^D$  is represented by one of the  $K$  centroids  $\mathbf{u}_k \in \mathbb{R}^D$  for  $k = 1, \dots, K$ . We assign  $\mathbf{x}$  to the centroid  $\mathbf{u}_{k^*}$  that has minimal distance.

**Definition 3.1** (The Clustering Problem). *Consider  $N$  data points in a  $D$ -dimensional space. Each data vector is denoted by  $\mathbf{x}_n$ ,  $n = 1, \dots, N$ . The goal is to partition the data into  $K$  clusters. More formally, we seek vectors  $\mathbf{u}_1, \dots, \mathbf{u}_K$  that represent the centroid of each cluster. A data point  $\mathbf{x}_n$  belongs to the cluster  $k$  if the Euclidean distance between  $\mathbf{x}_n$  and  $\mathbf{u}_k$  is smaller than the distance to any other centroid.*

We want to minimize the following cost function

$$J(\mathbf{U}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{U}\mathbf{Z}\|_2^2 = \sum_{n=1}^N \sum_{k=1}^K z_{k,n} \|\mathbf{x}_n - \mathbf{u}_k\|_2^2$$

The vectors  $\mathbf{z}_n$  are the assignments of data points to clusters, and we can have various constraints on  $\mathbf{Z}$ . For now, we consider *hard assignments*, i.e.,  $\mathbf{Z} \in \{0, 1\}^{K \times N}$  with  $\sum_k z_{k,n} = 1$  for all  $n$ . That is, one entry per column is set to 1.

We have two problems:

- Determine the set of cluster centroids  $\mathbf{u}_k$ .
- Compute the assignment matrix  $\mathbf{Z}$ .

Finding a solution to one of them (given the other) is analytically solvable.

We can write assignments (for hard assignments) down in two alternative notations. Either, we use a vector  $\hat{\mathbf{z}} \in \{1, \dots, K\}^N$  where every entry indicates (for each data point) to which cluster

index it is assigned to, e.g.,

$$\hat{z} = \begin{pmatrix} 2 \\ 3 \\ 4 \\ 2 \\ 2 \\ 1 \end{pmatrix}$$

Alternatively, we can use a matrix  $\mathbf{Z} \in \{0, 1\}^{K \times N}$ , e.g.,

$$\mathbf{Z} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \end{pmatrix}$$

### 3.1 The $K$ -means Algorithm

We *alternate* between two steps:

- Assigning data points to clusters. Let  $k^*(\mathbf{x}_n)$  denote the cluster index with the minimal distance between a cluster centroid and the data point:

$$k^*(\mathbf{x}_n) = \arg \min_k \left\{ \|\mathbf{x}_n - \mathbf{u}_1\|_2^2, \dots, \|\mathbf{x}_n - \mathbf{u}_K\|_2^2 \right\}$$

- updating the cluster centroids based on the data points assigned to it. We compute the mean/centroid of a cluster as

$$\mathbf{u}_k = \frac{\sum_{n=1}^N z_{k,n} \mathbf{x}_n}{\sum_{n=1}^N z_{k,n}} \quad \text{for all } k$$

Then, the  $K$ -means algorithm can be described as

1. Initiate with random choice for  $\mathbf{u}_1^{(0)}, \dots, \mathbf{u}_K^{(0)}$  (or use data points, or some other initialization) and set  $t = 1$ .
2. **Cluster assignment.** Solve for all  $n$

$$k^*(\mathbf{x}_n) = \arg \min_k \left\{ \|\mathbf{x}_n - \mathbf{u}_1^{(t-1)}\|_2^2, \dots, \|\mathbf{x}_n - \mathbf{u}_K^{(t-1)}\|_2^2 \right\}$$

3. **Centroid update.** The centroids are given by

$$\mathbf{u}_k^{(t)} = \frac{\sum_{n=1}^N z_{k,n}^{(t)} \mathbf{x}_n}{\sum_{n=1}^N z_{k,n}^{(t)}} \quad \text{for all } k$$

4. Increment  $t$  and repeat from step 2 until  $\|\mathbf{u}_k^{(t)} - \mathbf{u}_k^{(t-1)}\|_2^2 < \varepsilon$  for all  $k$  or  $t = t_{\text{finish}}$ .

The computational cost of each iteration is  $\mathcal{O}(KN)$  and convergence is guaranteed.  $K$ -means optimizes a non-convex objective and hence we can only guarantee to find a local minimum. The algorithms can be used to compress data (lossy, if  $K < N$  by only storing the centroids and the assignments of data points to clusters).

Problems of the approach include the following

- The non-convex objective guarantees only to find local minima, and the algorithm is sensitive to initializations. To some extent, this can be overcome by restarts.
- The optimal number of clusters  $K$  is unknown. This problem is referred to as *model order selection* or model validation.

### 3.2 Model Order Selection

We have a trade-off between two conflicting goals:

- **Data fit.** We want to predict the data well, e.g., for clustering

$$J(\mathbf{U}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{UZ}\|_2^2 = \sum_{n=1}^N \sum_{k=1}^K z_{k,n} \|\mathbf{x}_n - \mathbf{u}_k\|_2^2$$

should ideally be zero.

- **Complexity.** We want a model that is not very complex, which is often measured by the number of free parameters.

We can use the *negative log-likelihood* of data for  $K$ -means:

$$-\ell(\mathbf{U}, \mathbf{Z} | \mathbf{X}) = J(\mathbf{U}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{UZ}\|_2^2 = \sum_{n=1}^N \sum_{k=1}^K z_{k,n} \|\mathbf{x}_n - \mathbf{u}_k\|_2^2$$

We can use (at least) two heuristics for choosing the parameter  $K$ . We achieve a balance between data fit (measured by the cost function  $-\ell(\mathbf{U}, \mathbf{Z} | \mathbf{X})$ ) and complexity, which we measure by the number of free parameters  $\kappa(\mathbf{U}, \mathbf{Z}) = K \cdot D$ .

- **Akaike Information Criterion (AIC)**

$$\text{AIC}(\mathbf{U}, \mathbf{Z} | \mathbf{X}) = -\ell(\mathbf{U}, \mathbf{Z} | \mathbf{X}) + \kappa(\mathbf{U}, \mathbf{X})$$

- **Bayesian Information Criterion (BIC)**

$$\text{BIC}(\mathbf{U}, \mathbf{Z} | \mathbf{X}) = -\ell(\mathbf{U}, \mathbf{Z} | \mathbf{X}) + \frac{1}{2} \kappa(\mathbf{U}, \mathbf{X}) \cdot \ln N$$

We can then repeat the analysis for different values of  $K$  and choose the smallest AIC/BIC value.

### 3.3 Clustering Stability

The idea behind the notion of *stability* is the following concept. If we repeatedly sample data points from the same underlying model or of the same data generating process and apply a

clustering algorithm then a “good” algorithm should return clusterings that do not vary much from one sample to another. Using the  $K$ -means algorithm, this selection strategy reduces to choosing the optimal number of clusters for which stability is given.

We can use the following high-level stability test for a given set of data points and a given number of clusters:

1. Generate perturbed versions of the set (e.g., by adding noise or drawing sub-samples).
2. Apply the clustering algorithm on all versions.
3. Compute the pair-wise distances between all clusterings.
4. Compute the *instability* as the means distance between all clusterings.

This procedure can be repeated for different numbers of clusters and then one chooses the number that minimizes the instability.

To measure the distance between to clusterings that are defined on the same data points, we compute the distance between clusterings as follows:

1. Compute the distance between two by counting points on which the two clusterings disagree.
2. Repeat over all permutations of the cluster labels (as the same cluster might have different labels)
3. Choose the permutation with minimal distance.

