#### Data In, Facts Out: Automated Monitoring of Facts by FactWatcher

Chengkai Li Professor, Department of Computer Science and Engineering Director, Innovative Database and Information Systems Research (IDIR) Laboratory University of Texas at Arlington



### FactWatcher



Tuple t for new real
world event appended
to database

			$t_7$	Wesley	25
Constraint	Measure			. 1	
month=Feb	pts, ast, reb		t	na cons is in the	con
opp_team=Nets	ast, reb				
team= <i>Celtics</i> & opp_team= <i>Nets</i>	ast, reb	Genera	te fa	actual cla	aim
•••					

id	player	day	month	season	team	opp_team	pts	ast	reb
<i>t</i> <sub>1</sub>	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
<i>t</i> <sub>2</sub>	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
<i>t</i> <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
$t_4$	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
<i>t</i> <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
$t_6$	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
<i>t</i> <sub>7</sub>	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5

Find constraint-measure pair (C, M) such that t is in the contextual skyline

Wesley had 12 points, 13 assists and 5 rebounds on February 25, 1996 to become the first player with a 12/13/5 (points/assists/rebounds) in February.

http://en.wikipedia.org/wiki/Basketball

Fact Finding

#### Prominent streaks

Long consecutive subsequence of high values in a sequence

#### One-of-the-few objects

Qualifying statements that can only be made for very few objects

#### Situational facts



FactWatcher Finds Three Types of Facts (and can be Extended)

#### Domains

• sports, weather, crimes, transportation, finance, social media analytics

#### Examples from Real News Media

#### Prominent streaks

- "This month the Chinese capital has experienced 10 days with a maximum temperature in around 35 degrees Celsius – the most for the month of July in a decade." http://www.chinadaily.com.cn/china/2010-07/27/content\_11055675.htm
- "The Nikkei 225 closed below 10000 for the 12th consecutive week, the longest such streak since June 2009."

http://www.bloomberg.com/news/articles/2010-08-06/japanese-stocks-fall-for-second-day-this-week-on-u-s-jobless-claims-yen



FactWatcher Finds Three Types of Facts (and can be Extended)

Examples from Real News Media Situational facts, One-of-the-few objects

- "Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992."
- "The social world's most viral photo ever generated 3.5 million likes, 170,000 comments and 460,000 shares by Wednesday afternoon." http://www.cnbc.com/id/49728455



FactWatcher Demo http://idir.uta.edu/factwatcher/ https://vimeo.com/user48311227



FactWatcher	NBA 312	Weather	Nov 1, 1991 💶 📒	🗩 Apr 20, 2005 Speed <mark>1x 🚽</mark> 🜌
			Feb 22, 1998	

#### »LIVE UPDATE

[February 20, 1998] Todd Fuller had 1 assist, 3 steals and 1 block in the Golden State Warriors' defeat against the Denver Nuggets. It is one of the best performance made by him.

#### Presented In

SEARCH	michael jordan	FACT TYPE
	Michael Adonis Jordan	SITUATIONAL FACT
-A	Michael Jordan	PROMINENT STREAK
4	Michael Michael Jordan	ONE-OF-THE-FEW
	Michael Reggie Jordan	
<b>I</b> C 3	Michael Thomas Jordan	RANKING
<b>11</b>	[January 13, 1997] Horace Grant had 26 points and 6 assists in the Orlando Magic's victory	THIS
_	against the New Jersey Nets. It is one of the best performance made by him.	INTERESTINGNESS
<b>I</b> G 2	[January 13, 1997] After the Orlando Magic's win over the New Jersey Nets, for the first time in his career, Rony Seikaly had at least 20 points for 6 consecutive games, after today's game.	THIS POPULARITY
<b>1</b>	[January 13, 1997] Horace Grant had 26 points and 2 steals in the Orlando Magic's victory against the New Jersey Nets. It is one of the best performance made by him.	THIS
ı¢ 5	[January 13, 1997] Horace Grant had 26 points, 6 assists and 2 steals in the Orlando Magic's MORELIKET	TEAMS
	victory against the New Jersey Nets. It is one of the best performance made by him.	SEASONS
<b>1</b> 6 3	[January 13, 1997] After the Orlando Magic's victory against the New Jersey Nets, for the first time in his career. Rony Seikaly had at least 20 points and 8 rehounds for 6 consecutive games.	1996-97 (9) THIS
	after today's game.	1994-95 (5)
4	[January 13, 1997] Nick Anderson had 8 assists and 2 blocks in the Orlando Magic's win over the	1992-93 (1)
2	New Jersey Nets. It is one of the best performance made by him.	+MORE



