Master Recherche IAC Option 2 Robotique et agents autonomes

> Jamal Atif – Michèle Sebag LRI

> > Feb. 7th, 2014

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Overview

Feature selection

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

Non-linear Change of Representation

Reinforcement learning for Feature Selection

Starting point: gathering the data

| Patient | AGE x1 | SEX x2 | BMI x3 | BP x4 | | Response | | | | | |
|---------|-----------|-----------|-----------|----------|-----|---------------|---------------|----------|---------------|-----|-----|
| | | | | | 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

Reinforcement learning for Feature Selection

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

Reinforcement learning for Feature Selection

Latent Semantic Analysis

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

- 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�?



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�?

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 | m4 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 system 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 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 | m 2 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 \\ 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 compuser 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

Reinforcement learning for Feature Selection

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



900

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



◆□ > ◆□ > ◆臣 > ◆臣 > ○ ● ● ● ●

Example, Images



LLE

Overview

Feature selection

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

Non-linear Change of Representation

Reinforcement learning for Feature Selection

Feature Selection as a one-player game

Romaric Gaudel^{1,2,3} and Michèle Sebag^{1,2,3}

¹ Univ. Paris-Sud, LRI, UMR8623 ² CNRS ³ INRIA-Saclay

ICML, June 2010









900

Feature Selection

Optimization problem

 $\underset{F\subseteq\mathcal{F}}{\operatorname{argmin}}\operatorname{Err}\left(\mathcal{A}\left(F,D\right)\right)$

- \mathcal{F} : Set of features
- F: Feature subset
- D: Training data set
- \mathcal{A} : Machine Learning algorithm
- Err: Generalization error

- Feature Selection (FS)
 - Minimize the Generalization Error
 - Decrease the learning/use cost of models
 - Lead to more understandable models
- Bottlenecks
 - Combinatorial optimization problem: find $F \subseteq \mathcal{F}$
 - Unknown objective function: generalization error

4 3 5 4 3

Feature Selection: state of the art / drawbacks

- Filter approaches [1]
 - ★ No account for all feature interdependencies
- Wrapper approaches
 - Tackling combinatorial optimization [2,3,4]
 - Tractability vs. exhaustivity tradeoff
- Embedded approaches
 - Using the learned hypothesis [5,6]
 - Using a regularization term [7,8]
 - * Restricted to linear models [7] or linear combinations of kernels [8]

- [2] D. Margaritis Toward provably correct Feature Selection in arbitrary domains. NIPS'09
- [3] T. Zhang Adaptive forward-backward greedy algorithm for sparse learning with linear models. NIPS'08
- [4] M. Boullé Compression-based averaging of selective Naive Bayes classifiers. J. Mach. Learn. Res. 07

[5] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik Gene selection for cancer classification using Support Vector Machines. Mach. Learn. 2002

- [6] J. Rogers, and S. R. Gunn Identifying feature relevance using a Random Forest. SLSFS'05
- [7] R. Tibshirani Regression shrinkage and selection via the Lasso. Journal of the Royal Statistical Society 94
- [8] F. Bach Exploring large feature spaces with hierarchical Multiple Kernel Learning. NIPS'08 < 🚊 + < 🚊 + > 🚊 🔗 🔍

^[1] K. Kira, and L. A. Rendell A practical approach to feature selection. ML'92

The one-player game approach

- Goal
 - Find argmin $\operatorname{Err} (\mathcal{A}(F, D))$
- Exploration vs Exploitation tradeoff
 - Virtually explore the whole lattice
 - Gradually focus the search on most promising Fs
 - Use a frugal, unbiased assessment of F
- How? tractability vs. optimality tradeoff
 - Upper Confidence Tree (UCT) [1]
 - ★ UCT ⊂ Monte-Carlo Tree Search
 - UCT tackles tree-structured optimization problems



