

Multi-Task Data Mining toward Automating the KDD Process

Einoshin Suzuki

Kyushu University, Japan
suzuki@inf.kyushu-u.ac.jp

Simplified
version
e.g., movie

October 8, 2014 ICITEE 2014@Yogyakarta

Acknowledgement

Masamichi Shimura, Yves Kodratoff, Jan M. Zytkow, Shusaku Tsumoto, Yuu Yamada, Takeshi Watanabe, Hideto Yokoi, Katsuhiro Takabayashi, Masatoshi Jumi, Ning Zhong, Shin Ando, Daisuke Ikeda, Bin Tong, Thach Nguyen Huy, Hao Shao, Shigeru Takano, Asuki Kouno, Emi Matsumoto, Yutaka Deguchi, Daisuke Takayama, Vasile-Marian Scuturici, Jean-Marc Petit, Angdy Erna, Linli Yu, Kaikai Zhao, Wei Chen, Ryosuke Kondo and many others

for the more-than-a-dozen papers I am going to present in this talk.

Massive Data

Birth of data mining in the early 90's

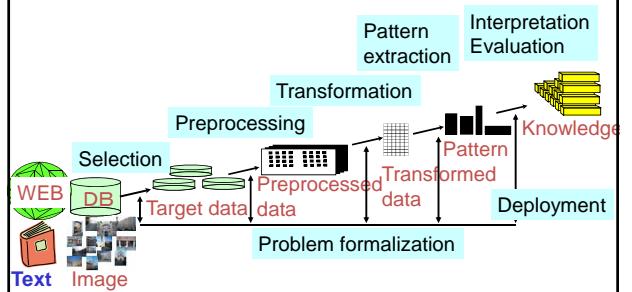


pixabay

KDD Process Model

KDD: Knowledge Discovery in Databases

Model modified from [Fayyad 1996]

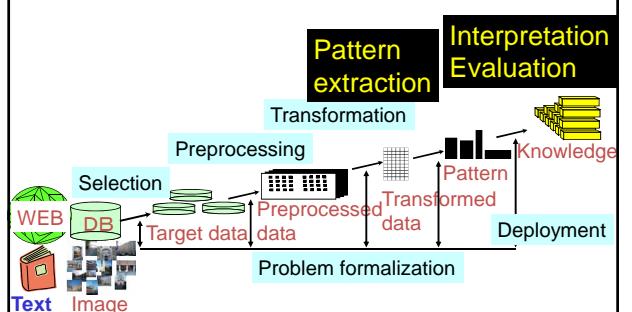


Toward Automating the KDD Process

*Analogy to mining:
impossible but a
must*

This talk: our 20-year efforts
Part 1: Exception rule discovery
Part 2: Medical data mining
Part 3: Multi-task learning
Part 4: Future direction

Part 1: Exception Rule Discovery



Rule Discovery

Fund type	5yrReturn	Diversit y	Beta	Stocks	Yield	Xpnse Rat%	Turnov er	Assets	Capgai n
Growth	belowS &P	Low	Under1	Over75 %	Under3 %	Low	Low	Large	10to20 %
Gth&In c	belowS &P	High	Under1	Over75 %	Under3 %	Low	High	Small	10to20 %
AggrGt h	belowS &P	High	Over1	Over75 %	Under3 %	Low	High	Small	Under 10%

FundType=Gth&Inc, Beta=under1 → Yield=over3%
(25 out of 30 examples)

Stocks=under75%, XpnsRat=low
→ 5yrReturn=belowS&P (15 out of 19 examples)

Too many rules are discovered
Rules are not related (i.e., fragmented)