#### Excellent Demo Award

LESS-

# COMPUTATION NO + JOURNALISM 14 SYMPOSIUM 44

#### http://idir.uta.edu/factwatcher/



#### How were these Facts Discovered in Current Systems?

#### Our (educated?) guess

- Experts monitor real-world events (e.g., watching an NBA game), have a gut-feeling, issue database queries, check out or not
- Prepared facts-to-be (e.g., Nowitzki only needs 477 more points to surpass O'Neal. Perhaps will happen around Christmas 2015)
- Predefined templates of facts/database queries
- Perhaps in-house systems/algorithms similar to FactWatcher







The World's Foremost Sports Statisticians and Historians

X

StatSheet

#### No. 1-Seeded Louisville Clips No. 4-Seeded Michigan 82-76, Wins NCAA Championship

Filed under Game Recap on April 9th, 2013

#### Share this recap

🔊 No. 1-Seeded Louisville Cl 🗴

Tweet Or ▲Like I One person likes this. Be the first of your friends.

#### NCAA Tournament 7th Round



#### Mon, Apr 08 2013, 10:23 PM EDT

Georgia Dome Atlanta, Georgia Attendance: 74,326 TV: CBS

🗲 🔶 🖸 🗋 thevilledaily.com/louisville-basketball/game-recap/no-1-seeded-louisville-clips-no-4-seeded-michigan-82-76-wins-ncaa-championship

#### Boxscore | Game Notes | Game Recap | StatSmack

No. 1-seeded Louisville got the win against No. 4-seeded Michigan 82-76 in the Championship Game of the NCAA Tournament on Monday, Apr. 8. The Cardinals were led by Peyton Siva, who got 18 points and six rebounds (5 Ast 4 Stl). Gorgui Dieng also had an outstanding outing, scoring eight points and adding eight rebounds (6 Ast 3 Blk). Michigan closes out its impressive season with a 31-8 overall record. The Wolverines got to the NCAA Tournament as an at-large team after falling to Wisconsin 68-59 in the Big Ten Tournament. In the regular season, they finished fourth in the Big Ten with a 12-6 conference record. In making the national championship game, Michigan knocked off No. 13-seeded South Dakota State 71-56 in the second round and No. 5-seeded Virginia Commonwealth 78-53 in the third round. Following that, the Wolverines got through No. 1-seeded Kansas 87-85 in the Sweet Sixteen, No. 3-seeded Florida 79-59 in the Elite Eight, and No. 4-seeded Syracuse 61-56 in the Final Four. For the Wolverines, Trey Burke got a game-high 24 points and four rebounds. Michigan (31-8) finished the regular season fourth in the Big Ten with a 12-6 record. Through their amazing run, Louisville got through No. 16-seeded North Carolina A&T 79-48 in the second round and No. 8-seeded Colorado State 82-56 in the third round. Following that, the Cardinals got through No. 12-seeded Oregon 77-69 in the Sweet Sixteen, No. 2-seeded Duke 85-63 in the Elite Eight, and No. 9-seeded Wichita State 72-68 in the Final Four.

#### • StatSeed: NCAA Automatic #1 Seed

- O X

☆ **=** 





More about Fan Satisfaction

#### Find another NCAA team:



#### Categories

#### Narrative

#### Science

\_ 🗆 🗙 Forbes Earnings Preview: / × C 🗅 www.forbes.com/sites/narrativescience/2013/05/02/forbes-earnings-preview-anadarko-petroleum-6/ Forbes -Popular Video New Posts Lists Search The Global 2000 Narrative Science xerox 🌒 Forbes Partner Narrative + Follow (50) INVESTING | 5/02/2013 @ 2:31PM | 488 views

0

f Share

13

😏 Tweet

1

in Share

0

Submit

0

Q +1

#### Forbes Earnings Preview: Anadarko Petroleum

By Narrative Science

+ Comment Now + Follow Comments

Analysts have become increasingly bullish on <u>Anadarko</u> <u>Petroleum</u> <u>APC +2 02%</u> (APC) in the month leading up to the company's first quarter earnings announcement scheduled for Monday, May 6, 2013. The consensus earnings per share estimate has moved up from 88 cents a share to the current expectation of earnings of 91 cents a share.

<u>Wall Street</u> projections are down 1.1% year-over-year, as the company reported earnings of 92 cents per share.

The consensus estimate has gone up, from 82 cents, over the past three months. Analysts are expecting earnings of \$4.04 per share for the fiscal year. Revenue is projected to be \$3.49 billion for the quarter, 1.2% above the year-earlier total of \$3.45 billion. For the year, revenue is projected to roll in at \$15.21 billion.

Revenue has declined for the third quarter in a row. The year-over-



Handling \$421 billion in accounts payables annually for companies like yours.



