

Machine Learning Strategies for Big Data Utilization:

Assembling via Statistical Soft Label

Yasuo Matsuyama

Department of Computer Science and Engineering Waseda University Tokyo 169-8555, Japan

On the copyright regulations:

> Data in this presentation are made by this presenter ©.

> Concepts from others are referred, and are used only for academic purposes.



1

List of contents

- 1. Big data concept at large.
- 2. Concept 1:

Perfectly measured data are incomplete yet to be completed.

3. Concept 2:

Good models mismatch well, and are therefore dependable.

4. Methods:

From simple methods to machine learning methods.

5. Case studies:

Similar-video retrieval, numerical data and GUI, brain signals for authentication.

- 6. Bagging and crowd sourcing.
- 7. Concluding remarks.

Big data concepts

> True big data has an inaccessible size \Rightarrow No way to manage.

 3V for big data: Data's {Volume, Velocity, Variety} need to be manageable (D. Laney, META Group, 2001).
 4V: Veracity (by Villanova University).

>5V: Value (by HRBOSS blog).

>So many requirements for human power alone.

Machine learning strategies: Novel methods inspired by unprecedented problems.

>The fusion with crowd sourcing will also be mentioned.



Concept 2: Good models

- Essentially, all models are wrong, but some are useful.
 (G. E. P. Box; In G. E. P. Box and N. R. Draper, Empirical model building and response surfaces, p. 424, John Wiley & Sons, 1987).
- (2) Good models mismatch well, and are therefore dependable. (The presenter of this talk, Y. Matsuyama).

[Example]





space of *m*-th order AR processes

R. M. Gray, A. Buzo, A. H. Gray, Jr., and Y. Matsuyama, Distortion measures for speech processing, IEEE Trans. vol. ASSP-28, pp. 367-376, 1980.

5

Case 1: Missing is a weight vector (least squares)

Don't underrate simple classic methods. ... Data size is big.

[Example : E. Hoerl & R. W. Kennard, Technometrics, vol.12, pp 69 - 82, 1970.]

 $\mathbf{x} \cdots \mathbf{position}$ vector

 $\{\mathbf{x}_i\}_{i=1}^N \cdots$ observed data

linear model $\dots \quad y = f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^{\mathrm{T}} \mathbf{x}$ RBF kernel model $\dots \quad y = f(\mathbf{x}; \mathbf{w}) = \sum_{i=1}^{N} w_i \exp(-\gamma^2 ||\mathbf{x} - \mathbf{x}_i||^2)$

Least squares learning :

Here,

 $\mathbf{X} = \begin{cases} [\mathbf{x}_1 \cdots \mathbf{x}_N] \cdots \text{ for the linear model (matrix)} \\ x_{ij} = \exp(-\gamma^2 || \mathbf{x}_i - \mathbf{x}_j ||^2) \cdots \text{ for the RBF kernel model} \end{cases}$

(But, still sensitive to outliers \Rightarrow many improved methods to the rubustness.)



Case 2: Missing is a boundary set (supervised)



Fig. Support vector machine with a soft margin.



Fig. An SVM example using LIBSVM of C-C. Chang and C-J. Lin of NTU.



Prof. V. Vapnik on 2013-11-26 in Tokyo. Photograph allowed by himself. No usage other than academic one is allowed (© Y. Matsuyama).

7



- > Y. Matsuyama, R. M. Gray, IEEE Trans. vol. IIT-27, pp.31-40, 1981 (tree + batch VQ).
- > Y. Matsuyama, R. M. Gray, IEEE Trans. vol. COM-30, pp. 711-720, 1982 (tree + inverse filter VQ).
- T. Kohonen, Leaning vector quantization, In The Handbook of Brain Theory and Neural Networks., pp. 537–540. MIT Press, Cambridge, MA, 1995 (successive mode, or stochastic gradient descent).



Case 3-2: Missing is a boundary set (unsupervised)

Vector quantization was unified as a competitive learning algorithm for composite cost functions.