Exception Rule Discovery [Suzuki, Shimura KDD 96, Suzuki KDD 97]

cured dead

antibiotics → cured
antibiotics, staphylococci → dead
(staphylococci ↗ dead)

$$\begin{cases} Y \rightarrow x \\ Y \wedge Z \rightarrow x' \\ Z \rightarrow x' \end{cases}$$
 where $\begin{cases} Y \equiv (a_1 = v_1) \wedge \dots \wedge (a_\mu = v_\mu) \\ Z \equiv (a_1' = v_1') \wedge \dots \wedge (a_v' = v_v') \\ x = (a'' = v_1'') \quad x' = (a'' = v_2'') \end{cases}$

a : attribute, v : value

Difficulties and Inventions

- High time complexity, $O(n^{2m+1})$ instead of $O(n^{m+1})$: branch-and-bound method, search pruning, bit-operation
- Many noisy patterns, weak exception or noise?: Analytical solution for the geometric problem

Pr(Y) ≥ θ₁^s,
Pr(x|Y) ≥ θ₁^f,
Pr(YZ) ≥ θ₂^s,
Pr(x'|YZ) ≥ θ₂^f,
Pr(x'|Z) ≤ θ₂^f

Confidence region + 5 constraints in a 7D-space

Analytical Solution

$$\begin{cases} 1 - \beta \sqrt{\frac{1 - \hat{p}(Y)}{n \hat{p}(Y)}} \hat{p}(Y) \geq \theta_1^s, \\ 1 - \beta \sqrt{\frac{1 - \hat{p}(Y, Z)}{n \hat{p}(Y, Z)}} \hat{p}(Y, Z) \geq \theta_2^s, \\ 1 - \beta \sqrt{\frac{\hat{p}(\bar{x}, Y)}{\hat{p}(x, Y) \left[(n + \beta^2) \hat{p}(Y) - \beta^2 \right]}} \hat{p}(x | Y) \geq \theta_1^f, \\ 1 - \beta \sqrt{\frac{\hat{p}(\bar{x}', Y, Z)}{\hat{p}(x', Y, Z) \left[(n + \beta^2) \hat{p}(Y, Z) - \beta^2 \right]}} \hat{p}(x' | Y, Z) \geq \theta_2^f, \\ 1 + \beta \sqrt{\frac{\hat{p}(\bar{x}', Z)}{\hat{p}(x', Z) \left[(n + \beta^2) \hat{p}(Z) - \beta^2 \right]}} \hat{p}(x' | Z) \leq \theta \end{cases}$$

β : value related to the dimension of the ellipsoid and δ

p : the ratio

Discovery from Meningitis Data [Suzuki, Tsumoto PAKDD 00]

High-quality data: aggregation at the patient level, collection by 2 physicians, one of whom is a data miner and the provider

Attribute in the conclusion (all)	#	Validness	Unexpectedness
CULT FIND	169	2.9	2.0
CT FIND	4	3.3	4.0
EEG FOCUS	36	3.3	3.0
	11	3.0	2.9

83 ≤ CSF PRO ≤ 121 → CULT FIND = F
+ FOCUL = + → CULT FIND = T
4/5 validness, novelty, unexpectedness, usefulness

Discovery from Blood Test Data [Suzuki, Tsumoto ws 00]

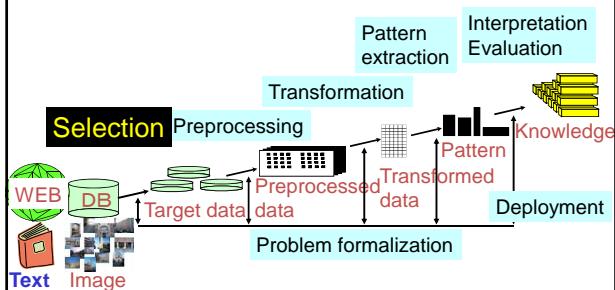
20919 examples (1 hospital, 1 year), 135 attributes, many missing data

Relatively "hard" data: example = medical test (no aggregation)

Sex = M, LCMs = Sensitive → PCG = Resistant
316/598 examples

Sex = M, LCMs = Sensitive, Ward = Outpatient → PCG = Sensitive
34/46 examples
Ward = Outpatient → PCG = Sensitive 1.6%