Ready For Real Business XEFOX

### Publications

- o Online Frequent Episode Mining. Xiang Ao, Ping Luo, Chengkai Li, Fuzhen Zhuang, and Qing He. ICDE 2015, pages 891-902.
- Data In, Fact Out: Automated Monitoring of Facts by FactWatcher. Naeemul Hassan, Afroza Sultana, You Wu, Gensheng Zhang, Chengkai Li, Jun Yang, and Cong Yu. VLDB 2014, pages 1557-1560. Demonstration description. (excellent demonstration award)
- Finding, Monitoring, and Checking Claims Computationally Based on Structured Data. Brett Walenz, You (Will) Wu, Seokhyun (Alex) Song, Emre Sonmez, Eric Wu, Kevin Wu, Pankaj K. Agarwal, Jun Yang, Naeemul Hassan, Afroza Sultana, Gensheng Zhang, Chengkai Li, Cong Yu. 2014 Computation+Journalism Symposium.
- Incremental Discovery of Prominent Situational Facts. Afroza Sultana, Naeemul Hassan, Chengkai Li, Jun Yang, Cong Yu. ICDE 2014, pages 112-123.
- Discovering General Prominent Streaks in Sequence Data. Gensheng Zhang, Xiao Jiang, Ping Luo, Min Wang, Chengkai Li. ACM TKDD, 8(2):article 9, June 2014.
- Discovering and Learning Sensational Episodes of News Events. Xiang Ao, Ping Luo, Chengkai Li, Fuzhen Zhuang, Qing He, and Zhongzhi Shi. WWW 2014, pages 217-218.
- o On "One of the Few" Objects. You Wu, Pankaj K. Agarwal, Chengkai Li, Jun Yang, Cong Yu. KDD 2012, pages 1487-1495.
- o Prominent Streak Discovery in Sequence Data. Xiao Jiang, Chengkai Li, Ping Luo, Min Wang, Yong Yu. KDD 2011, pages 1280-1288.



Incremental Discovery of Prominent Situational Facts. Afroza Sultana, Naeemul Hassan, Chengkai Li, Jun Yang, Cong Yu. ICDE 2014, pages 112-123.



"Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992." (http://espn.go.com/espn/elias?date=20130205)



"Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992." (http://espn.go.com/espn/elias?date=20130205)



"Paul George had 21 points, 11 rebounds and 5 assists to become the first Pacers player with a 20/10/5 (points/rebounds/assists) game against the Bulls since Detlef Schrempf in December 1992." (http://espn.go.com/espn/elias?date=20130205)



"The social world's most viral photo ever generated 3.5 million likes, 170,000 comments and 460,000 shares by Wednesday afternoon."

(http://www.cnbc.com/id/49728455/President Obama Sets New Social Media Record)



"The social world's most viral photo ever generated 3.5 million likes, 170,000 comments and 460,000 shares by Wednesday afternoon."

(http://www.cnbc.com/id/49728455/President Obama Sets New Social Media Record)



"The social world's most viral photo ever generated 3.5 million likes, 170,000 comments and 460,000 shares by Wednesday afternoon."

(http://www.cnbc.com/id/49728455/President Obama Sets New Social Media Record)



- •Stock Data: Stock A becomes the first stock in history with price over \$300 and market cap over \$400 billion.
- •Weather Data: Today's measures of wind speed and humidity are x and y, respectively. City B has never encountered such high wind speed and humidity in March.
- •Criminal Records: There were 50 DUI arrests and 20 collisions in city C yesterday, the first time in 2013.

Financial Analyst Journalists Scientists Citizens DiR &

id	player	day	month	season	team	opp_team	pts	ast	reb
$t_1$	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
<i>t</i> <sub>2</sub>	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
<i>t</i> <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
<i>t</i> <sub>4</sub>	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
<i>t</i> <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
<i>t</i> <sub>6</sub>	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
$t_7$	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5

Last tuple appended to table



id	player	day	month	season	team	opp_team	pts	ast	reb
$t_1$	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
<i>t</i> <sub>2</sub>	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
$t_4$	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
<i>t</i> <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
<i>t</i> <sub>6</sub>	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
<i>t</i> <sub>7</sub>	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5



id	player	day	month	season	team	opp_team	pts	ast	reb
$t_l$	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
<i>t</i> <sub>2</sub>	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
<i>t</i> <sub>4</sub>	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
<i>t</i> <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
$t_6$	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
<i>t</i> <sub>7</sub>			Feb.				12	13	5



id	player	day	month	season	team	opp_team	pts	ast	reb
$t_1$	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
<i>t</i> <sub>2</sub>	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
<i>t</i> <sub>4</sub>	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
<i>t</i> <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
$t_6$	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
<i>t</i> 7			Feb.				12	13	5

•Wesley had 12 points, 13 assists and 5 rebounds on February 25, 1996 to become the first player with a 12/13/5 (points/assists/rebounds) in February.

id	player	day	month	season	team	opp_team	pts	ast	reb
$t_1$	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
$t_2$	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t3	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
$t_4$	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
<i>t</i> <sub>6</sub>	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
<i>t</i> <sub>7</sub>				1995-96			12	13	5



id	player	day	month	season	team	opp_team	pts	ast	reb
$t_1$	Bogues	11	Feb.	1991-92	Homets	Hawks	4	12	5
$t_2$	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
<i>t</i> <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
<i>t</i> <sub>4</sub>	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
$t_6$	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
<i>t</i> <sub>7</sub>					Celtics	Nets		13	5

•Wesley had 13 assists and 5 rebounds on February 25, 1996 to become the second Celtics player with a 13/5 (assists/rebounds) game against the Nets.