(Y. Matsuyama, Harmonic competition: A self-organizing multiple criteria optimization, IEEE Trans. NN, vol. 7, pp. 652-668, 1996.)

$$D = \sum_{n=0}^{N-1} D_n$$

$$D_n = \sum_{m=0}^{M-1} \left(f_n + \sum_{k=0}^{K-1} \lambda_{nk} g_{nk} \right) \left(\prod_{l=0}^{L-1} h_{nl} \right) Q(\mathbf{x}_n, \mathbf{w}_m)$$

$$f(\mathbf{x}_{n}, \mathbf{w}_{m}) \cdots \text{ distance between } \mathbf{x}_{n} \text{ and } \mathbf{w}_{m}$$

$$g_{nk} \left(\{\mathbf{x}_{i}\}_{i=0}^{N-1}, \{\mathbf{w}_{j}\}_{j=0}^{M-1} \right) \cdots \text{ constraint}$$

$$h_{nl} \left(\{\mathbf{x}_{i}\}_{i=0}^{N-1}, \{\mathbf{w}_{j}\}_{j=0}^{M-1} \right) \cdots \text{ penalty for label conflict}$$

$$Q(\mathbf{x}_{n}, \mathbf{w}_{m}) \cdots \text{ lif } \mathbf{w}_{m} \text{ is the winner for the input } \mathbf{x}_{n}; \text{ otherwise 0.}$$



9

Case 3-3: Missing is an exemplar set (unsupervised)

Affinity propagation (similar to VQ, but an unspecified number of exemplars is found).

(B. J. Frey and D. Dueck, Clustering by passing messages between data points, Science, vol. 315, no. 5814, pp. 972-976, 2007.



Fig.. Exemplars found by affinity propagation.

- \succ Similarity matrix \Rightarrow exemplars and clusters associated with them.
- > Convergence is not yet proved rigorously (but, often effective).
- > For time series, we need its sophistication.

FIEEE IET ICALIP 2014

Case 4-1: Missing information is more general (unsupervised/supervised)

> EM algorithm (expectation-maximization algorithm):

This is regarded as a champion algorithm which estimates missing information. (A. P. Dempster, N. M. Laird, D. B. Rubin, Maximum likelihood from Incomplete data via the EM Algorithm, J. Royal Statistical Society. Series B, vol. 39, pp. 1-38, 1977).



Case 4-2: α -EM algorithm (EM algorithm is yet a special case) α -EM algorithm's basic equation: $L_{Y}^{(\alpha)}(\psi \mid \varphi) = E_{p_{X|Y,\varphi}}[L_{X}^{(\alpha)}(\psi \mid \varphi)] + \frac{1-\alpha}{2} \left\{ \frac{p_{Y|\psi}(y \mid \psi)}{p_{Y|\varphi}(y \mid \varphi)} \right\}^{\frac{1+\alpha}{2}} D_{X|Y}^{(\alpha)}(\varphi \parallel \psi)$ (1) We want to keep this them to be nonnegative. (3) Keep this term nonnegative



Y. Matsuyama,

The alpha-EM algorithm: Surrogate likelihood maximization using alphalogarithmic information measures, IEEE Trans. On Information Theory, vol. 49, No. 3, pp. 692-706, March, 2003.







Fig. Relationship of alpha-HMM and log-HMM



- Y. Matsuyama, Hidden Markov model estimation based on the alpha-EM algorithm: Discrete and continuous alpha-HMMs, Proceedings of International Joint Conference on Neural Networks, pp.809-816, San Jose, CA, 2011.
- > A. Rahimi's correction to L. Rabiner's flaws on scaling was generalized too.

Case 4-5: ICA (missing are hidden components).



- Fig. Rapid ICA
- R. Yokote, Y. Matsuyama, Rapid algorithm for independent component analysis, J. Signal and Info. Proc., vol. 3, pp. 275-285, 2013.
- Fast ICA was by A Hyvärinen; Fast and robust fixed-point algorithms for independent component analysis, IEEE Trans. NN, vol. 10, No. 3, 1999, pp. 626- 634., 1999.



15

Case studies: Unstructured data management

- Tremendous amount of data is accumulated continuously.
- Mostly, data are unstructured.
- Some are manually labeled, however, labels are usually personal.
- Structurization by machine learning strategies is effective and fashionable.
- New target problems ⇒ new methods
 ⇒ new theory



(1) Similar video retrieval: Numerical labelling

- > Video retrieval: One of the most familiar target to be structured.
- > The size of each element is already big.
- Lengths are different.



Fig. Similar video retrieval.



	Р	↑ -6	[、] -1	_← -7	⊱13	[、] -11	_← -17	_← -23	_← -29	← - 35
	А	↑ -12	[↑] †-7	` -2	8- →	∿ ← -14	ົ-13	`-1 3	_← -19	_← -25
i 	Ν	↑ -18	-12	[↑] -8	ົ2	, −8	`-13	`-15	[►] -7	-13 _م
	D	↑ -24	`-16	↑ -14	- 6	` -3	ہے ۔9	`15	∱ -13	ہ -8
Ļ	А	↑ -30	1 -22	`-1 7	↑ -12	` -7	` -5	` -5	_← -11	` -13

Fig. Needleman-Wunsch algorithm in bioinformatics.