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06 C > < - >











A >

Feature Selection as a one-player game

A Markov Decision Process

Set of features \mathcal{F} Set of states $\mathcal{S} = 2^{\mathcal{F}}$ Initial state \varnothing Set of actions $A = \{ \text{add } f, f \in \mathcal{F} \}$ Final state any state Reward function $V : \mathcal{S} \rightarrow [0, 1]$

• Ideally :
$$V(F) = \operatorname{Err} (\mathcal{A}(F, D))$$

 In practice: Fast unbiased estimate of Err (A (F, D))



Optimal Policy

Policy $\pi : S \to A$ Final state following a policy F_{π} Optimal policy

$$\pi^{\star} = \operatorname*{argmin}_{\pi} \operatorname{\mathsf{Err}}\left(\mathcal{A}\left(\mathsf{F}_{\pi},\mathsf{D}
ight)
ight)$$

Bellman's optimality principle $\pi^*(F) = \operatorname{argmin} V^*(F \cup F)$

 $\pi^{\star}(F) = \underset{f \in \mathcal{F}}{\operatorname{argmin}} V^{\star}(F \cup \{f\})$

with

$$V^{\star}(F) = \begin{cases} \mathsf{Err}(\mathcal{A}(F)) & \text{if } \mathit{final}(F) \\ \min_{f \in \mathcal{F}} V^{\star}(F \cup \{f\}) & \text{otherwise} \end{cases}$$



A (10) A (10)

 π^{\star} intractable \Rightarrow approximation using UCT







R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

ICML, June 2010 8 / 25

T > 4

< 6 k

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06: 🗆 🕨 (🗇 🔖 (🗟 🕨 (🗟) 🛬 (🕤)

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06: 🗆 🕨 (🗇 🗟 👘 (🗟 🕨 (🗟 👘 (🗟))

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06: 🗆 🕨 (🗇 🕨 (🗟 🕨 (🗟 🕨 (🧟 👘 (

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - ★ Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06 🗆 🕨 < 🗇 🛛 < 🚊 🗸 🥱 🔍

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - ★ Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06r 🗆 🕨 🐗 🗇 🖌 🧟 🕨 💈 🦘 ۹ 🤇

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - ★ Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06 🗆 🕨 4 🗇 💆 4 😇 💆 🍕
- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06(🗆 🕨 (🗇 💆 (😇) (😇) (

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06(🗆 🕨 (🗇 💆 (😇) (😇) (

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06 🗆 🕨 🖉 🕨 🍕 🐑 💈 🕤 🔍

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06 🗆 🕨 (🗇) (😇) (😇) (😇) (

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06: 🗆 🕨 (🗇) (😇) ()

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06: ロト く合い くきゃ くきゃ きゃう く

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06: 🗆 🕨 (🗇 🔖 (🗟) 🖉 ()

- Gradually grow the search tree
- Building Blocks
 - Select next action (bandit-based phase)
 - Add a node (leaf of the search tree)
 - Select next action bis (random phase)
 - ★ Compute instant reward
 - Update information in visited nodes
- Returned solution
 - Path visited most often



^[1] L. Kocsis, and C. Szepesvári Bandit based Monte-Carlo planning. ECML'06: 🗆 🕨 < 🗇 🔖 < 🚊 🔖 🚊 🕓

Multi-Arm Bandit-based phase

• Upper Confidence Bound (UCB1-tuned) [1]

• Select
$$\underset{a \in A}{\operatorname{argmax}} \hat{\mu}_a + \sqrt{\frac{c_e \log(T)}{t_a} \min\left(\frac{1}{4}, \hat{\sigma}_a^2 + \sqrt{\frac{c_e \log(T)}{t_a}}\right)}$$

- ★ T: Total number of trials in current node
- ★ t_a: Number of trials for action a
- ★ $\hat{\mu}_a$: Empirical average reward for action *a*
- ★ $\hat{\sigma}_a^2$: Empirical variance of reward for action *a*



[1] P. Auer, N. Cesa-Bianchi, and P. Fischer Finite-time analysis of the Multiarmed Bandit Problem: ML'02 💿 🔅 🔊 🔍



