



## Chapter 8. Classification: Ensembles

Meng Jiang

CSE 40647/60647 Data Science Fall 2017

Introduction to Data Mining

# Turing Award Recipients

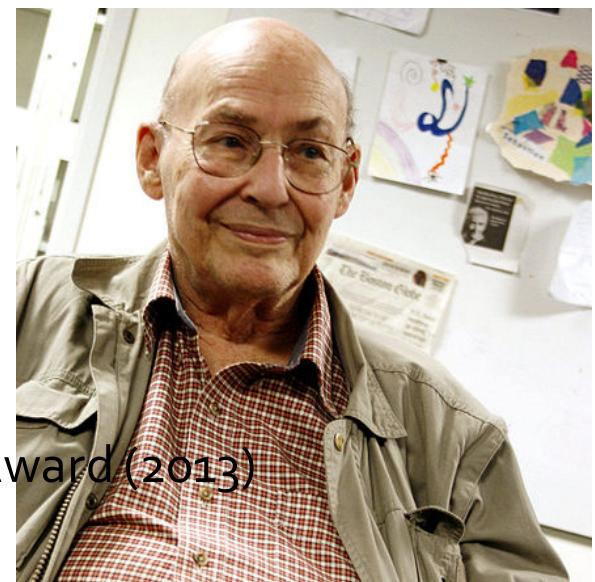
...	...	...
1998	Jim Gray	For seminal contributions to <b>database and transaction processing research</b> ...
1999	Frederick P. Brooks, Jr.	For landmark contributions to <b>computer architecture, operating systems, and software engineering</b> ...
2000	Andrew Chi-Chih Yao	<b>Theory of computation</b> , pseudorandom number generation, cryptography, and communication complexity...
2001	Ole-Johan Dahl Kristen Nygaard	<b>Object-oriented programming</b> , Simula I and Simula 67...
2002	Ronald L. Rivest, Adi Shamir, Leonard M. Adleman	For their ingenious contribution for making <b>public-key cryptography</b> useful in practice.
2003	Alan Kay	Contemporary <b>object-oriented programming languages</b> ...
2004	Vinton G. Cerf Robert E. Kahn	<b>Internetworking</b> , including the design and implementation of the Internet's basic communications protocols, TCP/IP...
2005	Peter Naur	<b>Programming language design</b> , ALGOL 60, compiler design...
2006	Frances E. Allen	<b>Optimizing compiler techniques</b> ...

# Turing Award Recipients (cont.)

2007	Edmund M. Clarke, E. Allen Emerson and Joseph Sifakis	For their roles in developing model checking into a <b>highly effective verification technology</b> , widely adopted in the hardware and software industries...
2008	Barbara Liskov	<b>Programming language and system design...</b>
2009	Charles P. Thacker	<b>Xerox Alto, the 1<sup>st</sup> modern PC, Ethernet and Tablet PC....</b>
2010	Leslie G. Valiant	Theory of computation...
2011	Judea Pearl	<b>Artificial intelligence through the development of a calculus for probabilistic and causal reasoning...</b>
2012	Silvio Micali Shafi Goldwasser	Transformative work that laid the complexity-theoretic foundations for the science of <b>cryptography</b> ...
2013	Leslie Lamport	Contributed to <b>distributed and concurrent systems...</b>
2014	Michael Stonebraker	Concepts underlying <b>modern database systems...</b>
2015	Martin E. Hellman Whitfield Diffie	Introduced <b>public-key cryptography</b> , the foundation for the most regularly-used security protocols on the Internet...
2016	Tim Berners-Lee	Invented the <b>World Wide Web</b> and the <b>first web browser...</b>