**Dimension space:**  $\mathcal{D}=\{d_1,\ldots,d_n\}$ 

**Measure space:**  $\mathcal{M} = \{m_1, \dots, m_s\}$ 

							-			
id	player	day	month	season	team	opp_team		pts	ast	reb
$t_1$	Bogues	11	Feb.	1991-92	Hornets	Hawks		4	12	5
<i>t</i> <sub>2</sub>	Seikaly	13	Feb.	1991-92	Heat	Hawks		24	5	15
<i>t</i> <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets		13	13	5
<i>t</i> <sub>4</sub>	Wesley	4	Feb.	1994-95	Celtics	Nets		2	5	2
<i>t</i> <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves		3	5	3
$t_6$	Strictland	3	Jan.	1995-96	Blazers	Celtics		27	18	8

append-only table



#### $\Box \text{Constraint}(C): d_1 = v_1 \land d_2 = v_2 \land \ldots \land d_n = v_n, v_i \in dom(d_i) \cup \{*\}$

■ team=*Celtics* ∧ opp\_team=*Nets* 

id	player	day	month	season	team	opp_team	pts	ast	rb
$t_1$	Bogues	11	Feb.	1991-92	Hornets	Hawks	1	12	5
$t_2$	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
<i>t</i> <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
<i>t</i> <sub>4</sub>	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t <sub>6</sub>	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8



# Constraint-Measure Pair (C, M): Combination of a constraint and measure subspace

(team=Celtics ^ opp\_team=Nets, {assists, rebounds})

id	player	day	month	season	team	opp_team	pts	ast	reb
$t_{I}$	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
$t_2$	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
t <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
<i>t</i> <sub>4</sub>	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t <sub>6</sub>	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8



**Contextual skyline:** skyline regarding (C, M)

# σ<sub>team=Celtics ∧ opp\_team=Nets</sub>(R), M={assists, rebounds} {t<sub>3</sub>}

id	player	day	month	season	team	opp_team	pts	ast	reb
$t_{I}$	Bogues	11	Feb.	1991-92	Hornets	Hawks		12	5
$t_2$	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
<i>t</i> <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
$t_4$	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
t <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
t <sub>6</sub>	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8



### FactWatcher



Tuple t for new real
world event appended
to database

			$t_7$	Wesley	25
Constraint	Measure			. 1	
month=Feb	pts, ast, reb		t	na cons is in the	con
opp_team=Nets	ast, reb				
team= <i>Celtics</i> & opp_team= <i>Nets</i>	ast, reb	Genera	te fa	actual cla	aim
•••					

id	player	day	month	season	team	opp_team	pts	ast	reb
<i>t</i> <sub>1</sub>	Bogues	11	Feb.	1991-92	Hornets	Hawks	4	12	5
<i>t</i> <sub>2</sub>	Seikaly	13	Feb.	1991-92	Heat	Hawks	24	5	15
<i>t</i> <sub>3</sub>	Sherman	7	Dec.	1993-94	Celtics	Nets	13	13	5
$t_4$	Wesley	4	Feb.	1994-95	Celtics	Nets	2	5	2
<i>t</i> <sub>5</sub>	Wesley	5	Feb.	1994-95	Celtics	Timberwolves	3	5	3
$t_6$	Strictland	3	Jan.	1995-96	Blazers	Celtics	27	18	8
<i>t</i> <sub>7</sub>	Wesley	25	Feb.	1995-96	Celtics	Nets	12	13	5

Find constraint-measure pair (C, M) such that t is in the contextual skyline

Wesley had 12 points, 13 assists and 5 rebounds on February 25, 1996 to become the first player with a 12/13/5 (points/assists/rebounds) in February.

