Data mining opportunities

Charles Elkan

 \diamond

University of California, San Diego



The charge given to us

 $\langle \bullet \rangle$

* "Highlight methods / models / techniques which may be useful in the analysis of climate data and understanding climate change."

What can computer science offer climate science?

12

ę.

63

8

(h)

Virtual and Artificial, but 58,000 Want Course

By JOHN MARKOFF Published: August 15, 2011

 \diamond

PALO ALTO, Calif. — A free online course at <u>Stanford University</u> on artificial intelligence, to be taught this fall by two leading experts from Silicon Valley, has attracted more than 58,000 students around the globe — a class nearly four times the size of Stanford's entire student body.

The two additional courses will be an introductory course on database software, taught by Jennifer Widom, chairwoman of the computer science department, and an introduction to machine learning, taught by Andrew Ng.

* Students* New thinking

What can climate science offer computer science?

05/20/2005 - Updated 11:09 AM ET

 $\langle \bullet \rangle$

Predicting hurricanes

Each December, Colorado State University's tropical storm research team forecasts the coming year's Atlantic hurricane season, which is June 1 through Nov. 30.

	(Winds of at least 39 mph)		(Winds of at least 74 mph)			
					(111 mph or more winds)	
	Forecast	Actual	Forecast	Actual	Forecast	Actual
1998	9	14	5	9	2	3
1999	14	12	9	8	4	5
2000	11	14	7	8	3	3
2001	9	15	5	9	2	4
2002	13	12	8	4	4	2
2003	12	-	8	(4)	3	

Source: Department of Atmospheric Science, Colorado State University

(b) Each year, a research team at Colorado State University predicts how many hurricanes (of three different intensity levels) will strike the eastern U.S. the next year. The predicted and actual numbers are attached. Explain how to use your answers to Problem 3 and part (a) to decide whether or not the forecasts have any value.

* "[In 2005] forecast 15 named storms, eight of them hurricanes." [USA Today, 2008]
* Actual: 28 storms and 15 hurricanes, including Katrina.

 \diamond

* "After the 2005 Atlantic hurricane season, Gray announced that he was stepping back from the primary authorship of CSU's tropical cyclone probability forecasts." [Wikipedia].

* "... forecasts that 16 named storms will form in the Atlantic ... 72% chance of a major hurricane striking land." [USA Today, May 2011]

* "Gray ... does not attribute global warming to anthropogenic causes, and is critical of those who do." [Wikipedia]

 $\langle \bullet \rangle$

The charge given to us

 $\langle \bullet \rangle$

* "Highlight methods / models / techniques which may be useful in the analysis of climate data and understanding climate change."

Standard versus new techniques

* Standard: clustering
* Each point is mapped to a single cluster center.
* New: topic modeling
* Each point is generated by *multiple* centroids.

 $\langle \bullet \rangle$

Standard: principal components analysis
Linear method to reduce dimensionality.
New: link prediction and collaborative filtering
Hidden dimensions are inferred to *predict* outcomes.

Topic modeling

In a topic model, each data point is a *combination* of prototypes.
Nonlinear, applicable to both discrete and continuous data.

Topic 1 $\bar{\theta_s} = .021$	Topic 2 $\bar{\theta_z} = .019$	Topic 3 $\bar{\theta_s} = .017$	Topic 4 $\bar{\theta_z} = .017$
Southwestern Energy Range Resources Cabot Oil & Gas EOG Resources Chesapeake Energy Pioneer Resources Devon Energy Peabody Energy Anadarko Petroleum	Penneys Macys Kohls Nordstrom Target Limited Lowes Home Depot	Capital One BNY Mellon Discover Northern Trust Janus JPMorgan Chase State Street Wells Fargo	Simon Property Kimco Realty Equity Residential AvalonBay Communities Apartment Investment Vornado Realty Trust Boston Properties Public Storage Host Hotels
Massey Energy	Abercrombie	T. Rowe Price	HCP Inc.

- * *Financial topic models*, G. Doyle and C. Elkan, Proc. NIPS Workshop on Applications for Topic Models: Text and Beyond, 2009.
- * Accounting for burstiness in topic models, G. Doyle and C. Elkan, Proc. Intl Conference on Machine Learning, 2009.



* A. K. Menon and C. Elkan, *Link prediction via matrix factorization*, Proc. ECML 2011.

Collaborative filtering

 $\langle \bullet \rangle$



* A. K. Menon and C. Elkan, *Predicting labels for dyadic data*. Data Mining and Knowledge Discovery (2):327-343, 2010.

Standard versus new techniques

***** Two more examples

 $\langle \bullet \rangle$

Learning a classifier from only positive examples

* Training a model to label land cover based on remote-sensing data



- * W. Li, Q. Guo, C. Elkan, *A Positive and Unlabeled Learning Algorithm for One-Class Classification of Remote-Sensing Data*. IEEE Transactions on Geoscience and Remote Sensing 49(2)717-725, 2011.
- * C. Elkan, K. Noto, *Learning classifiers from only positive and unlabeled data*, Proc. KDD 2008, pp. 213-220.



* Use two- or three- dimensional conditional random fields