Part 2: Medical Data Mining



Chronic Hepatitis Data

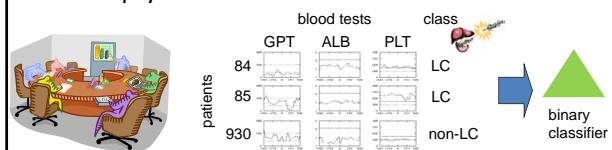
Collected for 20+ years in a university hospital

- Disease: virus inflame/harm liver cells
- Degree of fibrosis: index of progress F0 (no fiber) → F1 → F2 → F3 → F4 (liver cirrhosis: LC)
- Biopsy: to pick liver tissue for inspecting the degree of fibrosis (Invasive)
- Blood test: less invasive than biopsy
- Interferon: drug for killing virus. Has side effect

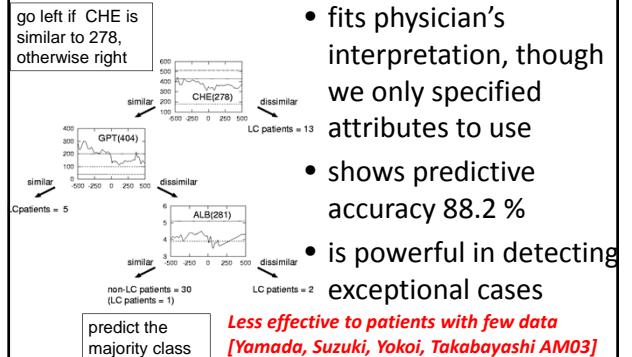


Selection/Formalization

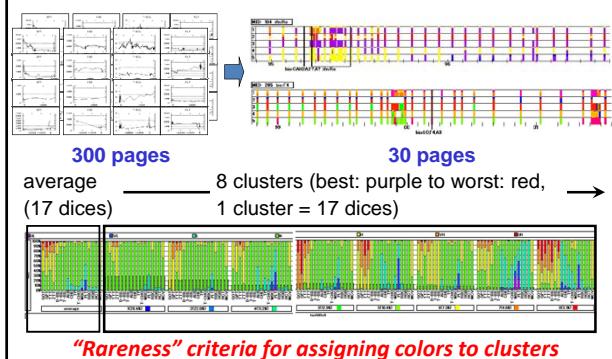
- Laborious KDD process: preprocessing and meetings (endless?)
- 1st Problem: liver cirrhosis prediction (binary classification) with 14 attributes.
- Selection: remove acute hepatitis and other diseases. Use each 500 days before and after the first biopsy. At least 10 tests.



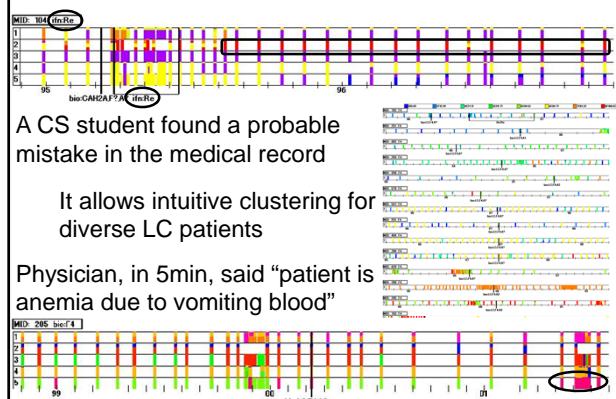
Time-Series Decision Tree [Yamada, Suzuki, Yokoi, Takabayashi, ICML 03]



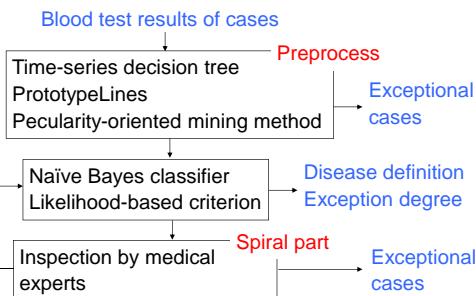
PrototypeLines [Suzuki, Watanabe, Yokoi, Takabayashi ICDM 03]



