M1 - Apprentissage

Michèle Sebag – Benoit Barbot LRI – LSV

10 février 2014

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ = ● ● ●

Overview

Feature selection

Linear Change of Representation Principal Component Analysis Random projection Linear Semantic Analysis

Non-linear Change of Representation

Propositionalisation



Starting point: gathering the data

	AGE	SEX	BMI	BP		Serum	Mea	surer	nents		Response
Patient	x1	$\mathbf{x}2$	x3	x4	x5	$\mathbf{x}6$	$\mathbf{x7}$	x8	$\mathbf{x9}$	x10	У
1	59	2	32.1	101	157	93.2	38	4	4.9	87	151
2	48	1	21.6	87	183	103.2	70	3	3.9	69	75
3	72	2	30.5	93	156	93.6	41	4	4.7	85	141
4	24	1	25.3	84	198	131.4	40	5	4.9	89	206
5	50	1	23.0	101	192	125.4	52	4	4.3	80	135
6	23	1	22.6	89	139	64.8	61	2	4.2	68	97
:	1	:	:	÷	÷	÷	:	:	÷	÷	1
441	36	1	30.0	95	201	125.2	42	5	5.1	85	220
442	36	1	19.6	71	250	133.2	97	3	4.6	92	57

Find features

Before learning: describe the examples

- ► Too poor a description \Rightarrow nothing possible
- Too rich \Rightarrow feature pruning is required

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Why?

- ML is not a well-posed problem
- Adding useless features (the captain's age) can deteriorate the hypotheses

Feature Selection, Position of the problem

Context

- Too many features wrt number of examples
 - Remove
 Feature Selection
 - Build new features
 - Project on few features
- A particular case, first-order logic:

Feature Selection Feature Construction

Dimensionality Reduction

Propositionalisation

The hidden goal: select or build features ?

- Feature Construction : build good features
- ... makes learning easier...
- Best features: good hypotheses.

When learning boils down to feature selection Bio-informatics



- 30 000 genes
- few samples (expensive)
- goal: find genes relevant to diseases, resilience,

Position of the problem

Goals

- Selection: find a subset of features
- Ranking: order features by increasing relevance

Formalization

Given $\mathcal{A} = \{a_1, ... a_d\}$. Define

$$\begin{array}{lll} \mathcal{F}:\mathcal{P}(\mathcal{A}) & \mapsto \mathbb{R} \\ \mathcal{A}\subset \mathcal{A} & \mapsto \textit{Err}(\mathcal{A}) = \text{ min error of hypotheses built from } \mathcal{A} \end{array}$$

Find $Argmin(\mathcal{F})$

Challenge

- A combinatorial optimization problem (2^d)
- \bullet An unknown optimization function ${\cal F}$

Feature selection: the filter approach

Univariate approach

- Given current solution \mathcal{A}
- Add a_i to \mathcal{A}
- Examine whether removing a_j is relevant

Backtrack = less greedy, better optima, much more expensive

Feature selection: the wrapping approach

Multivariate approach

Measure the quality of a feature subset: estimate $\mathcal{F}(a_{i1}, ... a_{ik})$

CONS

Expensive: an estimate = solving an ML problem.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

PROS

Better optima

Feature selection: embedded approach

Principle (beforehand)

An ML criterion which favors hypotheses with few features For instance: find w, $h(x) = \langle w, x \rangle$, = argmin

$$\sum_{i} (h(x_i) - y_i)^2 + ||w||_1$$

data fitting

favor w with many null coordinates

Principle – a posteriori Given

$$h(x) = \langle w, x \rangle = \sum_{j=1}^d w_j x_j$$

If $|w_j|$ small, the *j*-th feature is unimportant Remove and restart the learning.