In other words

$$d = \min_{\pi} \|\mathbf{Z} - \pi(\mathbf{Z}')\|_0$$

## 4 Soft Clustering

The term soft clustering is ambiguous, as it can refer to the algorithm or the model:

1. Algorithmic: soft  $K$ -means. In the  $K$ -means algorithms, a data point is assigned to one and only one cluster. A better estimation of the means would be to consider assigning to each cluster a probability that the data point belongs to it. This is known as the Gaussian mixture model (GMM). Note that the model still assumes that a data point belongs to only one cluster.
2. Model relaxation: One can also relax the “hard” constraint given by

$$z_{k,n} \in \{0, 1\} \quad \text{and} \quad \sum_{k=1}^K z_{k,n} = 1$$

and replace it with the following “soft” constraint:

$$z_{k,n} \in [0, 1] \quad \text{and} \quad \sum_{k=1}^K z_{k,n} = 1$$

This is known as *semi-NMF*.

Two modifications of our objective are possible:

1. **Gaussian mixture models.** Different clustering objective which arises from a generative probabilistic model. GMMs can not be written in matrix factorization form. The “softness” refers to the probability of hypotheses and, therefore, its nature is algorithmic.
2. **Semi non-negative matrix factorization.** Relax the constraints of the  $K$ -means objective to  $z_{k,n} \in [0, \infty)$  and drop the constraint on the sum  $\sum_{k=1}^K z_{k,n}$ .

The motivation is to have probabilistic clustering model. We have a generative model:

- Goal: Explain the observed data  $\{\mathbf{x}_n\}_{n=1}^N$  by a probabilistic model  $p(\mathbf{x})$ .
- We assume the parametric form of the model to be chosen a-priori and look at the GMM.
- The model has *parameters* that need to be learned in order to explain the observed data well. For instance, the  $K$ -means centroids  $\mathbf{u}_k$  are an example of such parameters.

## 4.1 Probability Basics

### 4.1.1 Categorical Distribution

A categorical distribution is a discrete distribution over  $K$  events that assigns the probability  $\pi_k$  to the  $k$ th event. There are two codings for the random variable; the vector or index.

- Index coding. The random variable takes the index of the event that occurred.
- Vector coding. The random variable is a vector of size  $K$  and if event  $k$  occurs then the  $k$ th element is set to one, all other to zero.

We will use the vector coding, where the sample space is the set of 1-of- $K$  encoded random vectors  $\mathbf{z}$  of dimension  $K$  with

$$\sum_{k=1}^K z_k = 1 \quad \text{and} \quad z_k \in \{0, 1\}$$

The probability mass function is defined as

$$p(\mathbf{z} \mid \boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_k}$$

where  $\pi_k$  represents the probability of seeing element  $k$ . We have

$$\pi_k \geq 0 \quad \text{and} \quad \sum_{k=1}^K \pi_k = 1$$

### 4.1.2 Gaussian Distribution

The probability density function of the Gaussian distribution is given by

$$p(x \mid \mu, \sigma) = \mathcal{N}(x \mid \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

The variable  $\mu$  stands for the mean of the distribution and is at the same place as the mode and the median.  $\sigma^2$  is the variance. The probability for  $X \in [a, b]$  is given as

$$P(a < X < b) = \int_a^b p(x) dx$$

### 4.1.3 Multivariate Gaussian Distribution

The multivariate Gaussian distribution is a generalization to higher dimension with sample space  $\mathcal{X} \subseteq \mathbb{R}^D$ . A random variable  $\mathbf{X} = (X_1, \dots, X_D)^T$  has a multivariate normal distribution if every linear combination of its components (i.e.,  $Y = a_1 X_1 + \dots + a_D X_D$  has a univariate normal distribution. Its definition is

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(\sqrt{2\pi})^D |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

### 4.1.4 Multidimensional Moment Statistics

The expectation is the vector of component expectations

$$\mathbb{E}[\mathbf{X}] = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_D])^T$$

The variance can then be generalized to the *covariance*:

$$\text{Cov}(X, Y) = \int_{\mathcal{X}} \int_{\mathcal{Y}} p(x, y)(x - \mu_X)(y - \mu_Y) dx dy = \mathbb{E}_{X,Y}[(x - \mu_X)(y - \mu_Y)]$$

For random variables  $X_1, \dots, X_D$  we record the covariances in the covariance matrix  $\boldsymbol{\Sigma}$  as  $\Sigma_{ij} = \text{Cov}(X_i, X_j)$ . This generalizes the notation of variance to sets of random variables for multiple dimensions.

### 4.1.5 Conditional Probability

**Bayes' rule.** The conditional probability of  $A$  given  $B$  (posterior) is given by

$$p(A | B) = \frac{p(B | A)p(A)}{p(B)}$$

We call  $p(A)$  prior,  $p(B | A)$  likelihood and  $p(B)$  evidence.

## 4.2 Introduction to Mixture Models

Mixture of  $K$  probability densities is defined as

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k p(\mathbf{x} | \boldsymbol{\theta}_k)$$

Each probability distribution  $p(\mathbf{x} | \boldsymbol{\theta}_k)$  is a *component* of the mixture and has its own parameter  $\boldsymbol{\theta}_k$ . Almost any continuous density can be approximated using a sufficient number of component



distributions. For a Gaussian component distribution the parameters  $\theta_k$  are given by the mean  $\mu_k$  and the covariance  $\Sigma_k$ .

Mixture models are constructed from (1) the component distributions of the form  $p(\mathbf{x} | \theta_k)$  and (2) the mixing coefficient  $\pi_k$  that give the probability of each component. In order for  $p(\mathbf{x})$  to be a proper distribution we have to ensure that

$$\pi_k \geq 0 \quad \text{and} \quad \sum_{k=1}^K \pi_k = 1$$

Thus, the parameters  $\pi_k$  define a categorical distribution which represents the probability of each component distribution.

### 4.3 Gaussian Mixture Models

The Gaussian mixture model (GMM) uses Gaussians as the component distributions; thus

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)$$

Given data points  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  we then aim at learning the parameters  $\mu_k$  and  $\Sigma_k$ , such that we approximate the data as closely as possible. This is equivalent to finding the parameters that maximize the likelihood of the data.

**Generative viewpoint.** Given the model parameters  $\Sigma$ ,  $\mu$ , and  $\pi$ , sample the data  $\mathbf{x}_n$  as follows:

1. Sample a cluster index  $k$  according to the probabilities  $\pi_k$ .
2. Sample a data point  $\mathbf{x}_n$  from the distribution  $p(\mathbf{x}_n | \mu_k, \Sigma_k)$ .

Parameter estimation is then to revert this process and go from the data to the parameters.

To estimate the parameters, we assume that the data points  $\mathbf{x}$  are independent and identically distributed (i.i.d.). The probability or likelihood of the observed data  $\mathbf{X}$ , given the parameters can thus be written as

$$p(\mathbf{X} | \pi, \mu, \Sigma) = \prod_{n=1}^N p(\mathbf{x}_n) = \prod_{n=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)$$

We aim at finding at finding the parameters that maximize the likelihood of the data:

$$(\hat{\pi}, \hat{\mu}, \hat{\Sigma}) \in \arg \max_{\pi, \mu, \Sigma} p(\mathbf{X} | \pi, \mu, \Sigma)$$

To simplify the expression, we take the logarithm, such that the product becomes a sum:

$$(\hat{\pi}, \hat{\mu}, \hat{\Sigma}) \in \arg \max_{\pi, \mu, \Sigma} \sum_{n=1}^N \ln \left[ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) \right]$$

This equation does not have a closed-form analytic solution, and here we present a solution based on *expectation maximization*. The intuition is as follows: if we know to which clusters the data points are assigned, then computing the maximum likelihood is easy. First, we need to introduce a latent (or hidden) variable for the assignments of data points to clusters:

- We define a  $K$ -dimensional binary random variable  $\mathbf{z}$  having a 1-of- $K$  representation with

$$\sum_{k=1}^K z_k = 1 \quad \text{and} \quad z_k \in \{0, 1\}$$

- The marginal distribution over  $\mathbf{z}$  is specified in terms of the mixing coefficients  $\pi_k$ , i.e.,

$$p(z_k = 1) = \pi_k$$

Because  $\mathbf{z}$  uses a 1-of- $K$  representation, we can write this distribution in the form