Data Cleaning, Clustering, Usability



Spiral Exception Discovery [Jumi, Suzuki, Ohsima, Zhong, Yokoi, Takabayashi ISMIS 2005]



Automatic Discovery [Jumi, Ohsima, Zhong, Yokoi, Takabayashi, Suzuki NGC j. 07]

No expert involved except the final evaluation. 14 + 21 blood tests. 140 LC patients + 106 non-LC patients (at least one value for each patient)

→ Discovered a hypothesis $\frac{\Pr(\text{AMY} = \text{high} \mid \text{LC})}{\Pr(\text{AMY} = \text{high} \mid \text{non-LC})} > 3$

The odd was below 3 in the original data so we paid no attention until our system removed exceptional cases

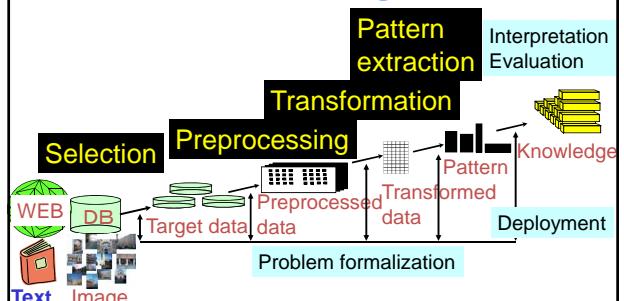
The hypothesis had been reported as the main result in Raffaele Pezzilli: "Serum Pancreatic Enzyme Concentrations in Chronic Viral Liver Diseases", Digestive Diseases and Sciences, Vol. 44, No. 2, pp. 350-355, 1999. (PMID: 10063922)

Anecdotes

- Medical criticism: Amy might be high due to operation and interferon → it is rare to do them before the first biopsy
- Astonishment: a pamphlet of MEXT says Amy is low → Dr. Takabayashi explained that the types of chronic hepatitis are different
- AI criticism: the discovery is just a good luck → Method shows an empirical evidence

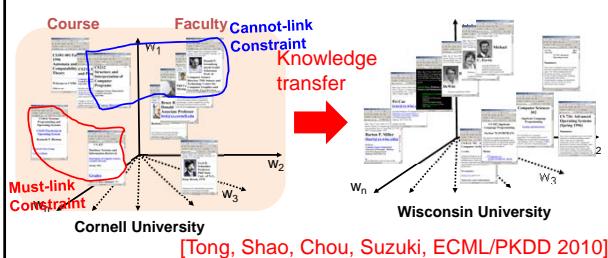
Difficulty of evaluating discovery methods
[Suzuki LNCS State-of-the-Art Survey 2002]

Part 3: Multitask Learning

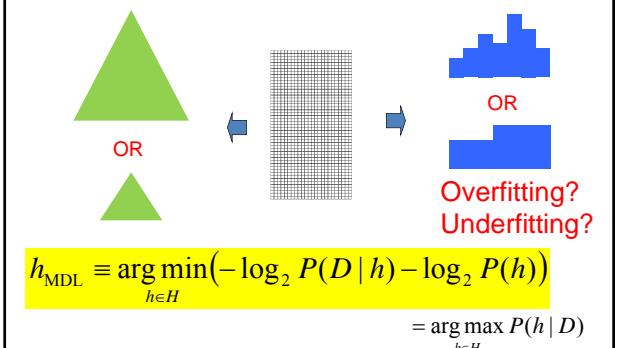


Multitask Learning

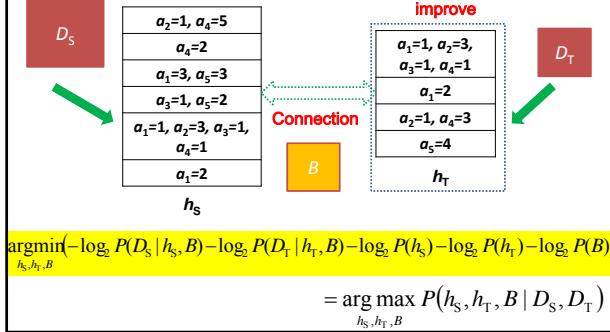
aims to improve the performance of learning algorithms by learning patterns for multiple tasks jointly (Weinberger's definition modified)