Filter approaches, 1

Notations

Training set:
$$\mathcal{E} = \{(x_i, y_i), i = 1..n, y_i \in \{-1, 1\}\}\$$

 $a(x_i) = value of feature a for example $(x_i)$$

Correlation

$$corr(a) = rac{\sum_i a(x_i).y_i}{\sqrt{\sum_i (a(x_i))^2 \times \sum_i y_i^2}} \propto \sum_i a(x_i).y_i = \langle a, y \rangle$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Limitations

Correlated features Non linear dependencies

Filter approaches, 2

Correlation and projection Repeat

Stoppiglia et al. 2003

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

• select a^* = feature most correlated to target

$$a^* = argmax\{\sum_i a(x_i)y_i, a \in \mathcal{A}\}$$

Project all other features on orthogonal space:

$$\begin{array}{lll} \forall b \in \mathcal{A} & b \rightarrow & b - \frac{\langle a^*, b \rangle}{\langle a^*, a^* \rangle} \ a^* \\ & b(x_i) \rightarrow & b(x_i) - \frac{\sum_j a^*(x_j)b(x_j)}{\sqrt{\sum_j a^*(x_j)^2}\sqrt{\sum_j b(x_j)^2}} a^*(x_i) \end{array}$$

Correlation and projection, cont

Project y on orthogonal space too

$$y \rightarrow y - \frac{\langle a^*, y \rangle}{\langle a^*, a^* \rangle} a^*$$
$$y_i \rightarrow y_i - -\frac{\sum_j a^*(x_j)y_j}{\sum_j a^*(x_j)^2} a^*(x_i)$$

- Until stopping criterion
 - Add random features $(r(x_i) = \pm 1)$ probe

When probes are selected, stop.

Limitations

does not work well when there are more than 6-7 relevant features (numerical noise).

Filter approaches, 3

Information gain

decision trees

$$p([a = v]) = Pr(y = 1 | a(x_i) = v)$$

$$Ql([a = v]) = -p([a = v]) \log p([a = v])$$

$$Ql(a) = \sum_{v} Pr(a(x_i) = v) Ql([a = v])$$



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Information gain, contd



Limitations

Myopic criterion Favors many-valued features Not well-suited to numerical features the XOR case

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Scores

in text mining, supervised learning Notations : c_i a class a_k a word or term

Criteria

- 1. Conditional probability
- 2. Mutual information
- 3. Chi-2
- 4. Relevance

$$P(c_i|a_k)$$

$$P(c_i, a_k)Log(\frac{P(c_i, a_k)}{P(c_i)P(a_k)})$$

$$\frac{(P(t,c)P(\neg t, \neg c) - P(t, \neg c)P(\neg t, c))^2}{P(t)P(\neg t)P(c)P(\neg c)}$$

$$\frac{P(t,c)+d}{P(\neg t, \neg c)+d}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Wrapper approaches

Principle: Generate and test

Given a list of candidate subsets $\mathcal{L} = \{A_1, .., A_p\}$

- Generate a new candidate A
- Compute $\mathcal{F}(A)$
 - learn h_A from $\mathcal{E}_{|A|}$
 - test h_A on a test set
- $\bullet \ \mathsf{Update} \ \mathcal{L}.$

Algorithms

- hill-climbing / multiple restart
- genetic algorithms
- genetic programming

$$=\hat{\mathcal{F}}(A)$$

Embedded approaches, 2

Principle

- Build a hypothesis
- Detect irrelevant features
- Prune them
- Iterate

Algorithm : SVM Recursive Feature Elimination Guyon et al. 03

- Linear SVM $\rightarrow h(x) = sign(\sum w_i.a_i(x) + b)$
- relevance(a_i) approx $|w_i|$
- Prune the bottom-k features
- Iterate.

Overview

Feature selection

Linear Change of Representation

Principal Component Analysis Random projection Linear Semantic Analysis

Non-linear Change of Representation

Propositionalisation

Dimensionality Reduction – Intuition

Degrees of freedom

- Image: 4096 pixels; but not independent
- ▶ Robotics: (# camera pixels + # infra-red) × time; but not independent

Goal

Find the (low-dimensional) structure of the data:

- Images
- Robotics
- Genes

Dimensionality Reduction

In high dimension

- Everybody lives in the corners of the space Volume of Sphere $V_n = \frac{2\pi r^2}{n} V_{n-2}$
- All points are far from each other

Approaches

- Linear dimensionality reduction
 - Principal Component Analysis
 - Random Projection
- Non-linear dimensionality reduction