http://en.wikipedia.org/wiki/Basketball



Т

					-		{tttttt_}	
id	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> <sub>3</sub>	<b>m</b> <sub>1</sub>	<i>m</i> <sub>2</sub>		((19023034903)	
$t_l$	$a_1$	$b_2$	$c_2$	10	15	$a_1$	$\mathfrak{b}_1$	$\mathbf{r}_1$
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10	$\{t_1, t_2, t_5\}$	$\{t_2, t_3, t_4, t_5\}$	$\{t_2, t_4, t_5\}$
t3	$a_2$	$b_1$	$c_2$	17	17			
<i>t</i> <sub>4</sub>	$a_2$	$b_1$	$c_1$	20	20	$\begin{bmatrix} a_1, b_1 \\ b_1 \\ b_2 \\ b_1 \\ b_1 \\ b_2 \\ b_1 \\ b_1 \\ b_2 \\ b_1 \\ b_2 \\ b_1 \\ b_2 \\ b_2 \\ b_1 \\ b_1 \\ b_2 \\ b_2 \\ b_2 \\ b_1 \\ b_2 \\ b_2 \\ b_2 \\ b_1 \\ b_2 \\ b_2$	(t .t )	$\{t, t, t\}$
<i>t</i> <sub>5</sub>	$a_1$	$b_1$	$c_1$	11	15			
							$a_1, b_1, c_1$	
				C	$l_1 = a_1 \wedge d$	$a=b_1 \wedge d_3 = c_1$	$\{t_2, t_5\}$	
							Lattice of $C^{t}$	5
Tup	le Sat	isfied	Const	raint a	⊄: If ∀	$d_i \in \mathcal{D}, C.d_i$	$=* \text{ or } C.d_i = t.d_i$	$l_i$ , t satisfies
С.							<b>i</b> Di	R 🗍

Μ	od	e]	lin	Ø
				0

Lattice of $C^{t_4}$
$a_1  a_2  b_1  c_1$
$a_{1}b_{1}a_{2}b_{1}a_{1}c_{1}a_{2}c_{1}b_{1}c_{1}$
Lattice of $C^{t_5}$ $a_1, b_1, c_1, a_2, b_1, c_1$



id	<i>d</i> <sub>1</sub>	$d_2$	$d_3$	<b>m</b> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_{I}$	$a_1$	$b_2$	$c_2$	10	15
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10
t <sub>3</sub>	$a_2$	$b_1$	<i>c</i> <sub>2</sub>	17	17
<i>t</i> <sub>4</sub>	<i>a</i> <sub>2</sub>	<b>b</b> <sub>1</sub>	<i>c</i> <sub>1</sub>	20	20
<i>t</i> <sub>5</sub>	<i>a</i> <sub>1</sub>	<b>b</b> 1	<i>c</i> <sub>1</sub>	11	15

id	<i>d</i> <sub>1</sub>	$d_2$	$d_3$	<b>m</b> <sub>1</sub>	$m_2$
$t_1$	$a_1$	$b_2$	$c_2$	10	15
$t_2$	$a_1$	$b_1$	$c_1$	15	10
<i>t</i> <sub>3</sub>	$a_2$	$b_1$	<i>c</i> <sub>2</sub>	17	17
<i>t</i> <sub>4</sub>	<i>a</i> <sub>2</sub>	<b>b</b> <sub>1</sub>	<i>c</i> <sub>1</sub>	20	20
<i>t</i> <sub>5</sub>	<i>a</i> <sub>1</sub>	<b>b</b> <sub>1</sub>	<i>c</i> <sub>1</sub>	11	15



Lattice Intersection:  $C^{t_4,t_5} = C^{t_4} \cap C^{t_5}$ 



## **Brute-Force Approach**

id	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> 3	<i>m</i> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_1$	$a_1$	$b_2$	$c_2$	10	15
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10
t <sub>3</sub>	$a_2$	$b_1$	$c_2$	17	17
<i>t</i> <sub>4</sub>	$a_2$	$b_1$	$c_1$	20	20
<i>t</i> <sub>5</sub>	$a_1$	$b_1$	<i>C</i> 1	11	15




id	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> 3	<i>m</i> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_1$	$a_1$	$b_2$	$c_2$	10	15
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10
t <sub>3</sub>	$a_2$	$b_1$	$c_2$	17	17
$t_4$	$a_2$	$b_1$	$c_1$	20	20
<i>t</i> <sub>5</sub>	$a_1$	$b_1$	<i>C</i> <sub>1</sub>	11	15





id	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> 3	<i>m</i> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_1$	$a_1$	$b_2$	$c_2$	10	15
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10
<i>t</i> <sub>3</sub>	$a_2$	$b_1$	$c_2$	17	17
<i>t</i> <sub>4</sub>	$a_2$	$b_1$	$c_1$	20	20
$t_5$	$a_1$	$b_1$	<i>C</i> 1	11	15





id	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> 3	<i>m</i> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_1$	$a_1$	$b_2$	$c_2$	10	15
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10
t <sub>3</sub>	$a_2$	$b_1$	$c_2$	17	17
<i>t</i> <sub>4</sub>	$a_2$	$b_1$	$c_1$	20	20
<i>t</i> <sub>5</sub>	$a_1$	$b_1$	<i>C</i> 1	11	15





id	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> 3	<i>m</i> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_1$	$a_1$	$b_2$	$c_2$	10	15
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10
t <sub>3</sub>	$a_2$	$b_1$	$c_2$	17	17
<i>t</i> <sub>4</sub>	$a_2$	$b_1$	$c_1$	20	20
$t_5$	$a_1$	$b_1$	<i>c</i> <sub>1</sub>	11	15