Extend UCT for Feature Selection: FUSE



R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

ICML, June 2010 11 / 25

A b

FUSE: bandit-based phase

Dealing with many arms

- Bottleneck
 - A many-armed problem (hundreds of features)
 - \Rightarrow need to guide UCT
- Ingredient 1: controlling the number of arms
 - Continuous heuristics [1]
 - ★ Use a small exploration constant c_e
 - Discrete heuristics [2,3]: Progressive Widening
 - ★ Consider only $\lfloor T^b \rfloor$ actions





- [1] S. Gelly, and D. Silver Combining online and offline knowledge in UCT. ICML'07
- [2] R. Coulom Efficient selectivity and backup operators in Monte-Carlo tree search. Computer and Games 2006

[3] P. Rolet, M. Sebag, and O. Teytaud Boosting Active Learning to optimality: a tractable Monte-Carlo, Billiard-based algorithm. ECML'09

R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

ICML, June 2010 12 / 25

FUSE: bandit-based phase

Sharing information among nodes

- Ingredient 2: sharing information among nodes
 - Rapid Action Value Estimation (RAVE) [1]
 - ★ RAVE(f) = average reward when $f \in F$



^[1] S. Gelly, and D. Silver Combining online and offline knowledge in UCT. ICML'07 + A = +

FUSE: random phase

Dealing with an unknown horizon

- Unknown best size of the feature subset
- Random phase policy
 - With probability 1 − $q^{|F|}$ stop
 Else add a uniformly selected feature
 |F| = |F| + 1 Iterate



< 回 > < 三 > < 三 >

Generalization error estimate

- Requisite
 - fast (to be computed 10⁴ times)
 - unbiased
- Proposed reward
 - k-NN like
 - + AUC criterion *
- Complexity: $\tilde{O}(mnd)$
 - d Number of selected features
 - n Size of the training set
 - *m* Size of sub-sample ($m \ll n$)



* Mann Whitney Wilcoxon test: $V(F) = \frac{|\{((x,y),(x',y')) \in \mathcal{V}^2, N_{F,k}(x) < \mathcal{N}_{F,k}(x'), y < y'\}|}{|\{((x,y),(x',y')) \in \mathcal{V}^2, y < y'\}|}$

Generalization error estimate

- Requisite
 - fast (to be computed 10⁴ times)
 - unbiased
- Proposed reward
 - k-NN like
 - + AUC criterion *
- Complexity: $\tilde{O}(mnd)$
 - d Number of selected features
 - n Size of the training set
 - *m* Size of sub-sample ($m \ll n$)



* Mann Whitney Wilcoxon test: $V(F) = \frac{|\{((x,y),(x',y')) \in \mathcal{V}^2, N_{F,k}(x) < \mathcal{N}_{F,k}(x'), y < y'\}|}{|\{((x,y),(x',y')) \in \mathcal{V}^2, y < y'\}|}$

Generalization error estimate

- Requisite
 - fast (to be computed 10⁴ times)
 - unbiased
- Proposed reward
 - k-NN like
 - + AUC criterion *
- Complexity: $\tilde{O}(mnd)$
 - d Number of selected features
 - n Size of the training set
 - *m* Size of sub-sample ($m \ll n$)



* Mann Whitney Wilcoxon test: $V(F) = \frac{|\{((x,y),(x',y')) \in \mathcal{V}^2, N_{F,k}(x) < \mathcal{N}_{F,k}(x'), y < y'\}|}{|\{((x,y),(x',y')) \in \mathcal{V}^2, y < y'\}|}$

Generalization error estimate

- Requisite
 - fast (to be computed 10⁴ times)
 - unbiased
- Proposed reward
 - k-NN like
 - + AUC criterion *
- Complexity: $\tilde{O}(mnd)$
 - d Number of selected features
 - n Size of the training set
 - *m* Size of sub-sample ($m \ll n$)



