

Advanced Machine Learning  
from Theory to Practice  
Lecture 4  
Dimension Reduction and Feature Design

F. d'Alche-Buc and E. Le Pennec

Fall 2015

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

# Dimension Reduction

## Dimension Reduction

- **Training data** :  $\mathcal{D} = \{\mathbf{X}_1, \dots, \mathbf{X}_n\} \in \mathcal{X}^n$  (i.i.d.  $\sim \mathbf{P}$ )
- Space  $\mathcal{X}$  of possibly high dimension.

### Dimension Reduction Map

- Construct a map  $\Phi$  from the space  $\mathcal{X}$  into a space  $\mathcal{X}'$  of **smaller dimension** :

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ \mathbf{X} &\mapsto \Phi(\mathbf{X})\end{aligned}$$

### Criterion

- Reconstruction error
- Distance preservation

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

# Dimension Reduction

## Reconstruction Error Approach

### Goal

- Construct a map  $\Phi$  from the space  $\mathcal{X}$  into a space  $\mathcal{X}'$  of **smaller dimension** :

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ \mathbf{X} &\mapsto \Phi(\mathbf{X})\end{aligned}$$

- Construct  $\tilde{\Phi}$  from  $\mathcal{X}'$  to  $\mathcal{X}$
  - Control the error between  $\mathbf{X}$  and its reconstruction  $\tilde{\Phi}(\Phi(\mathbf{X}))$
- 
- Canonical example for  $\mathbf{X} \in \mathbb{R}^d$  : find  $\Phi$  and  $\tilde{\Phi}$  in a parametric family that minimize

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\Phi}(\Phi(\mathbf{x}_i))\|^2$$

# Dimension Reduction

## Principal Component Analysis

- $\mathcal{X} \in \mathbb{R}^d$  and  $\mathcal{X}' = \mathbb{R}^{d'}$
- Affine model  $\mathbf{X} \sim m + \sum_{l=1}^{d'} \mathbf{X}'_l V^{(l)}$  with  $(V^{(l)})$  an orthonormal basis.
- Equivalent to :

$$\Phi(\mathbf{X}) = V^t(\mathbf{X} - m) \quad \text{and} \quad \tilde{\Phi}(\mathbf{X}') = m + V\mathbf{X}'$$

- Reconstruction error criterion :

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{X}_i - (m + VV^t(\mathbf{X}_i - m))\|^2$$

- **Explicit solution** :  $m$  is the empirical mean and  $V$  is any orthonormal basis of the space spanned by the  $d'$  first eigenvectors (the one with largest eigenvalues) of the empirical covariance matrix  $\frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - m)(\mathbf{X}_i - m)^t$ .

# Dimension Reduction

## Principal Component Analysis

### PCA Algorithm

- Compute the empirical mean  $m = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$
  - Compute the empirical covariance matrix  $\frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - m)(\mathbf{X}_i - m)^t$ .
  - Compute the  $d'$  first eigenvectors of this matrix :  $V^{(1)}, \dots, V^{(d')}$
  - Set  $\Phi(\mathbf{X}) = V^t(\mathbf{X} - m)$
- 
- Complexity :  $O(n(1 + d^2) + d'd^2)$
  - Interpretation :
    - $\Phi(\mathbf{X}) = V^t(\mathbf{X} - m)$  : coordinates in the restricted space.
    - $V^{(i)}$  : influence of each original coordinates in the  $i$ th new one.
  - **Scaling** : This method is not invariant to a scaling of the variables ! It is custom to normalize the variables (at least within groups) before applying PCA.



# Dimension Reduction

## Multiple Factor Analysis

- PCA assumes  $\mathcal{X} = \mathbb{R}^d$  !
- How to deal with categorical values ?
- MFA = PCA with clever coding strategy for categorical values.

### Categorical value code for a single variable

- Classical redundant dummy coding :

$$\mathbf{X} \in \{1, \dots, V\} \mapsto P(\mathbf{X}) = (\mathbf{1}_{\mathbf{X}=1}, \dots, \mathbf{1}_{\mathbf{X}=V})^t$$

- Compute the mean (i.e. the empirical proportion)  $\bar{P} = \frac{1}{n}P(\mathbf{X})$
- Renormalize  $P(\mathbf{X})$  by  $1/\sqrt{\bar{P}}$  :

$$P(\mathbf{X}) = (\mathbf{1}_{\mathbf{X}=1}, \dots, \mathbf{1}_{\mathbf{X}=V}) \mapsto P^r(\mathbf{X}) = \left( \frac{\mathbf{1}_{\mathbf{X}=1}}{\sqrt{\bar{P}_1}}, \dots, \frac{\mathbf{1}_{\mathbf{X}=V}}{\sqrt{\bar{P}_V}} \right)$$

- $\chi^2$  type distance !