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

Similarly, the conditional distribution of  $\mathbf{x}$  given a particular value for  $\mathbf{z}$  is a Gaussian

$$p(\mathbf{x} \mid z_k = 1) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

The conditional probability can be written as

$$p(\mathbf{x} \mid \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

We can use Bayes' rule to get the probability of assigning a data point to a cluster:

$$\begin{aligned} \gamma(z_k) &:= p(z_k = 1 \mid \mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x} \mid z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x} \mid z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \end{aligned}$$

#### 4.4 Expectation Maximization for GMM

*Expectation maximization* is a powerful method for finding maximum likelihood solutions for models with latent variables. Setting the derivatives of  $\ln p(\mathbf{X} \mid \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with respect to the means  $\boldsymbol{\mu}_k$  to zero and multiplying with  $\boldsymbol{\Sigma}_k$  we obtain:

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{k,n}) \mathbf{x}_n \quad N_k = \sum_{n=1}^N \gamma(z_{k,n})$$

The mean  $\boldsymbol{\mu}_k$  is obtained by taking a weighted mean of all the points in the data set. And further

$$\pi_k = \frac{N_k}{N}$$

Informally, we first choose some initial values for the means and mixing coefficients. Then, we alternate between two updates called the E and M step:

1. In the expectation step, the current values for the parameters are used to evaluate the posterior probability  $\gamma(z_{k,n})$ .

2. In the maximization step, these probabilities are used to re-estimate the means and mixing coefficients.

More formally, given a GMM, the goal is to maximize the likelihood function with respect to the parameters.

1. Initialize the means  $\boldsymbol{\mu}_k$ , and the mixing coefficients  $\pi_k$ . Set the  $\boldsymbol{\Sigma}_k$  to the given covariances.
2. **E-step.** Evaluate the responsibilities using the current parameter values:

$$\gamma(z_{k,n}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

3. **M-step.** Re-estimate the parameters using the current responsibilities.

$$\begin{aligned} \boldsymbol{\mu}_k^{\text{new}} &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{k,n}) \mathbf{x}_n \\ \pi_k^{\text{new}} &= \frac{N_k}{N} \quad N_k = \sum_{n=1}^N \gamma(z_{k,n}) \end{aligned}$$

4. Evaluate the log likelihood and check for convergence of either the parameters or the log likelihood.

## 5 Sampling

The basic idea is to *randomly* draw configurations from the probability distribution. Probable configurations appear frequently. This method is convenient, as

- solutions are in the hypothesis space,
- solutions obey to the probability distribution they are sampled from, and
- it works with difficult distributions (e.g., when normalization is infeasible).

### 5.1 Monte Carlo Methods

The first problem is to generate samples. That is, for a given probability distribution  $p(\mathbf{X})$ , generate samples  $\{\mathbf{x}^{(r)} | 1 \leq r \leq R\}$ . Instead of performing an integration we can evaluate a sum:

$$\frac{1}{R} \sum_{r=1}^R f(\mathbf{x}^{(r)}) \xrightarrow{R \rightarrow \infty} \int_{\Omega} f(\mathbf{x}) d^D \mathbf{x}$$

The second problem is expectation. Estimate the expectation of a given function  $f(\mathbf{x})$  under this distribution:

$$\mathbb{E}[f(\mathbf{x})] = \int p(\mathbf{x}) f(\mathbf{x}) d^D \mathbf{x}$$

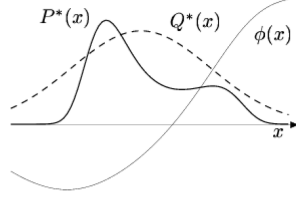


Figure 5: Importance Sampling.

### 5.1.1 Expectation from Samples

The empirical distribution is given as

$$\hat{p}(\mathbf{x}) = \frac{1}{R} \sum_{r=1}^R \delta(\mathbf{x} - \mathbf{x}^{(r)})$$

If we can solve the sampling problem, then we can solve the expectation problem by using the empirical mean

$$\hat{\mathbb{E}}[f(\mathbf{x})] = \int \hat{p}(\mathbf{x}) f(\mathbf{x}) d^D \mathbf{x} = \frac{1}{R} \sum_{r=1}^R f(\mathbf{x}^{(r)})$$

If  $\{\mathbf{x}^{(r)} \mid 1 \leq r \leq R\}$  are generated from  $p(\mathbf{x})$ , then the expectation of  $\hat{\mathbb{E}}[f(\mathbf{x})]$  is  $\mathbb{E}[f(\mathbf{x})]$  by the law of large numbers.

## 5.2 Importance Sampling

Importance sampling is a method for estimating the expectation of the function  $f(\mathbf{x})$ . Assume we can evaluate  $p^*(\mathbf{x})$  at any given point  $\mathbf{x}$ :

$$p(\mathbf{x}) = \frac{p^*(\mathbf{x})}{Z_p}$$

where  $Z_p$  is unknown. Assume that we cannot sample from  $p(\mathbf{x})$  but we have a simpler density  $q(\mathbf{x})$  from which we can sample.

$$q(\mathbf{x}) = \frac{q^*(\mathbf{x})}{Z_q}$$

is called the sampler density.

Some samples are under represented, other are over represented, so we correct by weights

$$w_r = \frac{p^*(\mathbf{x}^{(r)})}{q^*(\mathbf{x}^{(r)})}$$

which adjust the *importance* of each point. The expectation is thus

$$\hat{\mathbb{E}}[f(\mathbf{x})] = \frac{\sum_r w_r f(\mathbf{x}^{(r)})}{\sum_r w_r}$$

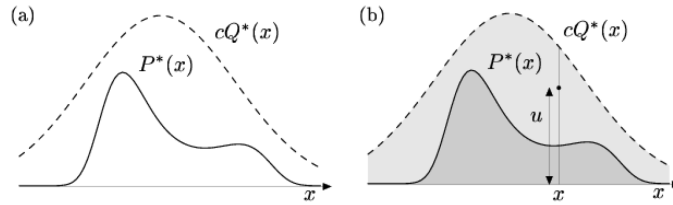


Figure 6: Rejection Sampling.

### 5.3 Rejection Sampling

Rejection sampling allows us to sample from relatively complex distributions. Assume that we can evaluate  $p^*(\mathbf{x})$  at any point  $\mathbf{x}$ ,

$$p(\mathbf{x}) = \frac{p^*(\mathbf{x})}{Z_p}$$

where  $Z_p$  is unknown. Assume we know a constant  $c$  such that

$$c \cdot q^*(x) > p^*(x) \quad \text{for all } x$$

The algorithm then is as follows:

1. Generate  $x$  from the proposal density  $q(x)$ .
2. Generate a uniformly distributed random variable  $u$  from the interval  $[0, c \cdot q^*(x)]$ .
3. If  $u > p^*(x)$  then  $x$  is rejected, otherwise it is added to our samples  $\{x^{(r)}\}$ .

Rejection sampling works best if  $q$  is a good approximation of  $p$ . If  $q$  is very different from  $p$ , then  $c$  may be very large (especially in high dimensions). Estimating the value of  $c$  may be hard, since we may not know where the modes of  $p^*$  are, or how high they may be.

## 6 Role Mining and Multi-Assignment Clustering

A simple way to manage access control is called *discretionary access control*. Assume we have  $N$  users and  $D$  permissions, then we use a binary access control matrix  $\mathbf{X} \in \{0, 1\}^{D \times N}$  where  $x_{di}$  indicates whether user  $i$  has access to resource  $d$ .

This simple idea has the problem of being tedious and expensive to maintain, and irregularities (i.e., potential errors) are very difficult to detect. Such irregularities might be special permissions that are given to someone on special-duty (need to be revoked later) or are simply mistakes.

