

Thematic classification for
EURO/IFORS conference
using expert model

Arsentiy Kuzmin, Alexander Aduenko, and Vadim Strijov

Moscow Institute of Physics and Technology
Department of Control and Applied Mathematics

EURO 2016, Poznan
05.06.2016

Construct a decision support system to assist the program committee and stream organizers make the forthcoming conference program

The goal:

- to construct a thematic model of the conference

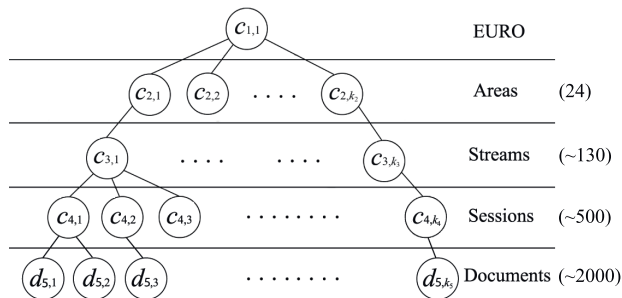
There given:

- historical expert thematic models of the previous conferences
- submitted abstracts for the forthcoming conference

The main idea:

- to calculate the similarity of a new abstract and each Stream,
- to show the most similar Streams to the Experts

EURO/IFORS conference hierarchical model



- 1 A group of experts is responsible for an Area,
- 2 participants submit their Abstracts to the collection,
- 3 the experts distribute the Abstracts over the Streams,
- 4 the Abstracts are organized into the Sessions.

Document-vector representation

Let $W = \{w_1, \dots, w_n\}$ be the terms dictionary of the collection.

There exists two different approaches to represent a word:

Vocabulary vector

$$C(w_j) \in \mathbb{R}^{|W|}, \quad |W| \approx 10^4$$

$$C(w_j) = e(j) = [\dots 0 \ 1 \ 0 \dots]^T$$

$$\|C(w_j) - C(w_i)\| = 2 \cdot \mathbb{I}(i \neq j)$$

Distance between words "logistic" and "transport" equals distance between "logistic" and "finance".

Document-vector representation:

$$\mathbf{x}_d = \sum_{w_j \in d} C(w_j)$$

word2vec, gloVec

$$C(w_j) \in \mathbb{R}^K, \quad |W| \approx 10^2$$

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

$$\|C(w_j) - C(w_i)\| \in \mathbb{R}$$

Words with similar contexts are closer, so

"logistic" \approx "transport", but
"logistic" \neq "finance".

Document-vector representation:
doc2vec, convolutional NN,
recursive NN.

The clustering quality function

Suppose F_0 is a mean intra-cluster similarity: $F_0 = \frac{1}{k_\ell} \sum_{i=1}^{k_\ell} S(c_{\ell,i}, c_{\ell,i})$,

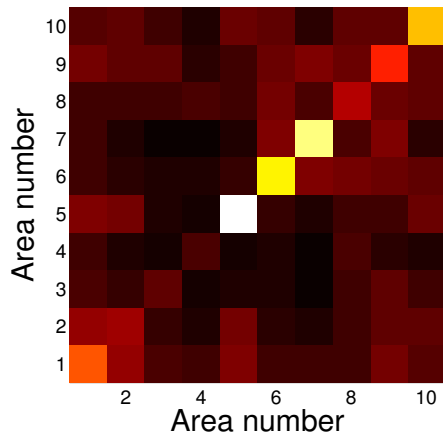
and F_1 is a mean inter-cluster similarity: $F_1 = \frac{2}{k_\ell(k_\ell - 1)} \sum_{i < j} S(c_{\ell,i}, c_{\ell,j})$

Clustering quality criterion

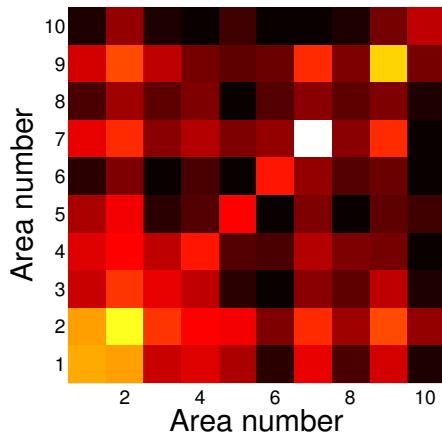
$$F = \frac{F_0}{F_1} \rightarrow \max$$

The expert model is the origin for the algorithmic thematic model.