Overfitting and MDL [Rissanen 78-] for Single-Task Learning



Extended MDL for Multitask Learning [Shao, Suzuki SDM 11]



Experimental setting

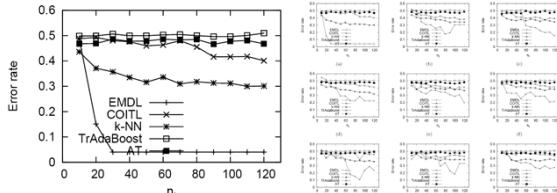
- Synthetic Datasets: Nine pairs (D_s, D_t) generated by the “target concepts” with 32 attributes. For example:

Index	μ	common	shortest	longest	avg.
(a)	3	1	1	3	1.33
(b)	3	1	2	3	2.33
(c)	3	2	2	3	2.33
(d)	4	1	2	4	3.24
(e)	5	1	3	4	3.60

- Real Datasets: Four data sets are used in the experiments in UCI repository. A pre-processing method [Y. Shi 09] is adopted on these data sets to split each data to the source and the target task.

Experimental Results (1)

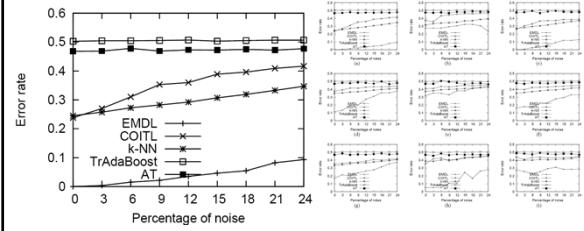
Results on synthetic datasets, for the different n_t of the target task, with 10% noise in both the source and the target tasks



Our method can find good features even n_t is small. By increasing n_t , EMDLP is able to achieve lower error rate.

Experimental Results (2)

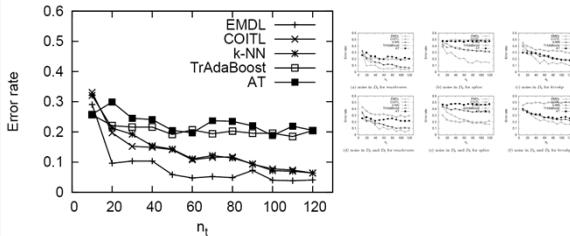
Results on synthetic datasets, for the different noise levels from 0% to 24% in both the source and the target tasks



Our EMDLP is robust to noise

Experimental Results (3)

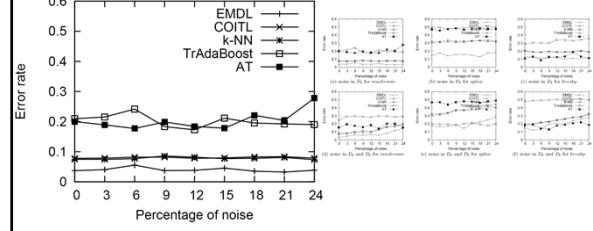
Results on real datasets, for the different n_t of the target task, with 10% noise in the source task or in both tasks.



Our method is able to achieve lower error rate with few labeled information available.

Experimental Results (4)

Results on real datasets, for the different noise levels from 0% to 24% in the source task or in both tasks.