Missing permissions are of course not a problem, but if a user has unnecessary permissions, someone might cause damage accidentally or on purpose. Thus, the *principle of least privilege* demands that users have only the privileges needed to fulfil their job. Detecting non-regular permissions and assessing their risk increases the security on IT systems.

### 6.1 Role-Based Access Control

Role-based access control (RBAC) simplifies access control and is, today, the method of choice. The basic entities of RBAC are:

- USERS, ROLES, OBJ, OPS, the set of users, roles, objects, and operations.
- PRMS  $\subseteq \{(obj, op) \mid obj \in \text{OBJ} \wedge op \in \text{OPS}\}$ , the set of permissions.
- UA  $\subseteq \text{USERS} \times \text{ROLES}$ , a user-role assignment relation.
- PA  $\subseteq \text{PRMS} \times \text{ROLES}$ , a permission-role assignment relation.
- UPA  $\subseteq \text{USERS} \times \text{PRMS}$ , a user-permission assignment relation.
- $\text{assigned\_users}(R) = \{u \in \text{USERS} \mid (u, R) \in \text{UA}\}$
- $\text{assigned\_permissions}(R) = \{p \in \text{PRMS} \mid (p, R) \in \text{PA}\}$

Each role defines a set of permissions. Users are assigned to a set of roles and get all permissions of these roles.

We define the boolean matrix product, denoted by  $\otimes$  as

$$\mathbf{X} = \mathbf{U} \otimes \mathbf{Z} \quad \iff \quad x_{di} = \bigvee_k (u_{dk} \wedge z_{ki})$$

### 6.2 Problem Definition

We are interested in finding good roles and role assignment based on a given user-permission matrix  $\mathbf{X}$ . That is, given  $\mathbf{X}$ , find  $\mathbf{U}$  and  $\mathbf{Z}$  such that

$$\mathbf{X} \approx \mathbf{U} \otimes \mathbf{Z}$$

with  $\mathbf{X} \in \{0, 1\}^{D \times N}$ ,  $\mathbf{U} \in \{0, 1\}^{D \times K}$ , and  $\mathbf{Z} \in \{0, 1\}^{K \times N}$ . Note that  $\mathbf{X}$  is typically noisy, i.e., it contains permission assignments without an underlying role.

There are three different definitions of the role-mining problem:

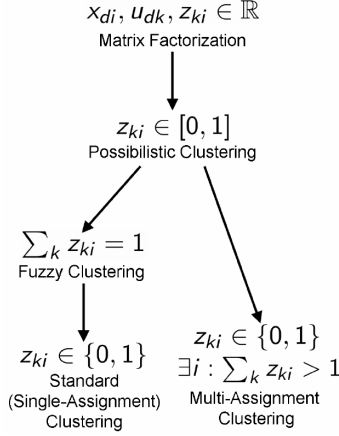


Figure 7: The matrix factorization landscape. Role mining is a multi-assignment clustering problem with  $x_{di}, u_{dk} \in \{0, 1\}$ .

1. **Min-noise approximation.** Given  $K$ , find matrices  $\mathbf{U} \in \{0, 1\}^{D \times K}$  and  $\mathbf{Z} \in \{0, 1\}^{K \times N}$  as

$$(\hat{\mathbf{U}}, \hat{\mathbf{Z}}) = \arg \min_{\mathbf{U}, \mathbf{Z}} \|\mathbf{X} - \mathbf{U} \otimes \mathbf{Z}\|_1$$

2.  **$\delta$ -approximation.** Given  $\delta \geq 0$ , find the minimal number of roles  $K$  and matrices  $\mathbf{U}$  and  $\mathbf{Z}$  such that

$$\|\mathbf{X} - \mathbf{U} \otimes \mathbf{Z}\|_1 < \delta$$

3. **Inference.** Assume that the given matrix  $\mathbf{X}$  has an underlying structure  $\mathbf{X}^S = \mathbf{U} \otimes \mathbf{Z}$  which was altered by a noise process  $\mathbf{X}^S \sim P(\mathbf{X}^S; \Theta^S)$ ,  $\mathbf{X} \sim P(\mathbf{X}; \mathbf{X}^S; \Theta^N)$ . Find the structure encoding in  $\mathbf{U}$  and  $\mathbf{Z}$ .

### 6.3 Evaluation Criteria for Boolean Matrix Decomposition

For each role mining problem definition, there is a set of evaluation criteria:

1. **Matrix reconstruction.** How accurately does the solution reconstruct the input matrix  $\mathbf{X}$ ?
2. **Number of sources.** How many sources  $K$  do you need to obtain a given accuracy?
3. **Inference quality.** How accurately can you infer the underlying structure?
4. **Generalization.** How precisely can you roles  $\mathbf{U}$  reconstruct a *new* matrix  $\mathbf{X}'$  (with the same distribution as  $\mathbf{X}$ )?
5. **Stability.** How consistent are the solutions of  $\mathbf{X}'$  and  $\mathbf{X}$ .

### 6.3.1 Evaluating a Matrix Reconstruction

To measure the discrepancy between  $\mathbf{X}$  and its inferred solution, we have several possibilities:

- **Deviation.** Compute (used to define the min-noise approximation):

$$\frac{1}{N \cdot D} \|\mathbf{X} - \mathbf{U} \otimes \mathbf{Z}\|_1$$

- **Coverage.** The ratio of true 1's which are reconstructed as 1's:

$$\text{Cov} := \frac{|\{(i, j) \mid \hat{x}_{ij} = x_{ij} = 1\}|}{|\{(i, j) \mid x_{ij} = 1\}|}$$

- **Deviating ones (zeros).** The number of deviating 1's (0's) is the number of matrix elements which are wrongly reconstructed divided by the total number of 1's (0's) in  $\mathbf{X}$ .

$$d_\xi := \frac{|\{(i, j) \mid \hat{x}_{ij} = 1 - \xi \wedge x_{ij} = \xi\}|}{|\{(i, j) \mid x_{ij} = \xi\}|} \quad \xi \in \{0, 1\}$$

Note that these measures have trivial optimal solutions, and if  $\mathbf{X}$  is noisy, one might not be interested in reconstructing every single bit.

### 6.3.2 Evaluation Inference Quality

Three more adequate (but more involved) measures are:

- **Parameter accuracy.** If you know the true parameters  $\Theta$  describing the roles, compare with your estimators  $\hat{\Theta}$ .

Compute the Hamming distance  $d_H$  between the estimated roles  $\hat{\mathbf{U}}$  and the true roles  $\mathbf{U}$ :

$$d_H(\mathbf{U}, \hat{\mathbf{U}}) := \min_{\pi \in P_K} \frac{1}{K \cdot D} \sum_{k=1}^K \sum_{d=1}^D |\hat{u}_{d,k} - u_{d,\pi(k)}|$$

where  $P_K$  is the set of all permutations over  $K$  elements. Note that the true source must be given, which is usually not the case.

- **Stability** of the clustering solution. Cluster two matrices  $\mathbf{X}, \mathbf{X}' \sim P(\mathbf{X}; \Theta)$ . Do you get similar clusterings?

In unsupervised learning, there is no fixed labelling of the clusters. Assume we have 2 clustering solutions  $\mathcal{K}_1$  and  $\mathcal{K}_2$  with clusters numbered  $1, \dots, K$  and they are defined by the elements they contain. Also, we know the costs  $c(k_1, k_2)$  for matching cluster  $k_1$  from  $\mathcal{K}_1$  to  $k_2$  from  $\mathcal{K}_2$ . To measure the match between  $\mathcal{K}_1$  and  $\mathcal{K}_2$  we need to find the optimal mapping  $\pi^*$  between the clusters, or more formally

$$\pi^* = \arg \min_{\pi} \left\{ \sum_k c(k, \pi(k)) \right\}$$

The permutation  $\pi^*$  can be efficiently found using the *Hungarian* method.