# Dimension Reduction

## Multiple Factor Analysis

- PCA becomes the minimization of

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \|P^r(\mathbf{X}_i) - (m + VV^t(P^r(\mathbf{X}_i) - m))\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{v=1}^V \left| \frac{\mathbf{1}_{\mathbf{X}_i=v} - (m' + \sum_{l=1}^{d'} V^{(l)t}(P(\mathbf{X}_i) - m')V^{(l,v)})}{\bar{P}_v} \right|^2 \end{aligned}$$

- Interpretation :
  - $m' = \bar{P}$
  - $\Phi(\mathbf{X}) = V^t(P^r\mathbf{X} - m)$  : coordinates in the restricted space.
  - $V^{(l)}$  appears as a probability profile.
- Complexity :  $O(n(1 + V^2) + d'V^2)$
- Link with Correspondence Analysis (CA)

### MFA Algorithm

- Redundant dummy coding of each categorical variable.
  - Renormalization of each block of dummy variable.
  - Classical PCA algorithm on the resulting variables
- 
- Interpretation as a reconstruction error with a rescaled/ $\chi^2$  metric.
  - Interpretation :
    - $\Phi(\mathbf{X}) = V^t(P^r(\mathbf{X}) - m)$  : coordinates in the restricted space.
    - $V^{(i)}$  : influence of each modality/variable in the  $i$ th new coordinates.
  - **Scaling** : This method is not invariant to a scaling of the continuous variables ! It is custom to normalize the variables (at least within groups) before applying PCA.

### PCA Model

- PCA : Linear model assumption

$$\mathbf{X} \simeq m + \sum_{l=1}^{d'} \mathbf{X}'_l V^{(l)}$$

- with
  - $V^{(l)}$  orthonormal
  - $\mathbf{X}'_l$  without constraints.
- Two directions of extension :
  - Other constraints on  $V$  (or the coordinates in the restricted space) : ICA, NMF, Dictionary approach
  - PCA on a non linear image of  $\mathbf{X}$  : kernel-PCA
- Much more complex algorithm !

# Dimension Reduction

## Non Linear PCA

### ICA (Independent Component Analysis)

- Linear model assumption

$$\mathbf{X} \simeq m + \sum_{l=1}^{d'} \mathbf{x}'_l V^{(l)}$$

- with
  - $V^{(l)}$  without constrains.
  - $\mathbf{x}'_l$  independent

### NMF (Non Negative Matrix Factorization)

- (Linear) Model assumption

$$\mathbf{X} \simeq m + \sum_{l=1}^{d'} \mathbf{x}'^{(l)} V^{(l)}$$

- with
  - $V^{(l)}$  non negative
  - $\mathbf{x}'_l$  non negative.

# Dimension Reduction

## Non Linear PCA

### Dictionary

- (Linear) Model assumption

$$\mathbf{X} \simeq m + \sum_{l=1}^{d'} \mathbf{x}'_l V^{(l)}$$

- with
  - $V^{(l)}$  without constrains
  - $\mathbf{X}'$  sparse (with a lot of 0)

### kernel PCA

- Linear model assumption

$$\Psi(\mathbf{X}) \simeq m + \sum_{l=1}^{d'} \mathbf{x}'_l V^{(l)}$$

- with
  - $V^{(l)}$  orthonormal
  - $\mathbf{X}'_l$  without constrains.

# Dimension Reduction

## Link with SVD

- Linear model assumption :

$$\mathbf{X} \simeq m + \sum_{l=1}^{d'} \mathbf{x}'_l V^{(l)}$$

- Vector rewriting

$$\mathbf{X}^t \simeq m^t + \mathbf{X}'^{t,t} V^t$$

- Matrix rewriting

$$\begin{pmatrix} \mathbf{X}_1^t - m^t \\ \vdots \\ \mathbf{X}_n^t - m^t \end{pmatrix} \simeq \begin{pmatrix} \mathbf{X}_1'^{t,t} \\ \vdots \\ \mathbf{X}_n'^{t,t} \end{pmatrix} V^t$$

- Low rank matrix factorization !
- Truncated SVD solution

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features



# Dimension Reduction

## Pairwise Distance

- Different point of view !
- Focus on pairwise distance  $d(\mathbf{X}_i, \mathbf{X}_j)$ .