Document representation comparison



Document vector - sum of it's word vectors $e(w_i)$.
F0/F1 = 1.98



Document vector - calculated with doc2vec approach.
F0/F1 = 1.31

Similarity function

Define the similarity function $s(\cdot, \cdot)$ between documents \mathbf{x}_i and \mathbf{x}_j as:

$$s(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_j}{\sqrt{\mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_i} \sqrt{\mathbf{x}_j^T \mathbf{\Lambda} \mathbf{x}_j}} = \mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_j, \text{ normalization: } \mathbf{x}_s \mapsto \frac{\mathbf{x}_s}{\sqrt{\mathbf{x}_s^T \mathbf{\Lambda} \mathbf{x}_s}},$$

where $\mathbf{\Lambda} = \text{diag}\{\lambda_{1,1}, \dots, \lambda_{n,n}\}$ is a term-importance matrix.

Define the similarity function $s(\cdot, \cdot)$ between the document \mathbf{x}_i and the cluster $c_{\ell,j}$ on the ℓ hierarchy level as:

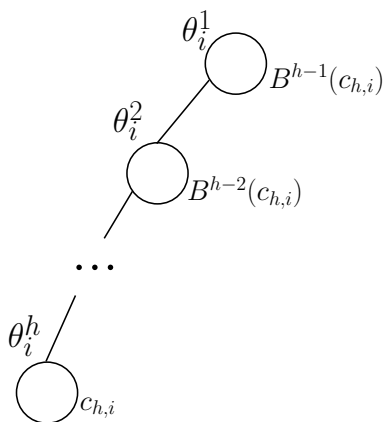
$$s(\mathbf{x}_i, c_{\ell,j}) = \mathbf{x}_i^T \mathbf{\Lambda} \bar{\mathbf{x}}_{\ell,i},$$

where $\bar{\mathbf{x}}_{\ell,i}$ is the mean vector of the cluster $c_{\ell,i}$.

Similarity between document and cluster of the h level

$$s(\mathbf{x}, c_{h,i}) = \sum_{j=0}^{h-1} \theta_i^{h-j} s(\mathbf{x}, B^j(c_{h,i})),$$

where θ_i^{h-j} is the weights parameter of the level $h-j$ for the cluster $c_{h,i}$, and B^j is the operator of the precedence that associate cluster $c_{h,i}$ with its predecessor on the level j .



Optimize weight parameters θ

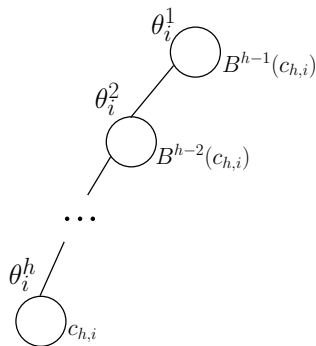
Find optimal parameters θ_i given expert hierarchy:

$$\theta_i^* = \arg \max_{\theta_i} \sum_{\mathbf{x} \in c_{h,i}} \sum_{j=0}^{h-1} \theta_i^{h-j} s(\mathbf{x}, B^j(c_{h,i})) + \mu \sum_{j=1}^h \left(\theta_i^j - \frac{1}{h} \right)^2,$$

$$\sum_{j=1}^h \theta_i^j = 1, \quad \theta_i^j \geq 0, \quad j \in \{1 \dots h\}.$$

Iterative procedure:

- find optimal $\mathbf{\Lambda}$ with fixed θ_i to maximize similarity $s(\mathbf{x}, B^j(c_{h,i}))$ (stay tuned);
- find optimal θ_i using fixed $\mathbf{\Lambda}$ and similarities;
- stop when $\Delta \theta_i$ and $\Delta \mathbf{\Lambda}$ for consecutive iterations are small.



Return for a new document all clusters sorted by similarity in descending order instead of the most similar one

Definition

Let $q \in S^{k_h}$ be the permutation of the level h clusters. The clusters in this permutation are sorted by the similarity to an object x in the descending order, k_h is the clusters quantity.

Example: $q = \{3, 1, \dots, 6\}$.

Definition

Let $R : \mathbb{R}^n \rightarrow S^{k_h}$ be the relevance operator. It maps the document $x \in \mathbb{R}^n$ to the permutation q .

Definition

Let $\text{pos}(q, j) : S^{k_h} \times \{1, 2, \dots, k_h\} \rightarrow \{1, 2, \dots, k_h\}$ be the position function. It returns the position of the given number in the permutation.

Example: $\text{pos}(q, 1) = 2$.