id	<i>d</i> <sub>1</sub>	<i>d</i> <sub>2</sub>	<i>d</i> 3	<i>m</i> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_1$	$a_1$	$b_2$	$c_2$	10	15
<i>t</i> <sub>2</sub>	$a_1$	$b_1$	$c_1$	15	10
t <sub>3</sub>	$a_2$	$b_1$	$c_2$	17	17
<i>t</i> <sub>4</sub>	$a_2$	$b_1$	$c_1$	20	20
<i>t</i> <sub>5</sub>	$a_1$	$b_1$	<i>C</i> 1	11	15





 $\mathbf{X}$ 

id	<i>d</i> <sub>1</sub>	$d_2$	<i>d</i> <sub>3</sub>	<i>m</i> <sub>1</sub>	<i>m</i> <sub>2</sub>
$t_1$	$a_1$	$b_2$	$c_2$	10	15
$t_2$	$a_1$	$b_l$	<i>c</i> <sub>1</sub>	15	10
t <sub>3</sub>	$a_2$	$b_1$	$c_2$	17	17
$t_4$	$a_2$	$b_1$	$c_1$	20	20
$t_5$	$a_1$	$b_1$	<i>c</i> <sub>1</sub>	11	15

 $\begin{array}{c|c} a_1 & \times \mathbf{b}_1 & \times \mathbf{c}_1 \\ a_1, \mathbf{b}_1 & a_1, \mathbf{c}_1 & \times \mathbf{b}_1, \mathbf{c}_1 \\ & a_1, \mathbf{b}_1, \mathbf{c}_1 \end{array}$ 

Total  $|R|^*(2^{|\mathcal{D}|+|\mathcal{M}|}-1)$  comparisons! Total 16 comparisons in this case!





Exhaustive comparison with every tuple
Under every constraint
Over every measure subspace



# **Experiment Setup**

#### NBA Dataset

- 317,371 tuples of NBA box scores from 1991-2004 seasons
- 8 dimension attributes
- 7 measure attributes
- Weather Dataset
  - 7.8 million tuples of weather forecast from different locations of six countries & regions of UK
  - 7 dimension attributes
  - 7 measure attributes



## **Discovered Facts**

- Lamar Odom had 30 points, 19 rebounds and 11 assists on March 6, 2004. No one before had a better or equal performance in NBA history.
- Allen Iverson had 38 points and 16 assists on April 14, 2004 to become the first player with a 38/16 (points/assists) game in the 2004-2005 season.
- Damon Stoudamire scored 54 points on January 14, 2005. It is the highest score in history made by any Trail Blazers.



# **Prominent Streaks**

Prominent Streak Discovery in Sequence Data. Xiao Jiang, Chengkai Li, Ping Luo, Min Wang, Yong Yu. KDD 2011, pages 1280-1288.

Discovering General Prominent Streaks in Sequence Data. Gensheng Zhang, Xiao Jiang, Ping Luo, Min Wang, Chengkai Li. ACM TKDD, 8(2):article 9, June 2014.



## Prominent Streaks

#### Prominent streaks stated in news articles:

"This month the Chinese capital has experienced 10 days with a maximum temperature in around 35 degrees Celsius – the most for the month of July in a decade."

"The Nikkei 225 closed below 10000 for the 12th consecutive week, the longest such streak since June 2009."

"He (LeBron James) scored 35 or more points in nine consecutive games and joined Michael Jordan and Kobe Bryant as the only players since 1970 to accomplish the feat."



#### Concepts Streak

Input: a sequence of values

Streak <[l, r], v> is a triple: left-end ( l ), right-end ( r ), minimum value in interval [l,r] 3 1 7 7 2 <u>5 4 6</u> 7 3

<[6, 8], 4>

Streak dominance relation

s1=<[l1, r1], v1> dominates s2=<[l2, r2], v2> if and only if r1 - l1 > r2 - l2, v1 >= v2 or r1 - l1 >= r2 - l2, v1 > v2

Prominent streaks (PS)

A streak is prominent if it is not dominated by any other streaks.





3 1 7 7 2 5 4 6 7 3



#### Prominent Streaks are Skyline Points in 2-d Space 3 1 7 7 2 5 4 6 7 3



Each streak can be viewed as a point in a 2-dimension space (length, minimal\_value). Prominent streaks are skyline points in this space.

Tasks

### Task 1: discovery

Find all prominent streaks in a sequence

#### Task 2: monitoring

Always keep prominent streaks up-to-date, when sequence grows (real-world sequences often grow)



## Solution Framework



Candidate Generation: Number Of Candidate streaks Brute-force

Quadratic (How many streaks are in a sequence of n values? These are all candidate streaks.)