- **Generalization ability.** How well can you explain new entries with your estimated roles?

To measure how well the inferred set of roles generalizes to a new data matrix  $\mathbf{X}'$ , we can select a random subset  $\mathcal{D}^* \subset \{1, \dots, D\}$  of the permissions (usually approx. half of the permissions). For every column  $i$  of  $\mathbf{X}'$  (i.e., for every new user)

1. Compute  $\hat{\mathbf{z}}_{\cdot, i}$  the optimal role combination of user  $i$ :

$$\hat{\mathbf{z}}_{\cdot, i} := \arg \min_{\mathbf{z}_i} \left\{ \sum_{d \in \mathcal{D}^*} |x'_{d, i} - \mathbf{u}_{d, \cdot} \otimes \mathbf{z}_{\cdot, i}| \right\}$$

2. Compute the Hamming distance between the predicted and the actual permissions of user  $i$

$$g_i := \frac{1}{D - |\mathcal{D}^*|} \sum_{d \notin \mathcal{D}^*} |x'_{d, i} - \mathbf{u}_{d, \cdot} \otimes \hat{\mathbf{z}}_{\cdot, i}|$$

The average over all values is the generalization ability  $G$ :

$$G := \frac{1}{N} \sum_i g_i$$

## 6.4 SVD for Role Mining

Naively taking the SVD of  $\mathbf{X}$  and then rounding the left and right singular values to boolean values is a bad idea and does not work well. However, it can be used for *denoising*:

1. Compute the singular value decomposition:

$$\mathbf{X} = \mathbf{U} \cdot \mathbf{S} \cdot \mathbf{V}^T$$

2. Discard columns  $K + 1, \dots, D$  of  $\mathbf{U}$  to get  $\mathbf{U}_{(K)}$ , discard rows and columns  $K + 1, \dots, N$  of  $\mathbf{S}$  to get  $\mathbf{S}_{(K)}$  and discard rows  $K + 1, \dots, N$  of  $\mathbf{V}$  to get  $\mathbf{V}_{(K)}$ .
3. Round the reconstruction

$$\mathbf{X}_{(K)} := \mathbf{U}_{(K)} \cdot \mathbf{S}_{(K)} \cdot \mathbf{V}_{(K)}^T$$

to get

$$\hat{\mathbf{X}}_{(K)} = (\mathbf{X}_{(K)} > t_{\mathbf{X}})$$

where  $t_{\mathbf{X}}$  is a threshold.

For this to work, we actually need three values:

- a threshold  $t_{\mathbf{X}}$  to be chosen depending on the data
- a *hyper-parameter*  $K$ , the number of dimensions
- an estimate for the reconstruction error

## 6.5 Cross-Validation

For choosing the right parameters, we need three disjoint data sets. If there is only one data set available, then one needs to split it into three parts:

- The *training data* is used to estimate the best value of the parameters  $\theta$ . For the application above,  $\theta = t_X$ , given a fixed  $K$ .
- The *validation data* is used to determine the best value of the hyper-parameters  $\theta^H$  (for the best value of  $\theta$ ). In our setting  $\theta^H = K$ .
- The *test data* is then used to estimate the reconstruction error.

On a high level, the algorithm for choosing a hyper-parameter works as follows:

1. Get three disjoint data (sub)sets.
2. For all possible values of the hyper-parameter  $\theta^H$ :
  - i) Fix  $\theta^H$ , and
  - ii) *train* your method on the training data to get  $\theta^*$ , and
  - iii) *validate* your method on the validation data.

3. Set

$$\theta^{H*} = \left( \begin{array}{l} \text{Value of } \theta^H \text{ that maximized the performance on the} \\ \text{validation data (with respective best value for } \theta \text{)} \end{array} \right)$$

Note that the test data is *not* used in the complete process.

Once the optimal hyper-parameter  $\theta^{H*}$  has been found, the following procedure can be used to evaluate the model:

1. Fix  $\theta^{H*}$ .
2. Train your method to get the optimal value  $\theta^*$  on the training data and the validation data.
3. Test your method with  $\theta^{H*}$  and  $\theta^*$  on the *test data*. This is the final result of your method.

It is (implicitly) assume that the data in the three parts has the *same distribution*. The division into three parts is random and reduces the number of training data. To avoid too much dependency on this choice, one often uses several splits. This is called *cross-validation*.

Cross-validation allows one to estimate how well a method will perform on a new, unseen data set. This procedure avoids overfitting, i.e., a too detailed adaptation to the training data which renders the method very specific to the training data. Therefore, the test data is only considered at the very end.

## 6.6 Simple Approach Using a Variant of $K$ -means

We can use the idea of  $K$ -means and adapt the algorithm to boolean data:

- Adapt distance to *Hamming distance* (0-norm):

$$R(\mathbf{U}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{U}\mathbf{Z}\|_0 = \sum_{i=1}^N \sum_{k=1}^K z_{ki} \|\mathbf{x}_i - \mathbf{u}_k\|_0$$

- Restrict centroids  $\mathbf{u}_k$  to boolean values and adapt the centroid update step:

$$u_{kd} = \text{median}(\{x_{dn} \mid z_{kn} = 1\}) \quad \text{for all } k \text{ and } d$$

However, this approach has several problems:

- The parameter  $K$  has to be given (use cross-validation)
- $K$ -means yields *disjoint* clustering, but in RBAC a user can have more than one role.

## 6.7 The ROLEMINER Algorithm

ROLEMINER is one of the first role mining approaches (heuristic) and comes in two variants; COMPLETEMINER and FASTMINER.

The algorithms uses the observation that users are described by sets of permissions, and that roles are defined as sets of permissions. The key idea is that a user's permissions are generated by a combination of roles, and thus all possible intersections of permission sets of all users should give the underlying components (the *atomic* roles, so to say). For instance, for users (represented as permission sets)  $\{1, 2, 5\}$  and  $\{2, 3, 5\}$  we would have the following candidate roles:  $\{1, 2, 5\}$ ,  $\{2, 3, 5\}$  and  $\{2, 5\}$ .

---

### Algorithm 1 The COMPLETEMINER algorithm

---

**Input:**  $x$

- 1: Extract all unique permission sets (columns)  $\mathbf{x}'$  from the input.
  - 2: Compute all unique set intersections between the members of  $\mathbf{x}'$  and add them to the set of  $k$  candidate roles  $\mathbf{u}$ .
  - 3: **for** each candidate in  $\mathbf{u}$  **do**
  - 4:   Count how many users have the permission set of the candidate as a subset.
  - 5: **end for**
  - 6: **return** The prioritized set of  $k$  roles  $\mathbf{u}$ , where the counts  $c = (c_1, \dots, c_k)$  give the priority.
- 

A variant of the COMPLETEMINER is FASTMINER, where only the set of *pairwise* intersections is computed. If  $t$  is the number of initial roles (i.e., the number of users with unique permission sets), then COMPLETEMINER has to compute  $2^t - t - 1$  set intersections, and FASTMINER only  $\frac{t(t-1)}{2}$ .

Clearly, there are various problems with the ROLEMINER algorithm:

- The algorithm is very sensitive to noise:
  - If only a few individual bits are noisy, the number of candidate roles gets much larger.
  - Slightly different noise leads to a completely different solution.
- No assignments of users to roles, only candidate for the columns of  $\mathbf{U}$  are given ( $\mathbf{Z}$  is neglected).

## 6.8 DBPSOLVER

DBPSOLVER approximately solves the *discrete basis problem*.

**Definition 6.1.** For a given boolean matrix  $\mathbf{X} \in \{0, 1\}^{D \times N}$  and a number  $K$  of basis vectors, find a boolean matrix  $\mathbf{U} \in \{0, 1\}^{D \times K}$  and a boolean matrix  $\mathbf{Z} \in \{0, 1\}^{K \times N}$  minimizing

$$\|\mathbf{X} - \mathbf{U} \otimes \mathbf{Z}\|_1$$

where  $\|\cdot\|_1$  is the  $L_1$  norm:  $\|M\|_1 = \sum_{i=1}^m \sum_{n=1}^n |M_{ij}|$

Note that the discrete basis problem and the min-noise role mining problem are identical.