* Mann Whitney Wilcoxon test: $V(F) = \frac{|\{((x,y),(x',y')) \in \mathcal{V}^2, N_{F,k}(x) < \mathcal{N}_{F,k}(x'), y < y'\}|}{|\{((x,y),(x',y')) \in \mathcal{V}^2, y < y'\}|}$

Generalization error estimate

- Requisite
 - fast (to be computed 10⁴ times)
 - unbiased
- Proposed reward
 - k-NN like
 - + AUC criterion *
- Complexity: Õ(mnd)
 - d Number of selected features
 - n Size of the training set
 - *m* Size of sub-sample ($m \ll n$)



* Mann Whitney Wilcoxon test: $V(F) = \frac{|\{((x,y),(x',y')) \in \mathcal{V}^2, N_{F,k}(x) < \mathcal{N}_{F,k}(x'), y < y'\}|}{|\{((x,y),(x',y')) \in \mathcal{V}^2, y < y'\}|}$

FUSE: update

- Explore a graph
 - \Rightarrow Several paths to the same node
- Update only current path



< 🗇 🕨

э

4 3 > 4 3

From UCT to Feature Selection



- End learner
 - Any Machine Learning algorithm
 - Support Vector Machine with Gaussian kernel in experiments

4 10 10 14









R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

ICML, June 2010 18 / 25

< 6 b

Experimental setting

- Questions
 - FUSE vs FUSE^R
 - Continuous vs discrete exploration heuristics
 - FS performance w.r.t. complexity of the target concept
 - Convergence speed
- Experiments on

| DATA SET | SAMPLES | FEATURES | PROPERTIES |
|-------------|---------|----------|--------------------|
| MADELON [1] | 2,600 | 500 | XOR-LIKE |
| ARCENE [1] | 200 | 10,000 | REDUNDANT FEATURES |
| COLON | 62 | 2,000 | "EASY" |

[1] Feature Selection Challenge. NIPS'03

< 口 > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Experimental setting

- Baselines
 - CFS (Constraint-based Feature Selection) [1]
 - Random Forest [2]
 - Lasso [3]
 - RAND^R: RAVE obtained by selecting 20 random features at each iteration
- Results averaged on 50 splits (10 × 5 fold cross-validation)
- End learner
 - Hyper-parameters optimized by 5 fold cross-validation

^[1] M. A. Hall Correlation-based Feature Selection for discrete and numeric class Machine Learning. ICML'00

^[2] J. Rogers, and S. R. Gunn Identifying feature relevance using a Random Forest. SLSFS'05

^[3] R. Tibshirani Regression shrinkage and selection via the Lasso. Journal of the Royal Statistical Society 94 🚊 🗠 🔍 🔾

Results on Madelon after 200,000 iterations



• Remark: FUSE^{*R*} = best of both worlds

- Removes redundancy (like CFS)
- Keeps conditionally relevant features (like Random Forest)

R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

Results on Arcene after 200,000 iterations



• Remark: FUSE^{*R*} = best of both worlds

- Removes redundancy (like CFS)
- Keeps conditionally relevant features (like Random Forest)

22/25

Results on Colon after 200,000 iterations



Remark

All equivalent

R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

ICML, June 2010 23 / 25

NIPS 2003 Feature Selection challenge

- Test error on the NIPS 2003 Feature Selection challenge
 - On an disjoint test set

| DATABASE | ALGORITHM | CHALLENGE | SUBMITTED | IRRELEVANT |
|----------|------------------------------|---------------------------------|-----------|------------|
| | | ERROR | FEATURES | FEATURES |
| MADELON | FSPP2 [1] | 6.22% (1 ^{<i>st</i>}) | 12 | 0 |
| | D-FUSE ^R | 6.50% (24 th) | 18 | 0 |
| | BAYES-NN-RED [2] | 7.20% (1 st) | 100 | 0 |
| ARCENE | D-FUSE ^R (ON ALL) | 8.42% (3 rd) | 500 | 34 |
| | D-FUSE ^R | 9.42% 500 (8 th) | 500 | 0 |