Criteria

- Complexity/Size
- Prior knowledge



e.g., relevant distance

Linear Dimensionality Reduction

Training set

unsupervised

$$\mathcal{E} = \{(\mathbf{x}_k), \mathbf{x}_k \in \mathbb{R}^D, k = 1 \dots N\}$$

Projection from \mathbb{R}^D onto \mathbb{R}^d

$$\begin{split} \mathbf{x} \in \mathbb{R}^D \to & h(\mathbf{x}) \in \mathbb{R}^d, \ d << D \\ & h(\mathbf{x}) = A \mathbf{x} \end{split}$$

s.t. minimize $\sum_{k=1}^N ||\mathbf{x}_k - h(\mathbf{x}_k)||^2$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへぐ

Principal Component Analysis

Covariance matrix S Mean $\mu_i = \frac{1}{N} \sum_{k=1}^N X_i(\mathbf{x}_k)$

$$S_{ij} = rac{1}{N}\sum_{k=1}^{N}(X_i(\mathbf{x}_k) - \mu_i)(X_j(\mathbf{x}_k) - \mu_j)$$

 $\mathsf{symmetric} \Rightarrow \mathsf{can} \ \mathsf{be} \ \mathsf{diagonalized}$

$$S = U\Delta U' \quad \Delta = Diag(\lambda_1, \dots \lambda_D)$$

Thm: Optimal projection in dimension *d* projection on the first *d* eigenvectors of *S*

Let u_i the eigenvector associated to eigenvalue λ_i $\lambda_i > \lambda_{i+1}$

$$h: \mathbb{R}^D \mapsto \mathbb{R}^d, h(\mathbf{x}) = <\mathbf{x}, u_1 > u_1 + \ldots + <\mathbf{x}, u_d > u_d$$



Sketch of the proof

1. Maximize the variance of
$$h(\mathbf{x}) = A\mathbf{x}$$

$$\sum_{k} ||\mathbf{x}_{k} - h(\mathbf{x}_{k})||^{2} = \sum_{k} ||\mathbf{x}_{k}||^{2} - \sum_{k} ||h(\mathbf{x}_{k})||^{2}$$

Minimize
$$\sum_{k} ||\mathbf{x}_{k} - h(\mathbf{x}_{k})||^{2} \Rightarrow \text{Maximize } \sum_{k} ||h(\mathbf{x}_{k})||^{2}$$

$$Var(h(\mathbf{x})) = \frac{1}{N} \left(\sum_{k} ||h(\mathbf{x}_{k})||^{2} - ||\sum_{k} h(\mathbf{x}_{k})||^{2} \right)$$

As

$$||\sum_{k} h(\mathbf{x}_{k})||^{2} = ||A\sum_{k} \mathbf{x}_{k}||^{2} = N^{2}||A\mu||^{2}$$

(ロ)、(型)、(E)、(E)、 E) の(の)

where $\mu = (\mu_1, \dots, \mu_D)$. Assuming that \mathbf{x}_k are centered $(\mu_i = 0)$ gives the result.

Sketch of the proof, 2

2. Projection on eigenvectors u_i of SAssume $h(\mathbf{x}) = A\mathbf{x} = \sum_{i=1}^{d} \langle \mathbf{x}, v_i \rangle v_i$ and show $v_i = u_i$. $Var(AX) = (AX)(AX)' = A(XX')A' = ASA' = A(U\Delta U')A'$ Consider d = 1, $v_1 = \sum w_i u_i$ $\sum w_i^2 = 1$ $remind \lambda_i > \lambda_{i+1}$

$$Var(AX) = \sum \lambda_i w_i^2$$

maximized for $w_1 = 1, w_2 = \ldots = w_N = 0$ that is, $v_1 = u_i$.