The idea of DBPSOLVER is the following. To cover most of the 1's in  $\mathbf{X}$  with as few basis vectors as possible, candidates should have 1's at dimensions (e.g., permissions) that are frequently 1 simultaneously.

DBPSOLVER computes a  $D \times D$  association matrix  $\mathbf{A}$  where  $a_{dd'}$  is the *pairwise association* between dimensions (e.g., permissions)  $d$  and  $d'$ :

$$a_{dd'} := a(d \Rightarrow d') = \frac{\langle \mathbf{x}_{d,\cdot}, \mathbf{x}_{d',\cdot} \rangle}{\langle \mathbf{x}_{d,\cdot}, \mathbf{x}_{d,\cdot} \rangle}$$

We have  $a_{dd'} = \hat{p}\{d' = 1 \mid d = 1\}$ . In role mining  $a_{dd'}$  is the empirical probability of having permission  $d'$  conditioned on having  $d$ .

The algorithm now works greedily as follows:

1. Compute the association matrix  $\mathbf{A}$ .
2. Compute  $\mathbf{A}' = (\mathbf{A} > \tau)$  by thresholding  $\mathbf{A}$  with a predefined  $\tau$ . Each column in  $\mathbf{A}'$  is a candidate basis vector. In role mining, it represents a candidate role.
3. Initialize  $\mathbf{U} = \{0\}^{D \times K}$  and  $\mathbf{Z} = \{0\}^{K \times N}$ .
4. For  $k = 1, \dots, K$ : substitute column  $k$  in  $\mathbf{U}$  with column  $t$  from  $\mathbf{A}'$  and modify  $\mathbf{z}_k$  such that

$$R(\mathbf{X}, \mathbf{Z}, \mathbf{U}, w^+, w^-) = w^+ \cdot \left| \{(i, d) \mid x_{di} = 1 \wedge (\mathbf{U}_{(t)} \otimes \mathbf{Z}_{(t)})_{di} = 1\} \right| - w^- \cdot \left| \{(i, d) \mid x_{di} = 0 \wedge (\mathbf{U}_{(t)} \otimes \mathbf{Z}_{(t)})_{di} = 1\} \right|$$

is maximized.

$w^+$  and  $w^-$  are predefined parameters that reward and penalize covering 0's and 1's with a 1.

DBPSOLVER has the following problems:

- The parameters  $\tau$ ,  $w^+$ ,  $w^-$  and  $K$  must be chosen (cross-validation).
- Greedy optimization does not find the *best* combination of basis vectors, but only picks a single best candidate at every step and ignores possible combinations with other vectors.
- Finding the best matrix approximation differs from discovering structure and noise! The best fit to  $\mathbf{X}$  might include noise (“overfitting”). This cannot be avoided and is intrinsic to the problem being solved (i.e., we asked the wrong question).

## 7 Sparse Coding

A signal and its representation are not the same thing, and there is an infinite number of possible representations.

Given a signal  $f$  and an orthonormal basis  $\{u_1, \dots, u_L\}$ , we can represent the signal by coefficients

$$z_l = \langle f, u_l \rangle$$

The signal  $f$  can then be reconstructed as

$$f = \sum_{l=1}^L z_l u_l = \sum_{l=1}^L \langle f, u_l \rangle u_l$$

In sparse coding, we hope that we get only few non-zero coefficients. We can compress by dropping some of the coefficients and only using certain basis functions for a subset  $\sigma$ :

$$\hat{f} = \sum_{k \in \sigma} z_k u_k$$

In this case, the reconstruction error is defined as

$$\|f - \hat{f}\|^2 = \sum_{k \notin \sigma} |\langle f, u_k \rangle|^2$$

since  $\langle u_k, u_l \rangle = 0$  for  $k \neq l$ .

Of course, the choice of basis is critical, and there is no choice of transform that is better than all others, as it depends on the type of the signal. For instance, the Fourier basis has global support and is good for “sine like” signals, but performs poor for localized signals. On the other hand, a Wavelet basis has local support but is poor for non-vanishing signals.

### 7.1 Comparison to Principle Component Analysis

We can compress a set of vectors  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  by performing PCA. That is, we compute the (centred) covariance matrix

$$\mathbf{\Sigma} = \frac{1}{n} (\mathbf{X} - \mathbf{M})(\mathbf{X} - \mathbf{M})^T$$

and compute its eigenvector decomposition, giving us  $\mathbf{U}$  and  $\mathbf{\Lambda}$ . We then choose the  $k$  eigenvectors corresponding to the largest eigenvalues as  $\mathbf{U}_k$ , and compress a given signal  $\mathbf{x}$  as

$$\hat{\mathbf{z}} = \mathbf{U}_k \mathbf{x}$$

*PCA basis:* Since  $\mathbf{U}_k$  is dependent on the data, we need to transmit both the basis  $\mathbf{U}$  and the elements of  $\hat{\mathbf{z}}$ . However, the basis is optimal for the given covariance.

On the other hand, using a *fixed basis*, both parties (compressor and decompressor) can agree upon a particular basis and then transmit only the non-zero elements of the transformed signal  $\hat{\mathbf{z}}$ .

## 7.2 Compressive Sensing

When we gather information (e.g., in photography, MRI technology or other applications), we do not want to collect huge amounts of data if we compress most of it away right afterwards. Instead, compressing the data while gathering it is called *compressive sensing* and helps to decrease acquisition time, required battery and storage.

Given an original signal  $\mathbf{x} \in \mathbb{R}^D$  which is assumed to be sparse in some orthonormal basis  $\mathbf{U}$  with  $K$  large coefficients in its sparse representation  $\mathbf{z}$ :

$$\mathbf{x} = \mathbf{U}\mathbf{z} \quad \text{such that } |\mathbf{I}_{\mathbf{z}} \gg 0| = K$$

The main idea is to acquire the set of  $M$  linear combinations  $y_k = \langle w_k, \mathbf{x} \rangle$  of the initial signal instead of the signal itself, where  $M \ll D$  and  $\mathbf{W} = \{w_1, \dots, w_M\}$  is a matrix with linear combinations coefficients

$$y_k = \langle w_k, \mathbf{x} \rangle \quad k \in M$$

$$\mathbf{y} = \mathbf{W}\mathbf{x} = \mathbf{W}\mathbf{U}\mathbf{z} \stackrel{\Theta = \mathbf{W}\mathbf{U}}{=} \Theta\mathbf{z}$$

Let all matrix elements  $w_{i,j}$  be independent and i.i.d. random variables from a Gaussian probability density function with mean zero and variance  $\frac{1}{D}$ . Surprisingly, regardless of the choice of orthonormal basis  $\mathbf{U}$  then  $M$

$$M \leq cK \log\left(\frac{d}{K}\right)$$

where  $c$  is some constant, is enough to obtain a stable solution for any  $K$ -sparse and compressible signal.

Reconstruction the initial signal of length  $D$  having  $M$  measurements means inverting the following system:

$$\mathbf{y} = \mathbf{W}\mathbf{x} = \mathbf{W}\mathbf{U}\mathbf{z} = \Theta\mathbf{z}$$

This problem seems to be *ill-posed*, as  $M \ll D$  we have many more unknowns than equations. We can look at the sparsest solution such that

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0 \quad \text{such that } \Theta\mathbf{z} = \mathbf{y}$$

This can be solved using the *matching pursuit* algorithm

## 7.3 Summary of Sparse Coding

Sparse coding works by coding via orthogonal transforms. That is, given an original signal  $\mathbf{x}$  and an orthogonal matrix  $\mathbf{A}$ , we compute the linear transformation (change of basis)  $\mathbf{z} = \mathbf{A}\mathbf{x}$  and truncate “small” values, giving us  $\hat{\mathbf{z}}$ . The inverse is then computed as  $\hat{\mathbf{x}} = \mathbf{A}^T \hat{\mathbf{z}}$  (since  $\mathbf{A}^{-1} = \mathbf{A}^T$ ).