Even under severe noise levels, our EMDLP is still able to obtain the contemplated h_t

Information Distance [Li et al. 92, 97]

Universal distance between string x and y

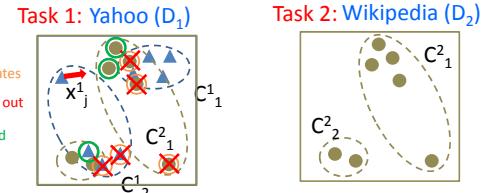
$$d(x, y) \equiv 1 - \frac{K(x) - K(x|y)}{K(xy)} \approx \frac{K(x|y) + K(y|x)}{K(xy)}$$

Kolmogorov complexity $K(x/y)$: length of the shortest program that returns y given x

K(xy): length of the shortest program that outputs xy
Not a computable function → Substitute K() by the compressed size

Can be used for single-task clustering

Extended Information Distance for Multi-Task Clustering [Nguyen Huy, Shao, Tong, Suzuki, ISMIS 11]



1. Clustering each domain separately
 2. Find candidate related instance set $S_{x_j^l} = \{C_1^1, C_1^2\}$
 3. Filter out unrelated instances $\Delta(x \mid y) = C(x) - C(x \mid y)$
 4. Calculate distance matrix $LDCDM(x_j^l, x_i^l) = \frac{C(x_j^l \mid S_{x_j^l}) + C(x_i^l \mid S_{x_j^l})}{C(x_j^l, x_i^l)}$

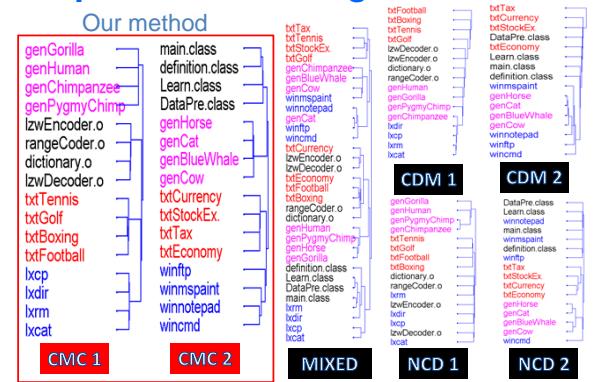
Go to Step 1

Experimental Setting

- **5 Data sets:**
 - Heterogeneous: 8 DNAs, 8 compiled, 8 text, 8 executive files.
 - DNA, Music, Language: Declaration of Human Rights [a] in 30 languages: 10 Europeans, 10 Americans and 10 Africans.
 - Document: 20 Newsgroup data [b].
 - **Comparison to 10 methods:** CDM [Keogh 04], NCD [Vitanyi 97], 4 Bernoulli methods [Steinbach 00], CLUTO [Karypis 98], Co-clustering [Dhillon 01], LSSMTC [Gu 09], MBC [Zhang 10].

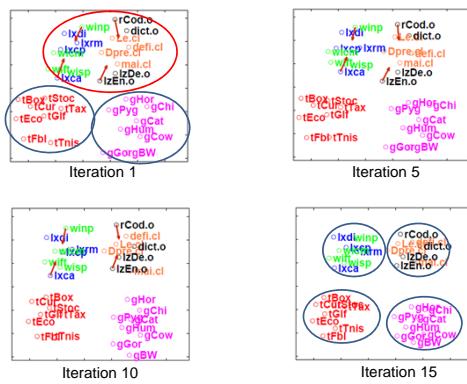
- [a] <http://www.un.org/en/documents/udhr/>
- [b] <http://people.csail.mit.edu/jrennie/20Newsgroups/>