# Marvin Minsky

- Marvin Lee Minsky (August 9, 1927 – January 24, 2016) was an American cognitive scientist concerned largely with research of **artificial intelligence** (AI), co-founder of the Massachusetts Institute of Technology's AI laboratory, and author of several texts concerning AI and philosophy.
- Awards
  - **Turing Award (1969)**
  - Japan Prize (1990)
  - IJCAI Award for Research Excellence (1991)
  - Benjamin Franklin Medal (2001)
  - Computer History Museum Fellow (2006)
  - BBVA Foundation Frontiers of Knowledge Award (2013)



Logical vs. Analogical  
or  
Symbolic vs. Connectionist  
or  
Neat vs. Scruffy

<https://web.media.mit.edu/~minsky/papers/SymbolicVs.Connectionist.html>

Marvin Minsky

In Artificial Intelligence at MIT, Expanding Frontiers, Patrick H. Winston (Ed.), Vol.1, MIT Press, 1990. Reprinted in AI Magazine, Summer 1991.

# Ensembles

“To solve really hard problems, we’ll have to use *several different representations*...

It is time to *stop arguing* over *which* type of pattern-classification technique *is best*...

Instead we should work at a higher level of organization and discover how to build *managerial systems* to exploit the *different virtues* and evade the *different limitations* of each of these ways of comparing things.” [Minsky, 1991]