Remarks

- Selected features: accurate
- Promising results

R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

^[1] K. Q. Shen, C. J. Ong, X. P. Li, E. P. V. Wilder-Smith Feature selection via sensitivity analysis of SVM probabilistic outputs. Mach. Learn. 2008

^[2] R. M. Neal, and J. Zhang Chap. High Dimensional Classification with Bayesian Neural Networks and Dirichlet Diffusion Trees. Feature extraction, foundations and applications, Springer 2006

Conclusion and Perspectives

- Contributions
 - Formalization of Feature Selection as a Markov Decision Process
 - Efficient approximation of the optimal policy (based on UCT)
 - ⇒ Any-time algorithm
 - Experimental results
 - State of the art
 - High computational cost (45 minutes on Madelon)
- Perspectives
 - Other end learners
 - Extend to Feature construction
 - ★ Inspired by [1]

R. Gaudel & M. Sebag (LRI)

^[1] F. de Mesmay, A. Rimmel, Y. Voronenko, and M. Püschel Bandit-based optimization on graphs with application to library performance tuning. ICML'09

Feature Selection as a one-player game

Romaric Gaudel^{1,2,3} and Michèle Sebag^{1,2,3}

¹ Univ. Paris-Sud, LRI, UMR8623 ² CNRS ³ INRIA-Saclay

ICML, June 2010









Filter approaches for Feature Selection

- Score features
- Select the best ones

Pro

Cheap

Cons

Cannot tackle all inter-dependencies between features

Filter approaches

- ANOVA (Analysis of Variance)
- RELIEFF [1]

[1] K. Kira, and L. A. Rendell A practical approach to feature selection. ML'92 <
D >
A B >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >
A E >

Wrapper approaches for Feature Selection

Test feature subsets

Really tackle the combinatorial problem

Pro

Look for the best solution

Cons

- Computationally expensive
- Wrapper approaches
 - Look ahead [1]
 - Mix forward/backward search [2]
 - Mix global/local search [3]

- [2] T. Zhang Adaptive forward-backward greedy algorithm for sparse learning with linear models. NIPS'08
- [3] M. Boullé Compression-based averaging of selective Naive Bayes classifiers. J. Mach, Learn. Res. 07: No. 2 🔊 🗠

^[1] D. Margaritis Toward provably correct Feature Selection in arbitrary domains. NIPS'09

Embedded approaches for Feature Selection

- Exploit the learned hypothesis
- And/Or modify the learning criterion to induce sparsity

Pro

• Based on relevance of features in the learned model

Cons

- Limited to linear models [1] or a linear combination of kernels [2]
- Possibly misled by feature interdependencies
- Embedded approaches
 - Lasso [1]
 - Multiple Kernel Learning [2]
 - Gini score on Random Forest [3]

- [2] F. Bach Exploring large feature spaces with hierarchical Multiple Kernel Learning. NIPS'08
- [3] J. Rogers, and S. R. Gunn Identifying feature relevance using a Random Forest. SLSES'05 💈 K K 🚊 K K 🚊 🖉 🔍

R. Gaudel & M. Sebag (LRI)

Feature Selection as a one-player game

ICML, June 2010 29 / 25

^[1] R. Tibshirani Regression shrinkage and selection via the Lasso. Journal of the Royal Statistical Society 94

FUSE: bandit-based phase

Dealing with many arms

Bandit-B Search Tree

- Bottleneck
 - A many-armed problem (hundreds of features)
 - \Rightarrow need to guide UCT
- Ingredient 1: controlling the number of arms
 - Discrete heuristics [1,2]: Progressive Widening
 - ★ Consider only $\lfloor T^b \rfloor$ actions
 - Continuous heuristics [3]
 - ★ Use a small exploration constant c_e
- Ingredient 2: sharing information among nodes
 - Rapid Action Value Estimation (RAVE) [3]