NLPS

Superlinear

LLPS

Linear



Given a sequence of n values, how many streaks are there? That's the number of candidate streaks the bruteforce algorithm will produce.

 $n^{*}(n+1)/2$ 

. . .

- Streaks starting at the 1<sup>st</sup> position: n
- Streaks start at the 2<sup>nd</sup> position: n-1

• Streaks start at the nth position: 1



## Local Prominent Streak

#### Local dominance relation

s1=<[l1, r1], v1> locally dominates s2=<[l2, r2], v2> if and only if a) s1 dominates s2 and b) [l1, r1]  $\supset$  [l2, r2] (i.e., s1 subsumes/contains s2 and s2 is thus a sub-sequence of s1). Local prominent streak (LPS)

A streak is locally prominent if it is not locally dominated by any other streaks.



# Important Properties

#### A prominent streak must be an LPS.

(2) LPS is small

The number of LPSs is less than or equal to the sequence length. (Hint: The number of LPSs getting min value at position k is at most 1.) Conclusion

LPS is an excellent set of candidate streaks, of linear size.

The candidate generation problem becomes finding local prominent streaks.

#### Local prominent streaks (candidates)









#### Linear LPS (LLPS) Method

Given a sequence  $p_1, p_2, ..., p_n, ...,$  maintain a list of growing streaks when scanning the sequence rightward.

- 1. At  $p_{1,}$  create a growing streak to include the 1<sup>st</sup> position.
- 2. After  $p_k$ , right-ends of growing streaks are all at position k.
- 3. At  $p_{k+1}$ , try to extend the growing streaks rightward. The growing streaks are partitioned by  $p_{k+1}$  into two groups: (3.a) Those growing steaks with s.v  $\leq p_{k+1}$ : extend s to include position k+1.

(3.b) Those growing streaks with s.v >  $p_{k+1}$ : include s into LPS; remove s from growing streaks as it won't grow anymore.

4. At  $p_{k+1}$ , do one of the following, based on which condition is satisfied:

(4.a) There was a growing streak with s.v =  $p_{k+1}$ : nothing more needs to be done for  $p_{k+1}$ .

(4.b) There were one or more growing streaks with  $s.v > p_{k+1}$ : create a new growing streak by extending the leftmost (longest) such growing streak s that satisfies  $s.v > p_{k+1}$  to include position k+1.

(4.c) There was no growing streak with s.v  $\geq p_{k+1}$ : create a new growing streak to include just position k+1.



## Linear LPS (LLPS) Method

Growing streaks share the same right-end. Their minimum values monotonically increase, if they are listed in the increasing order of left-ends.

Two different ways of illustrating the method at  $p_{10}$ , with all remaining growing streaks highlighted



# Linear LPS (LLPS) Method

Monitoring (keeping prominent streaks up-to-date) is simple:

- Keeps applying the algorithm when new values of the sequence are received.
- After the latest value at p<sub>n</sub>, it has found all LPSs ending before position n.
- Growing streaks ending at n either will eventually be LPSs or can be grown into LPSs ending at a future position after n.
- If prominent streaks till n are requested, compare all found LPSs and all growing streaks.



### Datasets In Experiments

name	length	# prominent streaks	description
Gold	1074	137	Daily morning gold price in US dollars, 01/1985-03/1989.
River	1400	93	Mean daily flow of Saugeen River near Port Elgin, 01/1988-12/1991.
Melb1	3650	55	The daily minimum temperature of Melbourne, Australia, 1981-1990.
Melb2	3650	58	The daily maximum temperature of Melbourne, Australia, 1981-1990.
Wiki1	4896	58	Hourly traffic to en.wikipedia.org/wiki/Main_page, 04/2010-10/2010.
Wiki2	4896	51	Hourly traffic to en.wikipedia.org/wiki/Lady_gaga, 04/2010-10/2010.
Wiki3	4896	118	Hourly traffic to en.wikipedia.org/wiki/Inception_(film),04/2010-10/2010.
SP500	10136	497	S&P 500 index, 06/1960-06/2000.
HPQ	12109	232	Closing price of HPQ in NYSE for every trading day, 01/1962-02/2010.
IBM	12109	198	Closing price of IBM in NYSE for every trading day, 01/1962-02/2010.
AOL	132480	127	Number of queries sent to AOL search engine in every minute over three months.
WC98	7603201	286	Number of requests to World Cup 98 web site in every second, 04/1998-07/1998.





(b) Prominent Streaks

#### Sample Prominent Streaks Melbourne daily min/max temperature between 1981 and 1990 (Melb1 & Melb2)

More than 2000 days with min temperature above zero



Traffic count of Wikipedia page of Lady Gaga (Wiki2) More than half of the prominent streaks are around Sep. 12th (VMA 2010) at least 2000 hourly visits lasting for almost 4 days



## General Prominent Streaks

Top-k, multi-dimensional and multi-sequence PS

"He (LeBron James) scored 35 or more points in nine consecutive games and joined Michael Jordan and Kobe Bryant as the only players since 1970 to accomplish the feat."