- Too simple: L-distance (Levenshtein distance) from another big data area is not applicable.
 - ✓ Only letter-wise.
 - ✓ Not context-aware.



- A more general distance measure ••• M-distance*.
- Each video can be numerically labelled using positions and expressions of exemplars.
- (*) Y. Matsuyama and M. Moriwaki

FIEEE IET ICALIP 2014

FIEEE IET ICALIP 2014



- > How do we express each as a vector.
- > You can find many effective ways.
- > We here show you an example by CSD.



Fig. A conceptual CSD (color structure descriptor).



(1-5) Similar video retrieval: Exemplars

- > Key frames \Rightarrow exemplars \Rightarrow their numbers are variable.
- > Each frame is a feature-extracted vector such that

 $\mathbf{x}_i = [x_{i1}, \cdots, x_{id}]^T, \ x_{ik} \ge 0, \quad \sum_{k=1}^d x_{ik} = 1.$ > Then, \mathbf{x}_i is in a simplex.

> Exemplars are in this space.







Fig. Exemplar frames reflecting time course.





- > A class of the harmonic competitive learning.
- Generalized vector k-means are found.
- > Number of competitive representatives is pre-specified.
- > Each exemplar is found as the nearest frame to the representative.

(1-8) How can we compare exemplar sets?

 $\underline{\text{M-distance}} = d\left(\{\text{exempla}_{i}, \text{context}_{i}\}_{i=1}^{N_{A}}, \{\text{exempla}_{j}, \text{context}_{j}\}_{j=1}^{N_{B}}\right)$

M-distance 1 (context-aware global alignment distance)

- > Naming is after the Levenshtein distance (edit distance).
- L-distance was a very special case:
 - ✓ Discrete alphabet
 - ✓ Identical or not identical · · · {0, 1}
 - ✓ Letter-to-letter (no ability reflect contexts)
 - ✓ Gap penalty is only letter-wise.
- M-distance 2 (context-aware local alignment distance) exists too.

(1-9) How can we compare exemplar sets? (cnt)

				Video B						
	ı —	→ j		exemplar 1	exemplar 2	exemplar 3	exemplar 4			
	↓ i			$E^{B}_{1} = 2$	$E^{B}_{2} = 3$	$E^{B}_{3} = 2$	$E^{B_{4}} = 1$			
			0	← - 0.4	← - 1	← - 1.4	← -1.6			
Video A	exemplar 1	$E^{A}_{1} = 2$	↑ - 0.4	× 1.054	٥.743 N	← 0.343	← 0.143			
	exemplar 2	$E_{2}^{A} = 2$	↑ -0.8	<u>م</u> 0.834	step 2.597	← 2.197	← 1.997			
	exemplar 3	$E_{3}^{A} = 3$	↑ -1.4	0.468	s 2.176	× 3.715	← 3.515			

Fig. Computation of M-distance for the global alignment.

(1-10) How do we compare exemplar sets in DB?



Fig. Exemplar set shuffling for similarity measurement: Sumo wrestlers and comic story tellers from NHK archive.

- Enquete has many queries.
- > Target persons are tired of answering \Rightarrow many blanks.
- > If we need to estimate blank items \Rightarrow EM family, NNMD.
- We put 0 to blank items.
- ➤ In this example, we consider a problem of GUI design.

Fig. Enlarged local user map: Icon positions are more uniform.

(with Mr. H, Kamiya and Dr. R. Yokote)

FIEEE IET ICALIP 2014

Fig. Doubly spherical GUI for all-recording TV database (with Mr. M. Maejima).

Fig. Doubly spherical GUI for NIRS brain signals with/without a task (with Mr. T. Horie).

(3-1) Security via brain signals

- > Big data per se \Rightarrow GUI \Rightarrow security by authentication.
- Bypassing the password
 (DARPA, the N.Y. Times, Mar. 17, 2012).
- Brain signals (NIRS, near infrared spectroscopy).

Fig. NIRS measurement for user authentication.

Fig. SVM used as a filter.

33

(4-2) Committee decision with crowd sourcing

Fig. Identification of human and nonhuman; e. g., Google reCAPTCHA.

Just one of these two queries are used for the authentication. The other is used for the improvement of machine's ability.

IEEE IET ICALIP 201

- Fig. Crowd sourcing example in bioinformatics.
 - (a) Foldit by University of Washington.
 - (b) 3D view is this presenter's own by using an NCBI's tool.

Note: These illustrations should be used only for academic uses.

- Common methods exist for a variety of big data.
- Machine learning gives one of such foundations. Seeds for new methods exist.
- Compatibility with human sensibility needs to be enhanced more.
- Committee-based methods with crowd sourcing will be quite powerful.