### Distance Preservation

- Construct a map  $\Phi$  from the space  $\mathcal{X}$  into a space  $\mathcal{X}'$  of **smaller dimension** :

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}'$$

$$\mathbf{X} \mapsto \Phi(\mathbf{X}) = \mathbf{X}'$$

- such that

$$d(\mathbf{X}_i, \mathbf{X}_j) \sim d'(\mathbf{X}'_i, \mathbf{X}'_j)$$

- Most natural criterion :

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |d(\mathbf{X}_i, \mathbf{X}_j) - d'(\mathbf{X}'_i, \mathbf{X}'_j)|^2$$

- $\Phi$  often defined only on  $\mathbf{D}...$

# Dimension Reduction

## Random Projection

### Random Projection Heuristic

- Draw at random  $d'$  unit vector (direction)  $U_i$ .
- Use  $\mathbf{X}' = U^t(\mathbf{X} - m)$  with  $m = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$
- **Property** : If  $\mathbf{X}$  lives in a space of dimension  $d''$ , then, as soon as,  $d' \sim d'' \log(d'')$ ,

$$\|\mathbf{X}_i - \mathbf{X}_j\|^2 \sim \frac{d}{d'} \|\mathbf{X}'_i - \mathbf{X}'_j\|^2$$

- Do not really use the data !

# Dimension Reduction

## Locally Linear Embedding

### LLE Heuristic

- For each point  $\mathbf{X}_i$ , define a neighborhood  $\mathcal{N}_i$  (either by a distance or a number of points).
- Compute some weights  $W_{i,j}$  such that

$$W_{i,j} = 0 \quad \text{if } \mathbf{X}_j \notin \mathcal{N}_i$$
$$\mathbf{x}_i \sim \sum_j W_{i,j} \mathbf{x}_j$$

- Find some  $\mathbf{X}'_i$  in a space  $\mathcal{X}'$  of **smaller dimension** such that

$$\mathbf{x}'_i \sim \sum_j W_{i,j} \mathbf{x}'_j$$

- LLE : use a least square metric for the fits.

### MDS Heuristic

- If  $d(x, y) = \|x - y\|^2$ , one can compute a Gram matrix

$$(\mathbf{X}_i - m)^t (\mathbf{X}_j - m)$$

for  $m = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$

- Match the *scalar* products :

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |(\mathbf{X}_i - m)^t (\mathbf{X}_j - m) - \mathbf{X}_i'^t \mathbf{X}_j'|^2$$

- Linear method :  $\mathbf{X}' = U^t (\mathbf{X} - m)$  with  $U$  orthonormal

- **Beware** :  $\mathbf{X}$  is unknown !

# Dimension Reduction

## MultiDimensional Scaling

- Resulting criterion : minimization in  $U^t(\mathbf{X}_i - m)$  of

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |(\mathbf{X}_i - m)^t(\mathbf{X}_j - m) - (\mathbf{X}_i - m)^t U U^t (\mathbf{X}_j - m)|^2$$

without knowing explicitly  $\mathbf{X}...$

- Explicit solution obtained through the eigendecomposition of the know Gram matrix  $(\mathbf{X}_i - m)^t(\mathbf{X}_j - m)$  by keeping only the  $d'$  largest eigenvalues.
- In this case, MDS yields the same result than the PCA (but with different inputs, distance between observation vs correlations) !

# Dimension Reduction

## MultiDimensional Scaling

- **Explanation** : Same SVD problem up to a transposition :

- MDS

$$\overline{\mathbf{X}}_{(n)}^t \overline{\mathbf{X}}_{(n)} \sim \overline{\mathbf{X}}_{(n)}^t U U^t \overline{\mathbf{X}}_{(n)}$$

- PCA

$$\overline{\mathbf{X}}_{(n)} \overline{\mathbf{X}}_{(n)}^t \sim U^t \overline{\mathbf{X}}_{(n)} \overline{\mathbf{X}}_{(n)}^t U$$

- Complexity : ACP  $O(d' d^2)$  vs MDS  $O(d' n^2)$ ...