To measure the performance, we can measure the *reconstruction error*  $\|\mathbf{x} - \hat{\mathbf{x}}\|$  or measure the *sparsity* of the coding vector  $\mathbf{z}$ , i.e.,  $\text{nnz}(\mathbf{z})$ .

## 7.4 Overcomplete Dictionaries

No single basis (e.g., Fourier or Wavelets) is optimally sparse for *all* signal classes. We could use more atoms (dictionary elements) than dimensions, i.e., make use of *overcompleteness* ( $L > D$ ).

We have two possibilities:

1. We use a union of orthogonal bases  $[U_1, \dots, U_B] \in \mathbb{R}^{D \times D \cdot D}$ :

$$\mathbf{x} = [U_1, \dots, U_B] \cdot \mathbf{z}$$

The coding vector  $\mathbf{z}$  is now also larger, but we hope that it is very sparsely populated.

2. We can also use a general overcomplete dictionary with some number  $L > D$  of atoms.

When we increase the overcompleteness factor  $\frac{L}{D}$  we increase (potentially) the sparsity of the coding, but also increase the linear dependency between atoms. We measure this linear dependency using the notion of *coherence*:

$$m(\mathbf{U}) = \max_{i \neq j} |\mathbf{u}_i^T \mathbf{u}_j|$$

We have  $m(\mathbf{B}) = 0$  for an orthogonal basis  $\mathbf{B}$ , and  $m([\mathbf{B}, \mathbf{u}]) \geq \frac{1}{\sqrt{D}}$  if an atom  $\mathbf{u}$  is added to  $\mathbf{B}$ .

#### 7.4.1 Signal Coding

Given an overcomplete dictionary  $\mathbf{U} \in \mathbb{R}^{D \times L}$ , solving the equation

$$\mathbf{x} = \mathbf{U} \cdot \mathbf{z}$$

is an ill-posed problem, as we have more unknowns than equations. We can add the constraint that we want to find the sparsest  $\mathbf{z}$ . We solve

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0 \quad \text{such that} \quad \mathbf{x} = \mathbf{U}\mathbf{z}$$

This is the same formulation as for compressive sensing.

#### 7.4.2 Noisy Observations

**Additive noise.** Assume each dimension is independently corrupted by zero-mean Gaussian noise with variance  $\sigma^2$ :

$$\mathbf{x} = \mathbf{U}\mathbf{z} + \mathbf{n} \quad \text{with} \quad n_d \sim \mathcal{N}(0, \sigma^2)$$

Then, we want to find

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0 \quad \text{such that} \quad \|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2 < \sigma$$

That is, we maximize the sparsity of  $\mathbf{z}$  while the residual (approximation error) remains below  $\sigma$ .

Alternatively, we can also solve

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2 \quad \text{such that} \quad \|\mathbf{z}\|_0 \leq K$$

That is, we minimize the residual while selecting less than  $K$  atoms from the dictionary.

**Approximate sparse coding.** We can also code the noise in terms of  $\mathbf{U}$ :

$$\mathbf{x} = \mathbf{U}\mathbf{z} + \mathbf{n} = \mathbf{U}\mathbf{z} + \mathbf{U}\mathbf{y} = \mathbf{U}(\mathbf{z} + \mathbf{y})$$

where  $\mathbf{y}$  is the coding of the noise  $\mathbf{n}$  in  $\mathbf{U}$ . The noise cannot be sparsely coded in any dictionary, therefore  $\mathbf{y}$  has many small coefficients. However, the few large coefficients are still due to  $\mathbf{z}$ .

**Geometry of sparse coding solution  $\mathbf{z}^*$ .** As shown in [Figure 8](#), the orthogonal projection of  $\mathbf{x}$  onto the subspace spanned by selected atoms  $\{\mathbf{u}_d \mid z_d^* \neq 0\}$  minimizes  $\|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2$ .

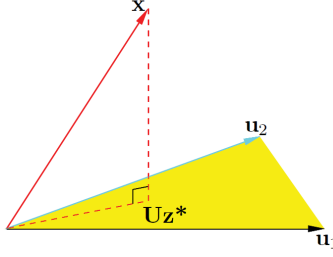


Figure 8: Geometry of sparse coding solution  $\mathbf{z}^*$ .

### 7.4.3 Matching Pursuit (MP)

We consider a greedy algorithm to approximate the NP hard problem iteratively, each time taking the action which is short-term optimal. Applied to sparse coding, we do the following:

1. Start with the zero vector  $\mathbf{z} = \mathbf{0}$  and residual  $\mathbf{r}^0 = \mathbf{x}$ .
2. At each iteration  $t$ , take a step in the direction of the atom  $\mathbf{u}_{d^*(t)}$  that maximally reduces the residual  $\|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2$ .

The atom selection in iteration  $t$  is done as follows:

$$d^*(t) = \arg \max_d |\langle \mathbf{r}^t, \mathbf{u}_d \rangle|$$

The objective of the algorithm is (as seen before):

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2 \quad \text{such that} \quad \|\mathbf{z}\|_0 \leq K$$

---

**Algorithm 2** The matching pursuit algorithm

---

**Objective:**  $\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{x} - \mathbf{U}\mathbf{z}\|_2$ , such that  $\|\mathbf{z}\|_0 \leq K$

1:  $\mathbf{z} \leftarrow \mathbf{0}; \mathbf{r} \leftarrow \mathbf{x}$

2: **while**  $\|\mathbf{z}\|_0 < K$  **do**

3:   Select the atom with the maximum absolute correlation to the residual:

$$d^* = \arg \max_d |\langle \mathbf{r}, \mathbf{u}_d \rangle|$$

4:   Update the coefficient vector and the residual:

$$z_{d^*} \leftarrow z_{d^*} + \mathbf{u}_{d^*}^T \mathbf{r}$$

$$\mathbf{r} \leftarrow \mathbf{r} - (\mathbf{u}_{d^*}^T \mathbf{r}) \mathbf{u}_{d^*}$$

5: **end while**

---

### 7.4.4 Support Recovery of MP

Assume  $\mathbf{x} = \mathbf{U}\mathbf{z}$  has a  $K$  sparse coding

$$\mathbf{x} = \sum_{d \in \Delta} z_d \mathbf{u}_d$$



Under which condition is MP successful in recovering the true *support*, i.e.,  $\Delta_{\text{MP}} = \Delta$ , where

$$\Delta_{\text{MP}} = \{d \mid z_d^* \neq 0\}$$

is the set of coding atoms for the MP solution  $\mathbf{z}^*$ . The *exact recovery condition* is

$$K < \frac{1}{2} \left( 1 + \frac{1}{m(\mathbf{U})} \right)$$

where  $m(\mathbf{U})$  is the coherence. The intuition is that if the coherence is small, then explaining a generating atom with other atoms is not sparse. Therefore, sparse coding recovers the support. We have a trade-off for increasing  $L$ ; this leads to sparser coding (smaller possible  $K$ ) but increases the coherence  $m(\mathbf{U})$ .