Principal Component Analysis, Practicalities

Data preparation

Mean centering the dataset

$$\mu_i = \frac{1}{N} \sum_{k=1}^N X_i(\mathbf{x}_k)$$

$$\sigma_i = \sqrt{\frac{1}{N} \sum_{k=1}^N X_i(\mathbf{x}_k)^2 - \mu_i^2}$$

$$z_k = (\frac{1}{\sigma_i} (X_i(\mathbf{x}_k) - \mu_i))_{i=1}^D$$

Matrix operations

Computing the covariance matrix

$$S_{ij} = \frac{1}{N} \sum_{k=1}^{N} X_i(z_k) X_j(z_k)$$

► Diagonalizing S = U'∆U might be not affordable... Complexity $\mathcal{O}(D^3)$

Random projection

Random matrix

define

$$egin{aligned} A: {
m I\!R}^D &\mapsto {
m I\!R}^d \quad A[d,D] \quad A_{i,j} \sim \mathcal{N}(0,1) \ & \ h({f x}) = rac{1}{\sqrt{d}} A{f x} \end{aligned}$$

Property: h preserves the norm in expectation

$$E[||h({\bf x})||^2] = ||{\bf x}||^2$$
 With high probability
$$1 - 2exp\{-(\varepsilon^2 - \varepsilon^3)\frac{d}{4}\}$$

$$|\mathbf{1} - \varepsilon)||\mathbf{x}||^2 \le ||\mathbf{h}(\mathbf{x})||^2 \le (1 + \varepsilon)||\mathbf{x}||^2$$

Random projection

Proof

$$h(\mathbf{x}) = \frac{1}{\sqrt{d}} A \mathbf{x}$$

$$E(||h(\mathbf{x})||^2) = \frac{1}{d} E \left[\sum_{i=1}^d \left(\sum_{j=1}^D A_{i,j} X_j(\mathbf{x}) \right)^2 \right]$$

$$= \frac{1}{d} \sum_{i=1}^d E \left[\left(\sum_{j=1}^D A_{i,j} X_j(\mathbf{x}) \right)^2 \right]$$

$$= \frac{1}{d} \sum_{i=1}^d \sum_{j=1}^D E[A_{i,j}^2] E[X_j(\mathbf{x})^2]$$

$$= \frac{1}{d} \sum_{i=1}^d \sum_{j=1}^D \frac{||\mathbf{x}||^2}{D}$$

$$= ||\mathbf{x}||^2$$

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ = ● ● ●

Random projection, 2

Johnson Lindenstrauss Lemma For $d > \frac{9 \ln N}{\varepsilon^2 - \varepsilon^3}$, with high probability $(1 - \varepsilon)||\mathbf{x}_i - \mathbf{x}_j||^2 \le ||h(\mathbf{x}_i) - h(\mathbf{x}_j)||^2 \le (1 + \varepsilon)||\mathbf{x}_i - \mathbf{x}_j||^2$

More:

http://www.cs.yale.edu/clique/resources/RandomProjectionMethod.pdf

Overview

Feature selection

Linear Change of Representation

Principal Component Analysis Random projection Linear Semantic Analysis

Non-linear Change of Representation

Propositionalisation

Latent Semantic Analysis

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

- 1. Motivation
- 2. Algorithm
- 3. Discussion

Example

- c1: <u>Human</u> machine <u>interface</u> for ABC <u>computer</u> applications
- c2: A <u>survey</u> of <u>user</u> opinion of <u>computer system</u> <u>response time</u>
- c3: The <u>EPS user interface</u> management <u>system</u>
- c4: System and <u>human system</u> engineering testing of <u>EPS</u>
- c5: Relation of <u>user</u> perceived <u>response time</u> to error measurement

- m1: The generation of random, binary, ordered <u>trees</u>
- m2: The intersection <u>graph</u> of paths in <u>trees</u>
- m3: <u>Graph minors</u> IV: Widths of <u>trees</u> and well-quasi-ordering
- m4: <u>Graph minors</u>: A <u>survey</u>

Example, cont

	c 1	c 2	c 3	c 4	c 5	m1	m 2	m3	m 4
human	1	0	0	1	0	0	0	0	0
interface	1	0	1	0	0	0	0	0	0
computer	1	1	0	0	0	0	0	0	0
user	0	1	1	0	1	0	0	0	0
system	0	1	1	2	0	0	0	0	0
response	0	1	0	0	1	0	0	0	0
time	0	1	0	0	1	0	0	0	0
EPS	0	0	1	1	0	0	0	0	0
survey	0	1	0	0	0	0	0	0	1
trees	0	0	0	0	0	1	1	1	0
graph	0	0	0	0	0	0	1	1	1
minors	0	0	0	0	0	0	0	1	1

LSA, 2

Motivations

- Context : bag of words
- Curse of dimensionality
- Synonymy / Polysemy

Goals

- Dimensionality reduction
- A good topology (distance, similarity)

Remark

- First solution: cosine similarity
- Why not ?