### MDS

- Apply this algorithm even if  $d(x, y) \neq \|x - y\|^2$  !
- **True distance minimization** : Simple gradient descent can be used (can be stuck in local minima).

- MDS : equivalent to PCA (but more expensive) if  $d(x, y) = \|x - y\|^2$  !
- ISOMAP : use a *localized* distance instead to limit the influence of very far point.

### ISOMAP

- For each point  $\mathbf{X}_i$ , define a neighborhood  $\mathcal{N}_i$  (either by a distance or a number of points) and let

$$d(\mathbf{X}_i, \mathbf{X}_j) = \begin{cases} 0 & \text{if } \mathbf{X}_j \notin \mathcal{N}_i \\ \|\mathbf{X}_i - \mathbf{X}_j\|^2 & \text{otherwise} \end{cases}$$

- Use the MDS algorithm with this modified distance

# Dimension Reduction

## Graph based

### Graph heuristic

- Construct a graph with weighted edges  $w_{i,j}$  measuring the *proximity* of  $\mathbf{X}_i$  and  $\mathbf{X}_j$  ( $w_{i,j}$  large if close and 0 if there is no information).
- Find the points  $\mathbf{X}'_i \in \mathbb{R}^{d'}$  minimizing

$$\frac{1}{n} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n w_{i,j} \|\mathbf{X}'_i - \mathbf{X}'_j\|^2$$

- Need of a constraint on the size of  $\mathbf{X}'_i$ ...
- Explicit solution through linear algebra :  $d'$  eigenvectors with smallest eigenvalues of the Laplacian of the graph  $D - W$ , where  $D$  is a diagonal matrix with  $D_{i,i} = \sum_j w_{i,j}$ .
- Variation on the definition of the Laplacian...



# Dimension Reduction

## t-Stochastic Neighbor Embedding

### SNE heuristic

- From  $\mathbf{X}_i \in \mathcal{X}$ , construct a set of conditional probability :

$$P_{j|i} = \frac{e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2}}{\sum_{k \neq i} e^{-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2}} \quad P_{i|i} = 0$$

- Find  $\mathbf{X}'_i$  in  $\mathbb{R}^{d'}$  such that the set of conditional probability :

$$Q_{j|i} = \frac{e^{-\|\mathbf{x}'_i - \mathbf{x}'_j\|^2 / 2\sigma_i^2}}{\sum_{k \neq i} e^{-\|\mathbf{x}'_i - \mathbf{x}'_k\|^2 / 2\sigma_i^2}} \quad Q_{i|i} = 0$$

is close from  $P$ .

- t-SNE** : use a Student-t term  $(1 + \|\mathbf{x}'_i - \mathbf{x}'_j\|^2)^{-1}$  for  $\mathbf{x}'_j$
- Minimize the Kullback-Leibler divergence  $(\sum_{i,j} P_{j|i} \log \frac{P_{j|i}}{Q_{j|i}})$  by a simple gradient descent (can be stuck in local minima).
- Parameters  $\sigma_i$  such that  $H(P_i) = -\sum_{j=1}^n P_{j|i} \log P_{j|i} = \text{cst.}$

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

### Feature Types

- Quantitative feature :
  - univariate, multivariate, *functional*
  - continuous, discrete
- Categorical feature :
  - Binary, nominal, ordinal
- List, relationship...

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

# Feature Design

## Single Quantitative Feature Transform

- **Idea** : a good prediction algorithm should be invariant to change of scale in the variables.

### Renormalization

- centering and standardization of a feature  $x$ ,
- mapping of  $x$  to  $[0, 1]$  if max and min are known
- Renormalization is *useless for purely linear methods...* but few methods are purely linear !
- **Idea** : Linear scale may not be the most natural one...

### Transformation

- log-scale instead of linear scale,
- Box-Cox transform,
- sigmoid, maxout, rectifier...
- application dependent transformation.

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - **Basis and Dictionary Learning**
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

- **Idea** : the behavior may not be *linear* in the feature.

#### Basis decomposition

- Replace a feature  $x$  by  $\varphi(x) = (\varphi_1(x), \dots, \varphi_B(x))$  where  $(\varphi_b)_{b=1}^B$  is a linearly independent family of functions.
  - In linear methods, use  $\langle \alpha, \varphi(x) \rangle$  instead of  $\alpha x$
  - Examples : polynomials, Fourier, wavelets...
- 
- Non parametric estimation technique...