Quality criteria $Q(R)$ and $AUCH(R)$

$Q(R)$ quality criterion

Denote $Q(R)$ by the average position of the expert cluster $z_{j,h}$ in the permutation $R(\mathbf{x}_j)$:

$$Q(R) = \frac{1}{|D|} \sum_{j=1}^{|D|} \text{pos}(R(\mathbf{x}_j), z_{j,h}).$$

$AUCH(R)$ quality criteria

$AUCH(R) \in [0, 1]$ is the area under the top curve for a histogram $\#\{\text{pos}(R(\mathbf{x}_j), z_{j,h}) \leq i\}$, where $i \in [1, k_h]$.

$$AUCH(R) = \sum_{i=1}^{k_L} \frac{\#\{\text{pos}(R(\mathbf{x}_j), z_{j,h}) \leq i\}}{k_h |D|}.$$

Optimize Λ using the expert thematic model

$$\Lambda^* = \arg \min_{\Lambda} Q(R_{SIM}, \mathbf{X}).$$

Terms significance

Denote by $\mathbf{p}_{\ell,j}$ the vector of j -th components of cluster vectors $\bar{\mathbf{x}}_{\ell,i}$ of the level ℓ .

$$\mathbf{p}_{\ell,j} = [\bar{x}_{\ell,1,j}, \dots, \bar{x}_{\ell,k_\ell,j}]^T, \quad \mathbf{p}_{\ell,j} \mapsto \frac{\mathbf{p}_{\ell,j}}{\sum_{i=1}^{k_\ell} p_{\ell,i,j}}$$

The word entropy

Define the entropy $I_\ell(w_j)$ of the word w_j for hierarchy level ℓ as:

$$I_\ell(w_j) = \sum_{i=1}^{k_\ell} -p_{\ell,i,j} \log(p_{\ell,i,j}).$$

Term w_j significance according to its entropy

$$\lambda_j = 1 + \alpha_\ell \log(1 + I_\ell(w_j)).$$

Optimization using the collection with the expert model

$$\alpha_\ell^* = \arg \min_{\alpha_\ell} Q(R).$$

The purpose of the experiment

Construct a thematic model of the conference EURO 2013

The collection D^1 :

- EURO 2010, $|D| = 1663$, 26 Areas, 113 Streams.
- EURO 2012, $|D| = 1342$, 26 Areas, 141 Streams.

The collection D^2 :

- EURO 2013, $|D| = 2313$, 24 Areas, 137 Streams.

We matched the Areas and the Streams from all collections:

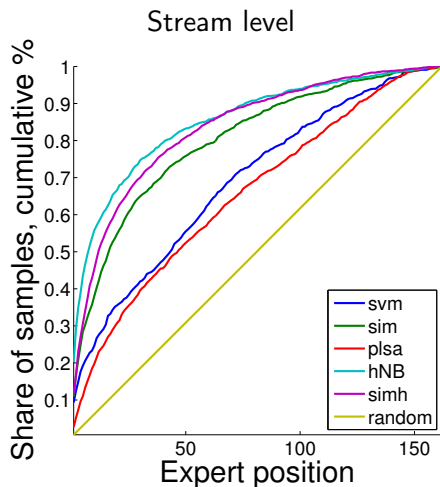
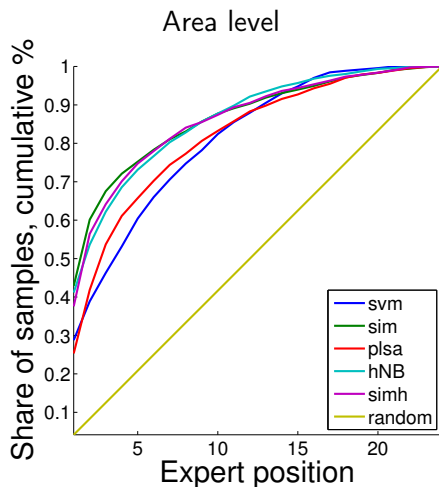
The unified structure has 24 Areas, 178 Streams.

15 out of 178 streams are present only in the EURO 2013.

Size of the expert dictionary:

- $|W| = 1675$ terms.

Quality comparison, EURO abstracts

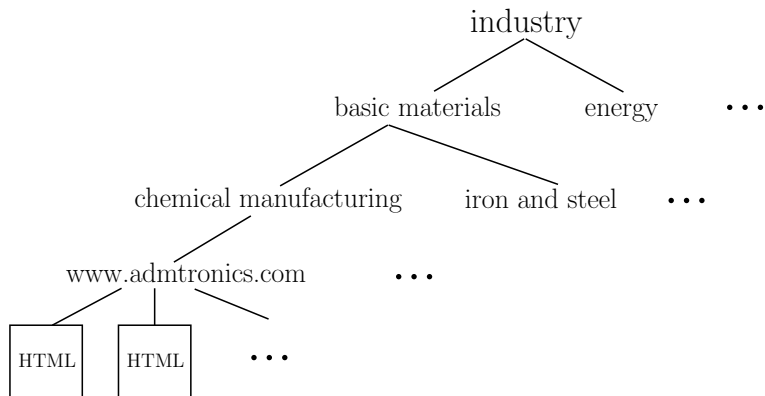


	svm	sim	plsa	hNB	simh
AUCH(R), Area level	0.80	0.85	0.81	0.84	0.85
AUCH(R), Stream level	0.69	0.78	0.65	0.84	0.83