More

```
http://lsa.colorado.edu
```

LSA, 3

Input

Matrix X = words \times documents



Principle

1. Change of coordinates concepts

2. Dimensionality reduction

Difference with Principal Component Analysis

from words and documents to

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

$LSA \equiv Singular Value Decomposition$

Input

Matrix X = words \times documents

$$X = U' S V$$

with \bullet U: change of word basis $m \times r$ $r \times d$

- V: change of document basis
- S: diagonal matrix

Dimensionality reduction

- S Order by decreasing eigenvalue
- S' = S cancel out all eigenvalues but the first (300) ones.

$$X' = U'S'V$$

 $m \times d$

 $r \times r$

Intuition

$$X=\left(egin{array}{cccccc} m_1 & m_2 & m_3 & m_4\ d_1 & 0 & 1 & 1 & 1\ d_2 & 1 & 1 & 1 & 0 \end{array}
ight)$$

 m_1 and m_4 are not present in the same documents, but are together with same words; "hence" they are somewhat related'... After SVD + Reduction,

$$X = \begin{pmatrix} m_1 & m_2 & m_3 & m_4 \\ d_1 & \epsilon & 1 & 1 & 1 \\ d_2 & 1 & 1 & 1 & \epsilon \end{pmatrix}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ



Singular value Decomposition of the words by contexts matrix

0.22	-0.11	0.29	-0.41	-0.11	-0.34	0.52	-0.06	-0.41
0.20	-0.07	0.14	-0.55	0.28	0.50	-0.07	-0.01	-0.11
0.24	0.04	-0.16	-0.59	-0.11	-0.25	-0.30	0.06	0.49
0.40	0.06	-0.34	0.10	0.33	0.38	0.00	0.00	0.01
0.64	-0.17	0.36	0.33	-0.16	-0.21	-0.17	0.03	0.27
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.30	-0.14	0.33	0.19	0.11	0.27	0.03	-0.02	-0.17
0.21	0.27	-0.18	-0.03	-0.54	0.08	-0.47	-0.04	-0.58
0.01	0.49	0.23	0.03	0.59	-0.39	-0.29	0.25	-0.23
0.04	0.62	0.22	0.00	-0.07	0.11	0.16	-0.68	0.23
0.03	0.45	0.14	-0.01	-0.30	0.28	0.34	0.68	0.18



Singular value Decomposition of the words by contexts matrix

・ロト ・ 一下・ ・ モト ・ モト・

æ





Singular value Decomposition of the words by contexts matrix

0.20	0.61	0.46	0.54	0.28	0.00	0.01	0.02	0.08
-0.06	0.17	-0.13	-0.23	0.11	0.19	0.44	0.62	0.53
0.11	-0.50	0.21	0.57	-0.51	0.10	0.19	0.25	0.08
-0.95	-0.03	0.04	0.27	0.15	0.02	0.02	0.01	-0.03
0.05	-0.21	0.38	-0.21	0.33	0.39	0.35	0.15	-0.60
-0.08	-0.26	0.72	-0.37	0.03	-0.30	-0.21	0.00	0.36
0.18	-0.43	-0.24	0.26	0.67	-0.34	-0.15	0.25	0.04
-0.01	0.05	0.01	-0.02	-0.06	0.45	-0.76	0.45	-0.07
-0.06	0.24	0.02	-0.08	-0.26	-0.62	0.02	0.52	-0.45