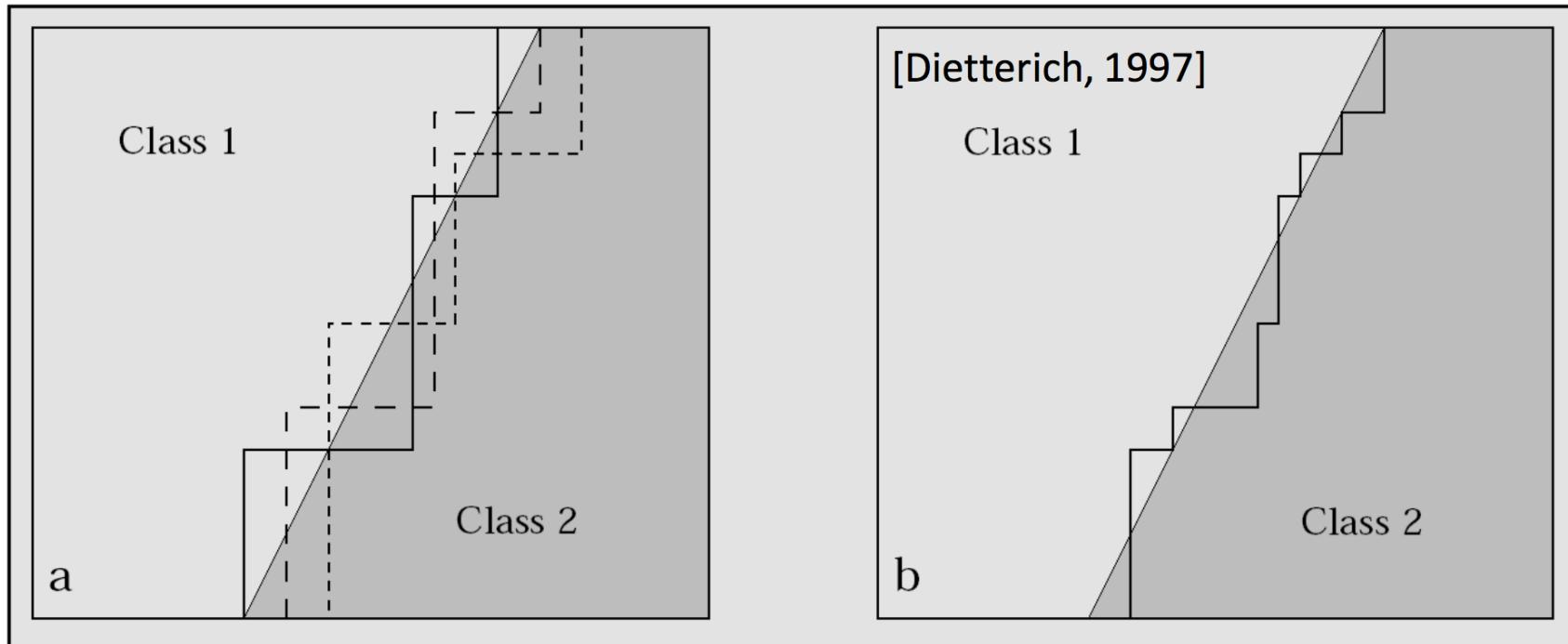


# Ensembles (cont.)

- An ensemble is a set of classifiers that learn a target function, and their *individual predictions are combined* (weighted or unweighted) to classify new examples
  - Each classifier should be *more accurate than by chance*, and independent of one another.
  - Usually *more accurate than a single classifier*.
- Ensembles generally improve the *generalization performance* of a set of classifiers on a domain.

# Ensemble Methods

- Ensemble methods
  - Use a *combination of models* to *increase accuracy*
  - Combine a series of  $k$  learned models,  $M_1, M_2, \dots, M_k$ , with the aim of creating an *improved model  $M^*$*



# Ensemble Methods (cont.)

- Popular ensemble methods
  - **Bagging:** averaging the prediction over a collection of classifiers
  - **Boosting:** weighted vote with a collection of classifiers
    - **AdaBoost (Adaptive Boosting):** adaptive in the sense that subsequent weak learners are tweaked in favor of those instances misclassified by previous classifiers.



# Bagging

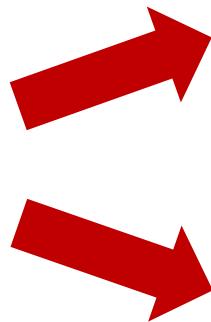
- Training
  - Given a data set  $D$  of  $d$  instances, a classifier model  $M_i$  is learned for a training set  $D_i$  of  $d$  instances that is *sampled with replacement* from  $D$  ( $i = 1 \dots k$ )
  - As a result of the *sampling-with-replacement* procedure, each classifier is trained on approximately **63.2%** of the training examples
  - For a dataset with  $d$  instances, each instance has a probability of  **$1 - (1 - 1/d)^d$**  of being selected at least once in the  $d$  samples.
    - For  $d \rightarrow \infty$ , this number converges to  $(1 - 1/e)$  or 0.632 [Bauer and Kohavi, 1999]

# Bagging (cont.)

- Classification: classify an unknown sample  $X$ 
  - Each classifier  $M_i$  returns its class prediction
  - The bagged classifier  $M^*$  *counts the votes* and assigns the class with the *most votes* to  $X$
- Accuracy: Proved improved accuracy in prediction
  - *Often significantly better* than a single classifier derived from  $D$

# Bagging (cont.)

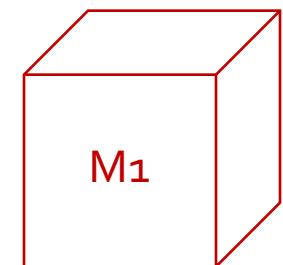
	Features	Label
1		
2		
3		
<i>Training D</i>		
4		
5		
6		