"Only player in NBA history to average at least 20 points, 10 rebounds and 5 assists per game for 6 consecutive seasons." (http://en.wikipedia.org/wiki/Kevin Garnett)

NLPS/LLPS extended to such general PSs



# Experiments On Multi-Sequence PSs

Table IX. Multi-sequence Prominent Streaks in Datast NBA1.

<b>n</b> 11	length	minimal value	players
	1	71	David Robinson
	2	51	Allen Iverson; Antawn Jamison
	4	42	Kobe Bryant
	9	40	Kobe Bryant
	13	35	Kobe Bryant
	14	32	Kobe Bryant
	16	30	Kobe Bryant
	17	27	Michael Jordan
	27	26	Allen Iverson
	34	24	Tracy McGrady
	45	21	Allen Iverson
	57	20	Allen Iverson
	74	19	Shaquille O'Neal
	94	18	Shaquille O'Neal
	96	17	Karl Malone
	119	16	Karl Malone
	149	15	Karl Malone
	159	14	Karl Malone
	263	13	Karl Malone
	357	12	Karl Malone
	527	11	Karl Malone
	575	10	Karl Malone
	758	7	Karl Malone
	858	6	Shaquille O'Neal
	866	2	Karl Malone
	932	1	John Stockton
	1185	0	Jim Jackson



## Experiments On Multi-Dim PSs

Table X. Data Sequences Used in Experiments on Multi-dimensional Prominent Streak Discovery.



Fig. 13. Experiments on Increasing Dimensionality.



## Experiments On General PSs

Table XIII. Data Sequences Used in Experiments on Top-5 Multi-sequence Multi-dimensional Prominent Streak Discovery.

name	# sequences	average length	# dimensions	# prominent streaks	description
NBA2	1185	290	6	10867	1991-2004 game log of all N- BA players (minutes, points, re- bounds, assists, steals, blocks)

Table XIV. Number of Candidate Streaks, Top-5 Multi-sequence Multi-dimensional Prominent Streak Discovery.

name	Baseline	NLPS	LLPS
NBA2	$9.41 \times 10^{7}$	$2.98 imes10^6$	$8.76  imes 10^{5}$

Table XV. Execution Time (in Milliseconds), Top-5 Multi-sequence Multi-dimensional Prominent Streak Discovery.

name	Baseline	NLPS	LLPS
NBA2	$1.39 \times 10^{7}$	$4.33 \times 10^{5}$	$1.14 \times 10^{5}$





R



Apply the LLPS algorithm to find prominent streaks in 4 1 3 5 4 7 3 6. Show local prominent streaks and growing streaks at every position.



# **One-of-the-few Objects**

On "One of the Few" Objects. You Wu, Pankaj K. Agarwal, Chengkai Li, Jun Yang, Cong Yu. KDD 2012, pages 1487-1495



## One-Of-The-Few Claims

- Do these claims really hold water?
- Karl Malone is ONE OF THE ONLY TWO players in NBA history with 25,000 points, 12,000 rebounds, and 5,000 assists in one's career.
- He is ONE OF THE ONLY THREE candidates who have raised more than 25% from PAC contributions and 25% from self-financing.
- How do we find truly interesting claims or individuals?



## X Is One-Of-K $\rightarrow$ X Is In K-Skyband

#### Claim

Karl Malone is ONE OF THE ONLY TWO players in NBA history with 25,000 points, 12,000 rebounds, and 5,000 assists in one's career.

#### General claim

Fewer than k objects dominate X in subspace of attributes  $S \subseteq \{A_1, A_2, ..., A_d\}$ 

*k*-skyband [Papadias et al. 2005] in *S* is the set of points each dominated by fewer than *k* other points in *S*  $1\frac{1}{5}$   $1\frac{1}{5}$ 

1-skyband : skyline

1-skyband 2-skyband



# Small K ≠Interesting

#### Subspaces are different

E.g., 2-skyand in {rebounds} vs. in {rebounds, assists}





# Small K ≠Interesting

#### Data distribution matters

E.g., 2-skyand in {points, rebounds} vs. in {rebounds, assists}



# Top-τ Skyband

#### k-Skyband

Using the same k for all subspaces doesn't work Asking user pick k for each subspace is infeasible Top- $\tau$  Skyband

- o User specifies a single parameter  $\tau$  to cap # skyband objects.
- For each subspace S, find its top- $\tau$  skyband, i.e., the largest k-skyband containing no more than  $\tau$  objects
- E.g., in {points, rebounds}:
- $\tau = 2 \rightarrow 1$ -skyband (size 2)
- $\tau = 6 \rightarrow 2$ -skyband (size 5; 3-skyband would be too big)