Singular value Decomposition of the words by contexts matrix

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

3.34 2.54



Singular value Decomposition of the words by contexts matrix

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで



Singular value Decomposition of the words by contexts matrix

	c1	c2	c3	v4	c 5	m1	m2	т3	m4
human	0.16	0.40	0.38	0.47	0.18	-0.05	-0.12	-0.16	-0.09
interface	0.14	0.37	0.33	0.40	0.16	-0.03	-0.07	-0.10	-0.04
computer	0.15	0.51	0.36	0.41	0.24	0.02	0.06	0.09	0.12
user	0.26	0.84	0.61	0.70	0.39	0.03	0.08	0.12	0.19
system	0.45	1.23	1.05	1.27	0.56	-0.07	-0.15	-0.21	-0.05
response	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
time	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
EPS	0.22	0.55	0.51	0.63	0.24	-0.07	-0.14	-0.20	-0.11
survey	0.10	0.53	0.23	0.21	0.27	0.14	0.31	0.44	0.42
trees	-0.06	0.23	-0.14	-0.27	0.14	0.24	0.55	0.77	0.66
graph	-0.06	0.34	-0.15	-0.30	0.20	0.31	0.69	0.98	0.85
minors	-0.04	0.25	-0.10	-0.21	0.15	0.22	0.50	0.71	0.62
		c 1	c 2	c3 c.	4 c 5	m 1	m 2	m3	m4
huma	n	c 1	c 2 0	c3 c	4 c 5 0	m 1 0	m 2 0	m3	m4
huma	n ace	c 1 1	c 2 0 0	c3 c 0 1 1 0	4 c 5 0	0 0	m 2 0 0	m3 0	m4 0
human interf compu	n ace uter	c 1 1 1	c 2 0 0 1	c3 c 0 1 1 0 0 0	4 c 5 0 0	m 1 0 0 0	m 2 0 0 0	m3 0 0	m4 0 0 0
human interf compu	n ace uter	c 1 1 1 0	c 2 0 1 1	c 3 c 0 1 1 0 0 0 1 0 1 0	4 c 5 0 0 1	m 1 0 0 0	m 2 0 0 0 0	m3 0 0 0 0	m 4 0 0 0 0
huma Interf compu user syster	n ace uter n	c 1 1 1 0 0	c 2 0 1 1 1	c 3 c 0 1 1 0 0 0 1 0 1 0 1 2	4 c5 0 0 1 0	m 1 0 0 0 0 0	m 2 0 0 0 0 0	m3 0 0 0 0 0	m 4 0 0 0 0 0 0
human Interf compu user system respo	n ace uter n nse	c 1 1 1 0 0 0 0	c 2 0 1 1 1 1 1	c3 c 0 1 1 0 0 0 1 0 1 2 0 0	4 c5 0 0 1 0 1	m1 0 0 0 0 0 0 0	m 2 0 0 0 0 0 0 0	m3 0 0 0 0 0 0 0	m4 0 0 0 0 0 0
human Interf compu user syster respo time	n ace uter n nse	c 1 1 1 0 0 0 0 0	c 2 0 1 1 1 1 1	c3 c 0 1 1 0 0 1 1 0 1 2 0 0 0 0	4 c5 0 0 1 0 1	m1 0 0 0 0 0 0 0 0	m2 0 0 0 0 0 0 0 0 0	m3 0 0 0 0 0 0 0 0 0	m 4 0 0 0 0 0 0 0 0
human interf compu user system respo time EPS	n ace uter n nse	c 1 1 1 0 0 0 0 0	c 2 0 1 1 1 1 1 1 0	c3 c 0 1 1 0 0 0 1 0 1 2 0 0 1 2 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1	4 c5 0 0 1 0 1 1 0	m1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	m 2 0 0 0 0 0 0 0 0 0	m3 0 0 0 0 0 0 0 0 0	m 4 0 0 0 0 0 0 0 0 0
huma Interf compu user syster respo time E PS surve	n ace uter n nse	c 1 1 1 0 0 0 0 0 0 0 0	c2 0 1 1 1 1 1 0 0	c3 c 0 1 1 0 0 0 1 2 0 0 1 2 0 0 1 1 0 0 1 1 0 0 1 1 0 0 0 0 1 1 0 0	4 c 5 0 0 1 1 1 0 0	m1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	m2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	m3 0 0 0 0 0 0 0 0 0 0	m 4 0 0 0 0 0 0 0 0 0 0
huma Interf compu user syster respo time E P S surve trees	n ace uter m nse y	c 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0	c 2 0 1 1 1 1 1 0 1 0 0	$\begin{array}{ccc} \mathbf{c3} & \mathbf{c}^{*} \\ \hline 0 & 1 \\ 1 & 0 \\ 0 & 0 \\ 1 & 0 \\ 1 & 2 \\ 0 & 0 \\ 0 & 0 \\ 1 & 1 \\ 0 & 0 \\ 0 \\ 0 & 0 \\ 0 \\ 0 & 0 \\ $	4 c5 0 0 1 1 0 0 0 0	m1 0 0 0 0 0 0 0 0 0 0 0 0	m2 0 0 0 0 0 0 0 0 0 0 0 0 0	m3 0 0 0 0 0 0 0 0 0 0 0 0	m 4 0 0 0 0 0 0 0 0 0 0 0 0 0
human interf compu user syster respo time EPS surve trees graph	n ace uter m nse y	c 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0	c 2 0 1 1 1 1 1 0 0 0 0 0	$\begin{array}{cccc} \mathbf{c3} & \mathbf{c} \\ \hline 0 & 1 \\ 1 & 0 \\ 0 & 0 \\ 1 & 0 \\ 1 & 2 \\ 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 1 & 1 \\ 0 & 0 \\ $	4 c5 0 0 1 1 1 0 0 0 0 0	m1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	m2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	m3 0 0 0 0 0 0 0 0 0 0 1 1	m 4 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