- **Extension** : the behavior may not be *linear* in some features.

#### Basis decomposition

- Replace some features  $x = (x_1, \dots, x_n)$  by  $\varphi(x) = (\varphi_1(x), \dots, \varphi_B(x))$  where  $(\varphi_b)_{b=1}^B$  is a linearly independent family of functions.
  - In linear methods, use  $\langle \alpha, \varphi(x) \rangle$  instead of  $\langle \alpha, x \rangle$
  - Positive definite kernel associated to the scalar product  $K(x, y) = \langle \varphi(x), \varphi(y) \rangle$  plays an important role.
  - $\varphi$  is often defined implicitly by a choice of  $K$ ...
- 
- **Kernel trick** : SVM but also (generalized) linear model as seen for example with *quadratic* logistic modeling.



- **Idea** : some combinations of the features may be more interesting than any single one.

#### Decomposition idea

- Obtain  $(x_1, \dots, x_K)$  as a linear combination of some vectors  $\varphi_1, \dots, \varphi_{K'}$  of dimension  $K$

$$\begin{pmatrix} x_1 \\ \vdots \\ x_K \end{pmatrix} \sim \sum_{k'=1}^{K'} x'_{k'} \varphi_{k'}$$

- Use  $(x'_1, \dots, x'_{K'})$  instead of  $(x_1, \dots, x_K)$
- **Variation** : use  $(x'_1, \dots, x'_{K''})$  with  $K'' \leq K'$
- Most important part : choice of  $\varphi$  !

### PCA

- $\varphi_1, \dots, \varphi_K$  are the eigenvectors of the empirical covariance matrix of  $(x_1, \dots, x_K)$
- The  $\varphi_k$  are obtained by minimizing recursively :

$$\varphi_k = \operatorname{argmin}_{\varphi} \frac{1}{n} \sum_{i=1}^n \min_{x' \in \mathbf{R}^k} \left\| (x_{1,n} \dots, x_{K,n})^t - \left( \sum_{k'=1}^{k-1} x'_{k'} \varphi_{k'} + x'_k \varphi \right) \right\|^2$$

- The coefficients  $x'_{k'}$  are obtained for a new feature by minimizing

$$\left\| (x_1 \dots, x_K)^t - \left( \sum_{k'=1}^{K'} x'_{k'} \varphi_{k'} \right) \right\|^2$$

- Change of basis ! Useless for purely linear method if all the coefficients are kept !

### NMF (force positive weights!)

- Fix  $K'$  and find vectors  $\varphi_1, \dots, \varphi_{K'}$  such that the vectors of feature are well approximated by a linear sum with positive weights
- Non trivial minimization problem to find  $\varphi_k$  :

$(\varphi_1, \dots, \varphi_{K'})$

$$= \operatorname{argmin}_{(\varphi_1, \dots, \varphi_{K'})} \frac{1}{n} \sum_{i=1}^n \min_{x' \in (\mathbf{R}^+)^{K'}} \left\| (x_{1,n} \dots, x_{K,n})^t - \left( \sum_{k'=1}^{K'} x'_{k'} \varphi_{k'} \right) \right\|^2$$

- The coefficients  $x'_{k'}$  are obtained for a new feature by minimizing over  $x' \in (\mathbf{R}^+)^{K'}$

$$\left\| (x_1 \dots, x_K)^t - \left( \sum_{k'=1}^{K'} x'_{k'} \varphi_{k'} \right) \right\|^2$$

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - **Categorical Feature Encoding**
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

- **Idea** : transform a categorical feature  $x$  in *quantitative* ones  $c$  (encoding)...

### Encodings

- Binary encoding :  $x \in \{C_0, C_1\}$ 
  - binary code :  $c(x) = \mathbf{1}_{x=C_1}$
  - symmetrized version :  $c(x) = 2\mathbf{1}_{x=C_1} - 1$
- Nominal variable :  $x \in \{C_1, \dots, C_V\}$ 
  - binary code :  $c(x) = (\mathbf{1}_{x=C_2}, \dots, \mathbf{1}_{x=C_V})$
  - $V - 1$  quantitative variables.
- Ordinal variables :  $x \in \{C_1, \dots, C_V\}$ 
  - binary code :  $c(x) = (\mathbf{1}_{x \geq C_2}, \dots, \mathbf{1}_{x \geq C_V})$
  - $V - 1$  quantitative variables.
  - Feature selection makes more sense with those feature.