Exp. Result: Heterogeneous Data

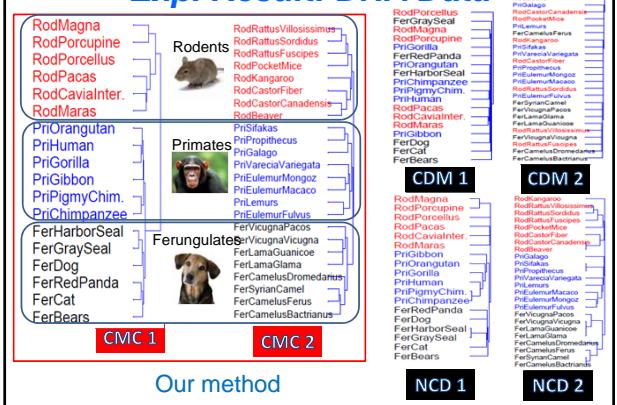


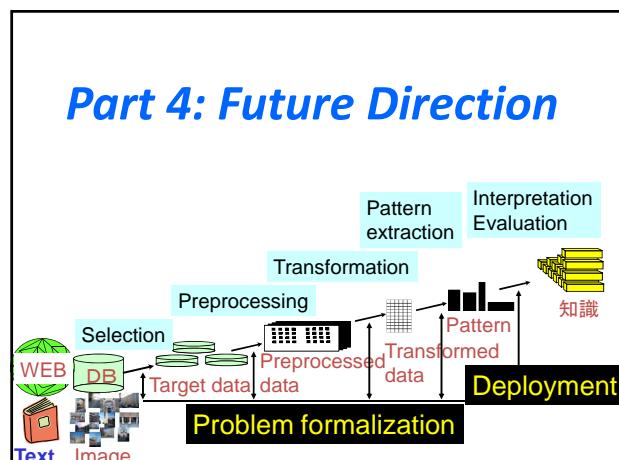
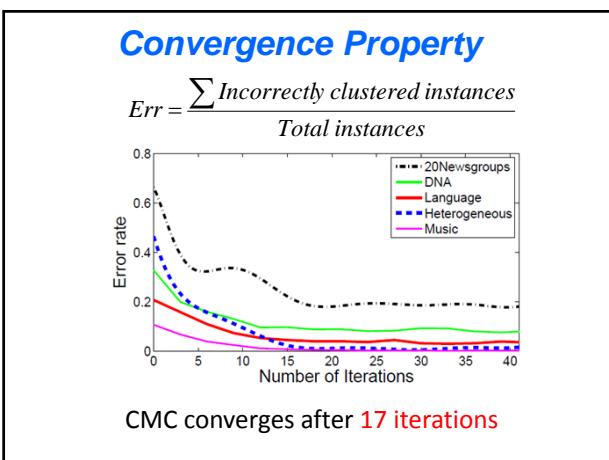
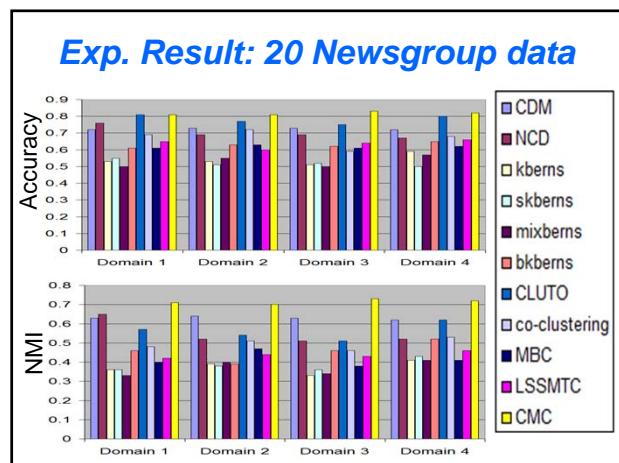
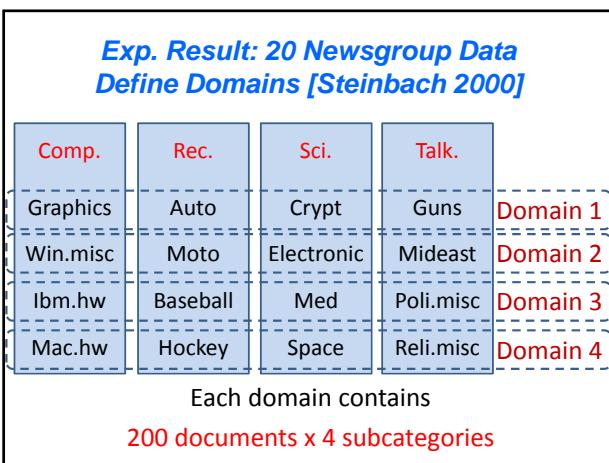
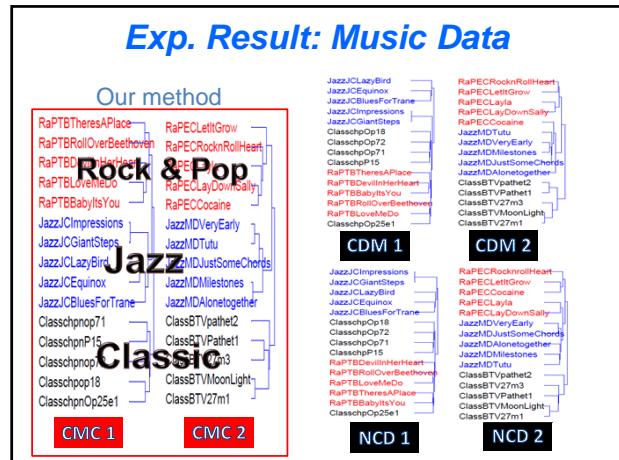
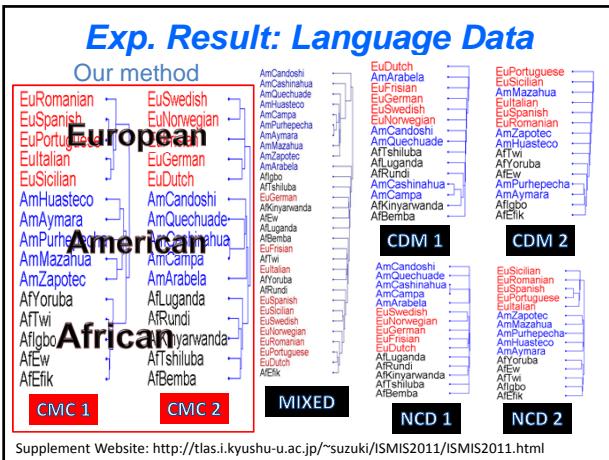
Supplement Website: <http://tlas.i.kyushu-u.ac.jp/~suzuki/ISMIS2011/ISMIS2011.html>