#### 7.4.5 Sparse Coding for Inpainting

To use sparse coding for image inpainting, we can sparse code the known parts of the images, and then predict the missing parts by reconstructing from the sparse code. An algorithm by Elad et al. (2005) works as follows:

1. Define the diagonal masking matrix  $\mathbf{M}$  where  $m_{d,d} = 1$  if the pixel  $d$  is known and  $m_{d,d} = 0$  if it is missing.
2. Sparse coding of the known parts in an overcomplete dictionary  $\mathbf{U}$ :

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0 \quad \text{such that} \quad \|\mathbf{M}(\mathbf{x} - \mathbf{U}\mathbf{z})\|_2 < \sigma \quad (1)$$

3. Image reconstruction using the mask:

$$\hat{\mathbf{x}} = \mathbf{M}\mathbf{x} + (\mathbf{I} - \mathbf{M})\mathbf{U}\mathbf{z}^*$$

The question that remains is how we can modify MP to solve [Equation 1](#). We can take an EM algorithm, as follows:

1. Assume initial sparse coding  $\mathbf{z}^0$  of the image.
2. Reconstruct the image using the mask

$$\hat{\mathbf{x}}^1 = \mathbf{M}\mathbf{x} + (\mathbf{I} - \mathbf{M})\mathbf{U}\mathbf{z}^0$$

3. Sparse coding of reconstructed image  $\mathbf{x}^1$ :

$$\mathbf{z}^1 = \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0 \quad \text{such that} \quad \|\mathbf{x}^1 - \mathbf{U}\mathbf{z}\|_2 < \sigma$$

4. Iterate steps two and three  $t$  times until convergence.
5. Image reconstruction using the mask:

$$\hat{\mathbf{x}} = \mathbf{M}\mathbf{x} + (\mathbf{I} - \mathbf{M})\mathbf{U}\mathbf{z}^t$$

## 7.5 Dictionary Learning

So far, we have seen two variants:

- Fixed orthonormal basis. The advantage is that coding is efficient and can be done by matrix multiplication  $\mathbf{z} = \mathbf{U}^T \mathbf{x}$ . However, the result is only sparse for a single class of signals.
- Fixed overcomplete basis. Now, sparse coding is possible for many signal classes, but finding the sparsest code requires an approximation algorithm and is problematic if the dictionary size  $L$  and the coherence  $m(\mathbf{U})$  are large.

Instead, we now consider *learning the dictionary*. This allows us to adapt the dictionary to our signal's characteristics. However, we have to solve a matrix factorization problem

$$\mathbf{X}_{D \times N} \approx \mathbf{U}_{D \times L} \cdot \mathbf{Z}_{L \times N}$$

subject to the sparsity constraint on  $\mathbf{Z}$  and atom norm constraint for  $\mathbf{U}$ .

We have the following objective

$$\arg \min_{\mathbf{U}, \mathbf{Z}} \|\mathbf{X} - \mathbf{U} \cdot \mathbf{Z}\|_F^2$$

where we use the Frobenius norm  $\|\mathbf{R}\|_F^2 = \sum_{i,j} r_{i,j}^2$  (sum-square error). This objective is not convex in both  $\mathbf{U}$  and  $\mathbf{Z}$  (local minima), but convex in either  $\mathbf{U}$  or  $\mathbf{Z}$  (unique minimum). We can use an iterative greedy minimization:

1. Coding step:  $\mathbf{Z}^{(t+1)} = \arg \min_{\mathbf{Z}} \|\mathbf{X} - \mathbf{U}^{(t)} \cdot \mathbf{Z}\|_F^2$ , subject to  $\mathbf{Z}$  being sparse.
2. Dictionary update step:  $\mathbf{U}^{(t+1)} = \arg \min_{\mathbf{U}} \|\mathbf{X} - \mathbf{U} \cdot \mathbf{Z}^{(t+1)}\|_F^2$  subject to  $\|\mathbf{u}_d\|_2 = 1$  for all  $d = 1, \dots, D$ .

### 7.5.1 Coding Step

The residual is column separable, i.e.,

$$\|\mathbf{R}\|_F^2 = \sum_{i,j} r_{i,j}^2 = \sum_j \|\mathbf{r}_j\|_2^2$$

Thus, we can perform  $N$  *independent* sparse coding steps. For all  $n = 1, \dots, N$ :

$$\mathbf{z}_n^{(t+1)} = \arg \min_{\mathbf{z}} \|\mathbf{z}\|_0 \quad \text{such that} \quad \|\mathbf{x}_n - \mathbf{U}^{(t)} \mathbf{z}\|_2 \leq \sigma \cdot \|\mathbf{x}_n\|_2$$

### 7.5.2 Dictionary Update

Unfortunately, the residual for the dictionary update step, i.e.,

$$\mathbf{U}^{(t+1)} = \arg \min_{\mathbf{U}} \|\mathbf{X} - \mathbf{U} \cdot \mathbf{Z}^{(t+1)}\|_F^2$$

is not separable in atoms (columns of  $\mathbf{U}$ ). Therefore, we approximate by updating one atom at a time. For all  $d = 1, \dots, D$

1. Fix all atoms except  $\mathbf{u}_d$ :

$$\mathbf{U} = [\mathbf{u}_1^{(t)}, \dots, \mathbf{u}_d, \dots, \mathbf{u}_D^{(t)}]$$

2. Isolate  $\mathbf{R}_d^{(t)}$ , the residual that is due to atom  $\mathbf{u}_d$ :

$$\begin{aligned} & \left\| \mathbf{X} - [\mathbf{u}_1^{(t)}, \dots, \mathbf{u}_d, \dots, \mathbf{u}_D^{(t)}] \cdot \mathbf{Z}^{(t+1)} \right\|_F^2 \\ &= \left\| \mathbf{X} - \left\{ \sum_{e \neq d} \mathbf{u}_e^{(t)} (\mathbf{z}_e^{(t+1)})^T + \mathbf{u}_d^{(t)} (\mathbf{z}_d^{(t+1)})^T \right\} \cdot \mathbf{Z}^{(t+1)} \right\|_F^2 \\ &= \left\| \mathbf{R}_d^{(t)} - \mathbf{u}_d (\mathbf{z}_d^{(t+1)})^T \right\|_F^2 \end{aligned}$$

3. Find  $\mathbf{u}_d^*$  that minimizes  $\mathbf{R}_d^{(t)}$ , subject to  $\|\mathbf{u}_d^*\|_2 = 1$ :

- $\mathbf{u}_d (\mathbf{z}_d^{(t+1)})^T$  is an outer product, i.e., a matrix.
- Minimizing the residual

$$\left\| \mathbf{R}_d^{(t)} - \mathbf{u}_d (\mathbf{z}_d^{(t+1)})^T \right\|_F^2$$

by approximating  $\mathbf{R}_d^{(t)}$  with rank-1 matrix  $\mathbf{u}_d (\mathbf{z}_d^{(t+1)})^T$ .

- Achieved by the SVD of  $\mathbf{R}_d^{(t)}$ :

$$\mathbf{R}_d^{(t)} = \tilde{\mathbf{U}} \mathbf{S} \tilde{\mathbf{V}}^T = \sum_i s_i \tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_i^T$$

- $\mathbf{u}_d^* = \tilde{\mathbf{u}}_1$  is the first left-singular vector.
- $\|\mathbf{u}_d^*\|_2 = 1$  is naturally satisfied.

### 7.5.3 Initialization

The approach is sensitive to the choice of  $\mathbf{U}^{(0)}$ ; the initial candidate solution is optimized locally and greedily until no progress is possible any more. Possible choices include:

1. **Random atoms.** Sampling  $\{\mathbf{u}_d^{(0)}\}$  on the unit sphere.
  - i) Sample  $D$  dimensional standard normal distribution:  $\mathbf{u}_d^{(0)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$ .
  - ii) Scale to unit length:

$$\mathbf{u}_d^{(0)} \leftarrow \frac{\mathbf{u}_d^{(0)}}{\|\mathbf{u}_d^{(0)}\|_2}$$

2. **Samples from  $\mathbf{X}$ .**

- i)  $\mathbf{u}_d^{(0)} \leftarrow \mathbf{x}_n$  where  $n \sim \mathcal{U}(1, N)$  is sampled uniformly.
- ii) Scale to unit length.

3. Use **fixed overcomplete dictionary**, e.g., overcomplete DCT.