Discussion

An application

Synonymy test



P. Turney



Number of Dimensions in LSA (log)

Setting the number of dimensions

Trial and error :-(

Remarks

Negation apparently does not matter More: Google hits

Some applications

- Educational Text Selection
- Essay Scoring
- Summary Scoring & Revision

Cross Language Retrieval

LSA – Principal Component Analysis

Similarities

- Input: matrix
- Diagonalizing
- Cancel all eigenvalues but the highest ones
- Projection on the corresponding eigenvectors

Differences

	ACP	LSA
Matrix	covariance attributs	words $ imes$ documents
d	2-3	100-300

Overview

Feature selection

Linear Change of Representation Principal Component Analysis Random projection Linear Semantic Analysis

Non-linear Change of Representation

Propositionalisation

Non-Linear Dimensionality Reduction



Conjecture

Examples live in a manifold of dimension $d \ll D$

Goal: consistent projection of the dataset onto \mathbb{R}^d Consistency:

- Preserve the structure of the data
- e.g. preserve the distances between points

Multi-Dimensional Scaling

Position of the problem

- Given $\{\mathbf{x}_1, \ldots, \mathbf{x}_N, \mathbf{x}_i \in \mathbb{R}^D\}$
- Given $sim(\mathbf{x}_i, \mathbf{x}_j) \in \mathbb{R}^+$
- Find projection Φ onto \mathbb{R}^d

$$\begin{array}{ll} x \in \mathbb{R}^D \to & \Phi(x) \in \mathbb{R}^d\\ sim(\mathbf{x}_i, \mathbf{x}_j) \sim & sim(\Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j)) \end{array}$$

Optimisation

Define X,
$$X_{i,j} = sim(\mathbf{x}_i, \mathbf{x}_j)$$
; X^{Φ} , $X_{i,j}^{\Phi} = sim(\Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j))$
Find Φ minimizing $||X - X'||$
Rq : Linear Φ = Principal Component Analysis
But linear MDS does not work: preserves all distances, while

only local distances are meaningful

Non-linear projections

Approaches

- Reconstruct global structures from local ones and find global projection
- Only consider local structures

Intuition: locally, points live in \mathbb{R}^d



LLE

Isomap

Tenenbaum, da Silva, Langford 2000 http://isomap.stanford.edu

Estimate $d(x_i, x_j)$

- ▶ Known if **x**_i and **x**_j are close
- Otherwise, compute the shortest path between x_i and x_j geodesic distance (dynamic programming)

Requisite

If data points sampled in a convex subset of \mathbb{R}^d , then geodesic distance \sim Euclidean distance on \mathbb{R}^d .