Exp. Result: Heterogeneous Data



Exp. Result: DNA Data





A photograph of the Discovery Robot, a small mobile robot with a camera mounted on top, positioned on a carpeted floor. The robot has a white body with black wheels and a red light on its front. To the left, the word 'movie' is written in yellow. To the right, there is text about the robot's capabilities and a reference to a paper by Zhang et al. A green bounding box highlights the robot.

Skeleton + color image, depth image, face image

Now we have
5 new Kinects!

<http://news.mynavi.jp/articles/2013/06/03/windows8Report/001.html>

Lifelong Learning [Ruvolo, Eaton ICML 13] and Application to 100-person Facial Expression Data [Tanoue, Suzuki dmc 14]

Sparse coding

$$e_T(\mathbf{L}) = \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{s}^{(t)}} \left\{ \frac{1}{n_t} \sum_{i=1}^{n_t} \mathcal{L} \left(f \left(\mathbf{x}_i^{(t)}; \mathbf{L}\mathbf{s}^{(t)} \right), y_i^{(t)} \right) + \mu \|\mathbf{s}^{(t)}\|_1 \right\} + \lambda \|\mathbf{L}\|_F^2$$

k	ELLA (%)	comparative method (%)
1	75	65
5	78	70
10	82	75
15	88	80

accuracy[%]

100-person data
[Erna, Yu, Zhao,
Chen, Suzuki AMT04]

- Outperforms single-task learning in 83/100 tasks
- Over 20% gain in 4 tasks

Effective knowledge transfer

Fall Risk Discovery [Deguchi, Suzuki ISMIS 14; Deguchi, Takayama, Takano, Scuturici, Petit, Suzuki PETRA 14 & Aml 14]

movie

Examples of discovered fall postures (clusters)

movie

Bilateral project with France (JSPS-CNRS)