- [2] R. Coulom Efficient selectivity and backup operators in Monte-Carlo tree search. Computer and Games 2006
- [3] S. Gelly, and D. Silver Combining online and offline knowledge in UCT. ICML'07 > < B > < > < > < > <

σ-RAVE

FUSE: bandit-based phase

Sharing information among nodes

Use of RAVE

- Discrete heuristics [1]
 - ★ When a new action allowed, add argmax RAVE(f)
- Continuous heuristics [2]
 - ★ Tradeoff UCB-RAVE

$$(1 - \alpha) \cdot \hat{\mu}_{F,f} + \alpha ((1 - \beta) \cdot \ell \text{-RAVE}(F, f) + \beta \cdot g \text{-RAVE}(f)) + \text{exploration term}$$

*
$$\alpha \searrow$$
 when $t_{F,f} \nearrow$

*
$$\beta \searrow$$
 when $\#\{f \in F_t, F \rightsquigarrow F_t\} \nearrow$

R. Gaudel & M. Sebag (LRI)

^[1] P. Rolet, M. Sebag, and O. Teytaud Boosting Active Learning to optimality: a tractable Monte-Carlo, Billiard-based algorithm. ECML'09

^[2] S. Gelly, and D. Silver Combining online and offline knowledge in UCT. ICML'07 K 🖉 K K 🚊 K K 🚊 K K

Feature stop

Dealing with an unknown horizon

- Any state can be final or not
 - ▶ Final(F) = "f_s ∈ F"
 - *f_s*: A virtual stopping feature
- RAVE(f_s)
 - g-RAVE $(f_s^{(d)}) = average \{V(F_t), |F_t| = d+1\}$
 - ★ V(F_t): Reward of Feature Subset F_t selected at iteration t
 - d: When RAVE(f_s) is used, d is set to the number of features in current state



< 🗇 ト

4 3 > 4 3

Sensitivity of FUSE to the Computational Effort Madelon



Remarks

- FUSE: not enough features
- FUSE^R: 10 times faster than RAND^R

ICML, June 2010 33 / 25
Experimental setting

• FUSE Hyperparameters

| | | | HOW TO RESTRICT EXPLORATION | | | |
|-----------|------|---------------------|-----------------------------|-----------------|---------------------|--|
| | | | DISCRETE | CONTINUOUS | | |
| | | | HEURISTICS | HEURISTICS | | |
| PARAMETER | k-NN | q | b | Ce | C, C/ | |
| VALUE | 5-NN | 1 – 10 ⁱ | 1/2 | 10 ⁱ | 10 ⁱ | |
| i | | $\{-1, -3, -5\}$ | | {-4, -2, 0, 2} | $\{-\infty, 2, 4\}$ | |

- Plot best results for FUSE^R
- Preliminary results
 - FUSE is limited to deal with deep search tree
 - ▶ FUSE coincides with the beginning of the FUSE^R curve

4 3 5 4 3

< 6 k

Best hyperparameters

| HEURISTICS | ANY | Disc. | CONTINUOUS | |
|--------------|----------------------|-------|--------------------|-------------------------|
| PARAMETER | q | b | Ce | C, C/ |
| TESTED VALUE | 1 – 10 ⁱ | 1/2 | 10 ⁷ | 10 ⁷ |
| i | $\{-1, -3, -5\}$ | | $\{-4, -2, 0, 2\}$ | $\{-\infty, 2, 4\}$ |
| | $1 - 10^{-1}$ | 1/2 | | |
| ARCENE | $1 - 10^{-1}$ | | 10 ⁻² | ANY |
| | $1 - 10^{-3}$ | | 10 ⁻⁴ | ALMOST ANY |
| MADELON | 1 - 10 ⁻³ | 1/2 | | |
| | $1 - 10^{-1}$ | | 10 ⁻² | $\{(10^2,0),(10^4,0)\}$ |
| COLON | 1 - 10 ⁻⁵ | 1/2 | | |
| | $1 - 10^{-5}$ | | ΑΝΥ | ALMOST ANY |

æ

イロト イヨト イヨト イヨト