### General coding scheme

- Matrix representation :

$$c(x)^t = M \begin{pmatrix} \mathbf{1}_{x=C_0} \\ \vdots \\ \mathbf{1}_{x=C_V} \end{pmatrix}$$

with  $M \in \mathcal{M}_{V-1,V}$  chosen such that its column are linearly independent.

- Examples for  $V = 2$  :

$$M = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad M = \begin{pmatrix} -1 & 1 & -1 \\ -1 & -1 & 1 \end{pmatrix}$$

$$M = \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix}, \quad M = \begin{pmatrix} -1/3 & 2/3 & -1/3 \\ -1/3 & -1/3 & 2/3 \end{pmatrix}, \dots$$

- Choice of the coding matters for feature selection

- **Idea** : Capture interaction between features by encoding them jointly.

Full interaction between  $x = (x_1, \dots, x_K)$

- Matrix encoding :

$$c(x)^t = M \begin{pmatrix} \mathbf{1}_{x=(C_{1,1}, \dots, C_{1,K-1}, C_{1,K})} \\ \mathbf{1}_{x=(C_{1,1}, \dots, C_{1,K-1}, C_{V_K,K})} \\ \mathbf{1}_{x=(C_{1,1}, \dots, C_{2,K-1}, C_{1,K})} \\ \vdots \\ \mathbf{1}_{x=(C_{V_1,1}, \dots, C_{V_{K-1},K-1}, C_{V_K,K})} \end{pmatrix}$$

with  $M \in \mathcal{M}_{\prod_{k=1}^K v_k - 1, \prod_{k=1}^K v_k}$  such that the lines are linearly independent...

- **Example** of code for  $x = (x_1, x_2)$  with respectively 2 and 3 categories :

$$c(x)^t = \begin{pmatrix} 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{1}_{x=(1,1)} \\ \mathbf{1}_{x=(1,2)} \\ \mathbf{1}_{x=(1,3)} \\ \mathbf{1}_{x=(2,1)} \\ \mathbf{1}_{x=(2,2)} \\ \mathbf{1}_{x=(2,3)} \end{pmatrix}$$

- Interaction of order  $K$  leads to a very large number of categories !



- Modified coding scheme to impose a **hierarchical structure** :

$$c(x)^t = M \begin{pmatrix} \mathbf{1}_{x_1=C_{1,1}} \\ \vdots \\ \mathbf{1}_{x_K=C_{V_K,K}} \\ \mathbf{1}_{(x_1,x_2)=(C_{1,1},C_{1,2})} \\ \vdots \\ \mathbf{1}_{(x_{K-1},x_K)=(C_{V_{K-1},K-1},C_{V_K,K})} \\ \vdots \\ \mathbf{1}_{x=(C_{1,1},\dots,C_{1,K-1},C_{1,K})} \\ \vdots \\ \mathbf{1}_{x=(C_{V_1,1},\dots,C_{V_{K-1},K-1},C_{V_K,K})} \end{pmatrix}$$

with  $M \in \mathcal{M}_{\prod_{k=1}^K V_k - 1, \sum_{k=1}^K \sum_{i_1 < \dots < i_{k'} < \dots < i_k} \prod_{i_k'} V_{i_k'}}$  such that :

- its lines are linearly independent
- the first  $(\sum_{i_1 < \dots < i_{k'} < \dots < i'_K} \prod V_{i_k'}) - 1$  below are equal to zeros after the column  $\sum_{k=1}^K \sum_{i_1 < \dots < i_k} \prod V_{i_k}$

# Feature Design

## Categorical Feature Interaction

- **Example** with  $x = (x_1, x_2)$  with respectively 2 and 3 categories

$$c(x)^t = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{1}_{x_1=1} \\ \mathbf{1}_{x_1=2} \\ \mathbf{1}_{x_2=1} \\ \mathbf{1}_{x_2=2} \\ \mathbf{1}_{x_3=3} \\ \mathbf{1}_{x=(1,1)} \\ \mathbf{1}_{x=(1,2)} \\ \mathbf{1}_{x=(1,3)} \\ \mathbf{1}_{x=(2,1)} \\ \mathbf{1}_{x=(2,2)} \\ \mathbf{1}_{x=(2,3)} \end{pmatrix}$$

- Restriction possible to a given order of interaction by using only the first lines!
- Coding variant for ordered categories.
- **Extension** to interaction between a quantitative feature  $x_1$  and a categorical feature  $x_2$  through the mapping

$$(x_1, x_2) \rightarrow (x_1, x_1 c(x_2))$$

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

- **Idea** : Going from quantitative feature to binary features...
- Used in practice mainly for computing reason.