	Features	Label
7		
8		
9		
<i>Test D<sub>test</sub></i>		
10		

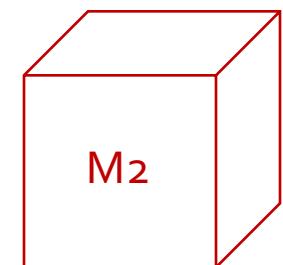
	Features	Label
3		
4		
5		
4		
5		
6		

*Training D<sub>1</sub>*



	Features	Label
1		
2		
5		
2		
5		
6		

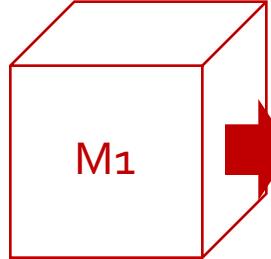
*Training D<sub>2</sub>*



# Bagging (cont.)

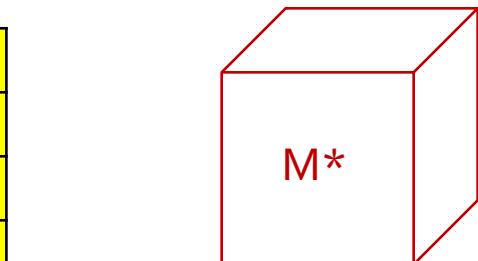
	Features	Label
3		
4		
5		
4		
5		
6		

*Training D<sub>1</sub>*



	Features	Label
7		
8		
9		
10		

*M<sub>1</sub> on D<sub>test</sub>*

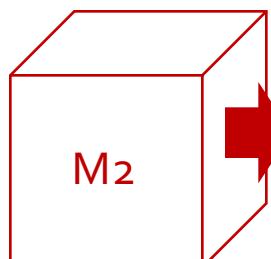


	Features	Label
7		
8		
9		
10		

*Majority voting on D<sub>test</sub>*

	Features	Label
1		
2		
5		
2		
5		
6		

*Training D<sub>2</sub>*



	Features	Label
7		
8		
9		
10		

*M<sub>2</sub> on D<sub>test</sub>*

*DT<sub>1</sub> (ID<sub>3</sub>): If the object is red, it is an apple not a banana.*

*DT<sub>2</sub> (ID<sub>3</sub>): If the object is round, it is an apple not an banana.*

# Majority Voting

- *Example:* (Course Project) Entity type recognition
  - Features: Triggers in contextual words
    - The probability that the technical term is a “method” (“problem”, “dataset”, “metric”, etc.).
  - Each context setting as a classifier: contextual words.

<http://www.meng-jiang.com/teaching/TypingDemo.zip>

# Boosting

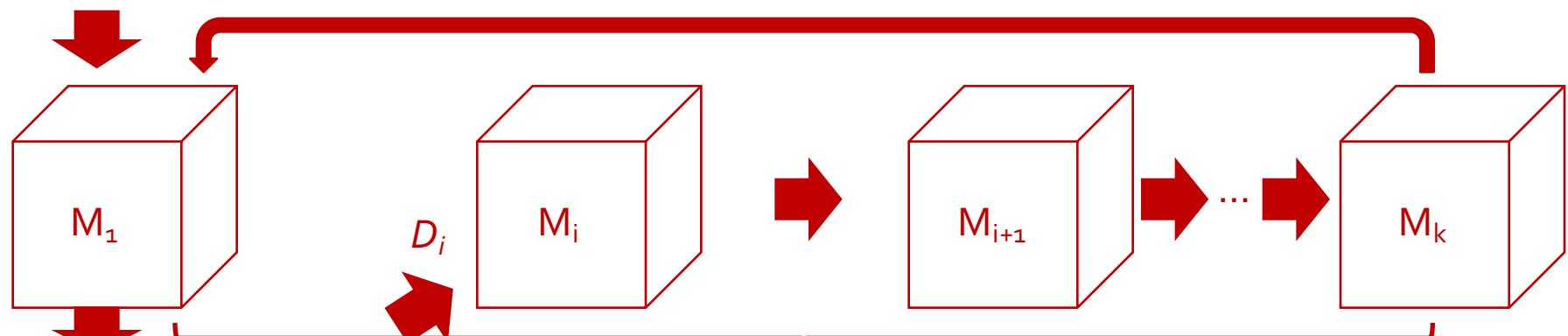
- Training
  - *Weights* are assigned to each training instance
  - A series of  $k$  classifiers is *iteratively* learned
  - After a classifier  $M_i$  is learned, the *weights* are updated to allow the subsequent classifier,  $M_{i+1}$ , to pay more attention to the *training* instances that were *misclassified* by  $M_i$
- Classification
  - The final  $M^*$  *combines the votes* of each individual classifier, where the *weight* of each classifier's vote is a function of its *accuracy* on classifying training instances
- Comparing with Bagging: Boosting tends to have *greater accuracy*, but it risks *overfitting* the model to misclassified data