Conference program validation for EURO/INFORMS abstract collection

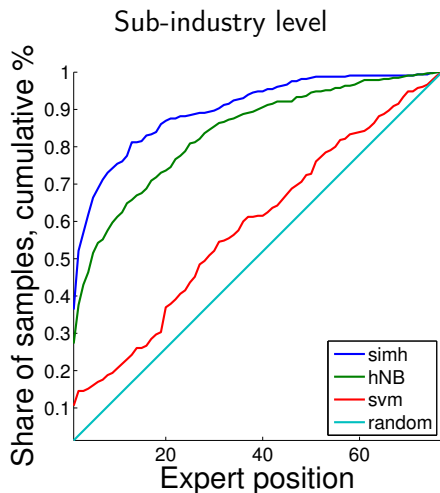
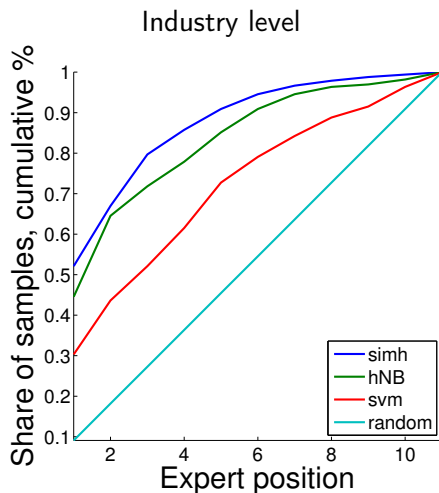
Paste title and abstract here	Search results (page 1 of 18)
<p>Title:</p> <input type="text" value="Hierarchical thematic model visualizing algorithm"/>	<p>Area: Emerging Applications of OR Stream: Models of Embodied Cognition <input type="button" value="Select"/></p>
<p>Abstract:</p> <p>The talk is devoted to the problem of the thematic hierarchical model construction. One must to construct a hierarchcal model of a scientific conference abstracts using machine learning clustering approach, to check the adequacy of the expert models and to visualize hierarchical differences between the algorithmic and expert models. An algorithms of hierarchcal thematic model constructing is developed. It uses the notion of terminology similarity to construct the model. The obtained model is visualized as the plane graph.</p>	<p>Area: OR in Health, Life Sciences & Sports Stream: Medical Decision Making <input type="button" value="Select"/></p>
<p><input type="button" value="Clear"/> <input type="button" value="Search"/></p>	<p>Area: Discrete Optimization, Geometry & Graphs Stream: Graphs and Networks <input type="button" value="Select"/></p>
	<p>Area: Data Science, Business Analytics, Data Mining Stream: Machine Learning and its Applications <input type="button" value="Select"/></p>
	<p>Area: Discrete Optimization, Geometry & Graphs Stream: Boolean and Pseudo-Boolean Optimization <input type="button" value="Select"/></p>
	<p>Area: Discrete Optimization, Geometry & Graphs Stream: Geometric Clustering <input type="button" value="Select"/></p>
	<p>Area: Multiple Criteria Decision Making and Optimization Stream: Preference Learning <input type="button" value="Select"/></p>
	<p>Area: Multiple Criteria Decision Making and Optimization Stream: Innovative Software Tools for MCDA <input type="button" value="Select"/></p>

The documents collection: industry sector web sites



- Collection: $|D| = 1036$ web sites classified by experts into 11 industries and 78 sub-industries.
- Dictionary: $|W| = 18000$ terms.
- Preprocessing: combine all web pages of a site, remove all html-tags

Quality comparison, Industry sector web sites



	svm	hNB	simh
AUCH(R), Industry level	0.70	0.80	0.83
AUCH(R), Sub-industry level	0.59	0.83	0.89

- The weighted hierarchical similarity function is proposed.
- The entropy-based method to calculate terms significance matrix Λ and optimize level weights θ is proposed.
- The relevance operator is proposed.
- The hierarchical similarity approach shows the same quality as the strongest base line with the expert dictionary.
- The hierarchical similarity approach shows significantly better quality with the automatically created dictionary.