### Quantization/Binarization strategy

- Construct a finite size quantifier for the feature.
- Code the feature by a binary code of its quantized version.
- Construction of a quantifier :
  - Binning : use of a (regular) histogram with  $V$
  - $V$ -means : use the centers as *keywords* and assign a point to its nearest keyword.
  - Use  $V - 1$  quantiles
  - Vector coding...
- Extreme case :  $V = 2$  binarization !

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

- **Idea** : Reduce the number of values of a nominal feature with categories in a large set  $\mathcal{D}$  ?

### Hashing

- Construction of a *hashing* function  $H : \mathcal{D} \rightarrow \{1, \dots, V\}$  and use the hashed value instead of the original one.
  - The hashing function should be *as injective as possible...* in a probabilistic sense..
- 
- Design of such hashing function is a real art !
  - One can sometimes find hashing function such that  $d(H(C_k), H(C_{k'})) \sim d'(C_k, C_{k'})...$

- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

- **Idea** : Group features in a non linear way to reduce the number of features.

### Pooling

- **Quantitative features** : replace a subset or a list by
    - its min/max,
    - its range,
    - its average (linear...),
    - ...
  - **Nominal features** : replace a subset or a list by
    - its histogram
    - its most frequent values
    - ...
- 
- Discretization/Binarization can be use before and after...



- 1 Dimension Reduction
  - Reconstruction error
  - Distance preservation
- 2 Feature Design
  - Renormalization
  - Basis and Dictionary Learning
  - Categorical Feature Encoding
  - Quantization and Binarization
  - Hashing
  - Pooling
  - Application Specific Features

- Use of **application field specific** features :
  - Lots of know-how by experts of the field
  - (Almost) no price of using them if good feature selection

### Two examples

- Text and bag of words
  - Image and SIFT (Scale Invariant Feature Transform)
- 
- Deep Neural Networks seem not to require this part...

- How to transform a **text** into a vector of **categorical features** ?

### Bag of Words strategy

- Make a *list* of words,
- Compute a *weight* for each words.

### List building

- Make a list of all used words with their number of occurrence
- Gather the words in the list having the same stem (stemming)
- Hash the stem using a word specific hashing function (MurmurHash with 32bits for instance)
- Compute the histogram  $h_w(d)$

### Weight computation

- Compute the histogram  $h_w(d)$
- Apply the the tf-idf transform to the histogram

$$\text{tf-idf}_w(d) = \text{idf}_w \times \text{tf}_w(d)$$

with  $\text{idf}$  a corpus dependent weight

$$\text{idf}_w = \log \frac{n}{\sum_{i=1}^n \mathbf{1}_{h_w(d_i) \neq 0}}$$

and  $\text{tf}_w(d)$  the frequency within the document  $d$

$$\text{tf}_w(d) = \frac{h_w(d)}{\sum_w h_w(d)}.$$

- Most classical text preprocessing!

- How to transform an **image** into a vector of **features**?

### SIFT Strategy

- Compute a local descriptor based on local gradient.
- Agregate those measurements by histograms.

### SIFT Local Descriptor

- Compute a local scale and a principal orientation
- Compute the gradient at that scale, its norm and its angle with the principal direction
- Quantize this angle with 16 different values (binning)
- On each subwindow of  $4 \times 4$  subwindow grid oriented with the principal and of size the local scale, compute the sum of the gradient norm for each angle bin renormalized by the number of points in the subwindow.
- Use the  $4 \times 4 \times 16$  values as the local descriptor.

### SIFT based representation

- Compute the SIFT descriptor at each point of a regular grid
  - Assign each SIFT descriptor to the closest *keyword* obtained by *K*-means on the whole corpus (typically with  $K \sim 2000$ )
  - Compute a normalized histogram of the counts of those keywords
  - Use this 2000 values as the image descriptor
- 
- Used to be state-of-the-art, now lags behind DNN...