	Features	Label
3		
4		
5		
4		
5		
6		

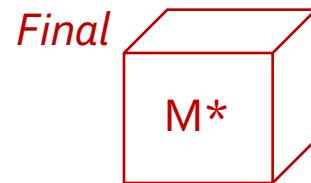
	Features	Label	Weight
1			1.0
2			1.0
3			1.0
4			1.0
5			1.0
6			1.0

	Features	Label
3		
...		
...		

Iteratively



	Predicted	Label	Weight
3	Yes	No	3.0
4	Yes	Yes	1.0



# How were Chinese middle school students learning English?

2012年中考题 41

Stop \_\_\_\_\_ about the traffic.

Just think about what we can do to improve it.

- A. complain
- B. to complain
- C. complaining
- D. complained

2012年中考题 68.

Jessica solved the physics problems without any help.  
(保持句意不变)

Jessica \_\_\_\_\_ the physics problems without any help.

我的反思

答案: B C

[stop to do sth. 停下来去做某事]

[stop doing sth. 停止做某事]

[或只能接 doing sth.]

practice, finish, give up,  
be worth, enjoy, mind,

look forward to, be busy,  
spend time/money doing sth.

答案: worked out

(注意时态)

Work out<sup>①</sup> = figure out 算出  
(物) 算出, 制定出

(中考大纲) 需求(语法)

② 锻炼

# AdaBoost (Adaptive Boosting)

- Given a set of  $d$  class-labeled instances,  $(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_d, y_d)$
- Initially, all the *weights* of instances are set the same ( $1/d$ )
- Generate  $k$  classifiers in  $k$  rounds. At round  $i$ ,
  - Instances from  $D$  are *sampled with replacement* to form a training set  $D_i$  of the same size
  - Each instance's chance of being selected is based on its *weight*
  - A classification model  $M_i$  is derived from  $D_i$
  - Its *error rate* is calculated *using  $D_i$  as a "test set"*
  - If an instance is misclassified, its *weight* is increased, otherwise it is decreased

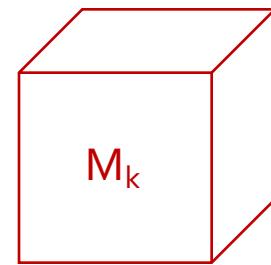
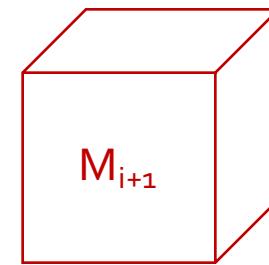
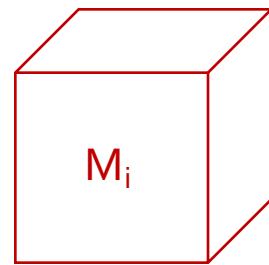
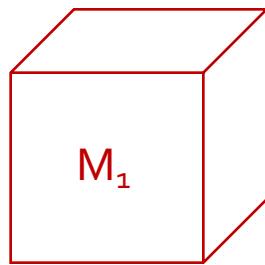
# AdaBoost (cont.)

	Features			Label	Weight
1					1.0
2	<i>Training D</i>				1.0
3					1.0
4					1.0
5					1.0
6					1.0



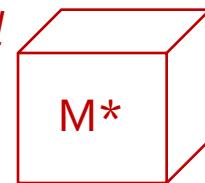
	Features			Label	Weight
1					3.0
2					2.1
3					0.5
4					1.5
5					1.0
6					1.0

Error rate (if high)



Vote (then small)