General case

- Given $d(\mathbf{x}_i, \mathbf{x}_j)$, estimate $< \mathbf{x}_i, \mathbf{x}_j >$
- Project points in \mathbb{R}^d

Isomap, 2



Locally Linear Embedding

Roweiss and Saul, 2000 http://www.cs.toronto.edu/~roweis/lle/

Principle

 Find local description for each point: depending on its neighbors



Local Linear Embedding, 2

Find neighbors

For each \mathbf{x}_i , find its nearest neighbors $\mathcal{N}(i)$

Parameter: number of neighbors

Change of representation

Goal Characterize **x**_i wrt its neighbors:

$$\mathbf{x}_i = \sum_{j \in \mathcal{N}(i)} w_{i,j} \mathbf{x}_j \quad ext{ with } \sum_{j \in \mathcal{N}(i)} w_{ij} = 1$$

Property: invariance by translation, rotation, homothety **How** Compute the local covariance matrix:

$$C_{j,k} = < x_j - x_i, x_k - x_i >$$

Find vector w_i s.t. $Cw_i = 1$

Local Linear Embedding, 3

Algorithm Local description: Matrix W such that

 $\sum_{j} w_{i,j} = 1$

$$W = argmin\{\sum_{i=1}^{N} ||\mathbf{x}_i - \sum_j w_{i,j}\mathbf{x}_j||^2\}$$

Projection: Find $\{z_1, \ldots, z_n\}$ in \mathbb{R}^d minimizing

$$\sum_{i=1}^{N} ||z_i - \sum_j w_{i,j} z_j||^2$$

Minimize ((I - W)Z)'((I - W)Z) = Z'(I - W)'(I - W)Z

Solutions: vectors z_i are eigenvectors of (I - W)'(I - W)

• Keeping the *d* eigenvectors with lowest eigenvalues > 0

Example, Texts



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

Example, Images



LLE

Overview

Feature selection

Linear Change of Representation Principal Component Analysis Random projection Linear Semantic Analysis

Non-linear Change of Representation

Propositionalisation

Relational domains





Relational learning

PROS

Inductive Logic Programming

Use domain knowledge

CONS

Covering test \equiv subgraph matching

Data Mining exponential complexity

Getting back to propositional representation: propositionalization

West - East trains



▲□▶ ▲圖▶ ▲厘▶ ▲厘▶ 厘 の��

Linus (ancestor)

```
Lavrac et al, 94
```

$$\begin{array}{lll} \textit{West}(a) \leftarrow & \textit{Engine}(a,b), \textit{first_wagon}(a,c), \textit{roof}(c), \textit{load}(c,\textit{square},3)...\\ \textit{West}(a') \leftarrow & \textit{Engine}(a',b'), \textit{first_wagon}(a',c'), \textit{load}(c',\textit{circle},1)... \end{array}$$

West	Engine(X)	First Wagon(X,Y)	Roof(Y)	$Load_1(Y)$	$Load_2$ (Y)
а	b	С	yes	square	3
a'	b'	c'	no	circle	1

Each column: a role predicate, where the predicate is determinate linked to former predicates (left columns) with a single instantiation in every example

Stochastic propositionalization

Kramer, 98 Construct random formulas \equiv boolean features SINUS - RDS

http://www.cs.bris.ac.uk/home/rawles/sinus http://labe.felk.cvut.cz/~zelezny/rsd

- Use modes (user-declared) modeb(2,hasCar(+train,-car))
- Thresholds on number of variables, depth of predicates...
- Pre-processing (feature selection)





DB Schema

Propositionalization

RELAGGS

Database aggregates

- average, min, max, of numerical attributes
- number of values of categorical attributes

Apprentissage par Renforcement Relationnel







Contexte variable

- Nombre de robots, position des robots
- Nombre de camions, lieu des secours

Besoin: Abstraire et Generaliser

Attributs

- Nombre d'amis/d'ennemis
- Distance du plus proche robot ami

Distance du plus proche ennemi