Final





# AdaBoost (cont.)

- *Error rate:*  $err(\mathbf{X}_j)$  is the misclassification error of instance  $\mathbf{X}_j$ . Classifier  $M_i$ 's error rate is the sum of the weights of the misclassified instances:

$$error(M_i) = \sum_j^d w_j \times err(\mathbf{X}_j)$$

- The *weight* of classifier  $M_i$ 's vote is

$$\log \frac{1 - error(M_i)}{error(M_i)}$$

# Summary

- ❑ Ensembles
- ❑ Bagging
- ❑ Boosting
  - ❑ AdaBoost

# An Example

“SetExpan: Corpus-Based Set Expansion via Context  
Feature Selection and Rank Ensemble”

by Shen et al. at UIUC

ECML PKDD 2017

SKOPJE, MACEDONIA  
18-22 SEPTEMBER

THE EUROPEAN CONFERENCE ON MACHINE LEARNING &  
PRINCIPLES AND PRACTICE OF KNOWLEDGE DISCOVERY IN DATABASES



# Macedonia



# Skopje



# An Example (cont.)

“**SetExpan: Corpus-Based Set Expansion via Context Feature Selection and Rank Ensemble**”

by Shen et al. at UIUC

<http://mickeystroller.github.io/resources/ECMLPKDD2017.pdf>

# References

- C. Apte and S. Weiss. Data mining with decision trees and decision rules. Future Generation Computer Systems, 13, 1997
- P. K. Chan and S. J. Stolfo. Learning arbiter and combiner trees from partitioned data for scaling machine learning. KDD'95
- A. J. Dobson. An Introduction to Generalized Linear Models. Chapman & Hall, 1990.
- R. O. Duda, P. E. Hart, and D. G. Stork. Pattern Classification, 2ed. John Wiley, 2001
- U. M. Fayyad. Branching on attribute values in decision tree generation. AAAI'94.
- Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. J. Computer and System Sciences, 1997.
- J. Gehrke, R. Ramakrishnan, and V. Ganti. Rainforest: A framework for fast decision tree construction of large datasets. VLDB'98.
- J. Gehrke, V. Gant, R. Ramakrishnan, and W.-Y. Loh, BOAT -- Optimistic Decision Tree Construction. SIGMOD'99.
- T. Hastie, R. Tibshirani, and J. Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer-Verlag, 2001.
- T.-S. Lim, W.-Y. Loh, and Y.-S. Shih. A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms. Machine Learning, 2000

# References (cont.)

- J. Magidson. The Chaid approach to segmentation modeling: Chi-squared automatic interaction detection. In R. P. Bagozzi, editor, *Advanced Methods of Marketing Research*, Blackwell Business, 1994
- M. Mehta, R. Agrawal, and J. Rissanen. SLIQ : A fast scalable classifier for data mining. EDBT'96
- T. M. Mitchell. *Machine Learning*. McGraw Hill, 1997
- S. K. Murthy, Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey, *Data Mining and Knowledge Discovery* 2(4): 345-389, 1998
- J. R. Quinlan. Induction of decision trees. *Machine Learning*, 1:81-106, 1986.
- J. R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
- J. R. Quinlan. Bagging, boosting, and c4.5. AAAI'96.
- R. Rastogi and K. Shim. **Public: A decision tree classifier that integrates building and pruning**. VLDB'98
- J. Shafer, R. Agrawal, and M. Mehta. **SPRINT : A scalable parallel classifierfor data mining**. VLDB'96
- J. W. Shavlik and T. G. Dietterich. **Readings in Machine Learning**. Morgan Kaufmann, 1990
- P. Tan, M. Steinbach, and V. Kumar. **Introduction to Data Mining**. Addison Wesley, 2005
- S. M. Weiss and C. A. Kulikowski. **Computer Systems that Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems**. Morgan Kaufman, 1991
- S. M. Weiss and N. Indurkhya. **Predictive Data Mining**. Morgan Kaufmann, 1997
- I. H. Witten and E. Frank. **Data Mining: Practical Machine Learning Tools and Techniques**, 2ed. Morgan Kaufmann, 2005