



#### INTEGRATED VEHICLE HEALTH MANAGEMENT

#### Event Cube: An Organized Approach for Mining and Understanding Anomalous Aviation Events

#### PI Name: Jiawei Han

Aviation Safety Program Technical Conference November 17-19, 2009 Washington D.C.

## Outline





### **Problem Statement**



#### Aviation Safety Reporting System

## How to mine ASRS for analyzing patterns and causes of aviation anomalies?

Light	Anomaly	Phase	Report
Night	equipment : critical	intermediate altitude	we hit what appeared to be a flock of large white birds. Some of them seem to be ducks
Daylight	equipment : critical	intermediate altitude	our flight encountered a small flock of ducks
Daylight	equipment : less severe	ground: parked	there were no bird strikes noticed while taxiing

January 1981 – December 2006

- 40,000 reports in 2006
- 3,300 reports /month
- 160+/working day

- What are possible causes of landing problems? (cause analysis)
- Is there any difference in causal factors for landing at night and daylight time? (comparative analysis)
- Where did landing problems most often happen? (ranking)

Source: http://asrs.arc.nasa.gov/overview/summary.html



 General mission: Develop advanced technology for mining text data and interrelated, multidimensional datasets for effective and efficient IVHM

#### • Approach:

- Develop Event Cube technology for effective analysis of causal and contributing factors in aviation anomalies
- Develop a set of data mining and text analysis methods for systematic, multidimensional analysis of ASRS database
- Further extend such event cube and text mining methods for effective analysis of other text data sets of NASA's interest

#### **IVHM Milestones Addressed**





**3.2.4** An automated capability to diagnose the causal factors of anomalous operations in real or emulated data of large, fleet-wide or airspace heterogeneous data sources.

**3.1.3** An anomaly detection method that has the ability to <u>detect</u> at least 3 known anomalies in real or emulated data of large, fleetwide heterogeneous data sources.

#### Potentially,

**3.3.4** Forecasting technology that has the ability to <u>predict</u> at least 3 known anomalies in real or emulated data of large, fleetwide heterogeneous data sources.

# Approach: Event Cube

Multidimensional OLAP, Ranking, Cause Analysis,

**Support Topic Summarization/Comparison** ..... Topic Topic turbulence Encounter birds undershoot **Event Cube Deviation** overshoot Representation AX SJC MIA AUS CA Location Location Light Anomaly Phase Report Multidimensional Night equipment : intermediate ...we hit what appeared to be a flock of large white critical altitude birds. Some of them seem to be ducks.. ...our flight encountered a small flock of ducks... Daylight equipment : intermediate **Text Database** critical altitude Daylight eauioment : less around: ... there were no bird strikes noticed while taxiing... severe parked

Analysis

2009 Aviation Safety Program Technical Conference

--- ---



- 1. Topic/Text Cube for comparative analysis
- 2. Cell ranking and summarization
- 3. Causal factor annotation

### 1. Text/Topic Cube: General Idea





## Topic Cube: Multidimensional Text Analysis by Topic Modeling





- Classification, clustering and summarization of each cell
  - Cluster the documents in one cell to reveal aspects of anomalous events
  - Summarization of documents in a cell
- Comparative analysis of anomalous events in contexts, e.g.,
  - In what kind of weather condition, does a flight tend to have the problem of landing without clearance?
  - In Place A and Time B, which kind of anomaly event occurs most often?
- Keyword search in text cube
  - Ranking relevant cells

. . .

- Ranking cuboids (in the order of discriminativeness of keywords)
- Cause analysis based on text/topic cube
  - E.g. what are the main causal factors of each anomalous event in different situations?

#### How can we efficiently do all these at large-scale?

#### Aggregation Speeds Up TopicCube Construction





## Sample Result : Topic Cube for Topic Content Comparison



#### landing without clearance: daylight vs. night

Context	Word	p(w θ)	WINDS ALOFT AT PATTERN ALT OF
	Tower	0.075	1000 FI MSL, WERE MUCH SIRONGER
	Pattern	0.061	SAY. THE PATTERNS AND LNDGS
dovlight	Final	0.060	WERE DIFFICULT FOR MY STUDENT
daylight	Runway	0.053	AND THERE WAS LIGHT TURB ON THE
	Land	0.052	DOWNWIND
	Downwind	0.039	
	Tower	0.035	I LISTENED TO HWD ATIS AND
	Runway	0.029	FOUND THE TWR CLOSED AND AN
	Light	0.027	INTENSITY LIGHTS FOR RWY 28L
night	Instrument Landing System	0.015	WERE INOP. BROADCASTING IN THE BLIND AND LOOKING FOR THE TWR BEACON AND LOW INTENSITY
	Beacon	0.014	LIGHTS AGAINST A VERY BRIGHT
			BACKGROUND CLUTTER OF STREET
			ILIGHTS ETC

#### Sample Result: TopicCube for Topic Coverage Comparison



#### Comparison of distributions of anomalies in FL, TX, and CA



#### **Comparative Analysis of Shaping Factors**





2009 Aviation Safety Program Technical Conference

#### Text Cube for Comparative Analysis of Sub-Events





#### Leverage Sequential Pattern Mining: More Meaningful Units



- Anomaly1 = aircraft equipment problem : critical
- Anomaly2 = inflight encounter : weather
- Anomaly3 = conflict : nmac

Pattern	Support			
	Anomaly I	Anomaly2	Anomaly3	
LNDG UNEVENTFUL	П	0	0	
LANDED WITHOUT INCIDENT	12	0	0	
SHUT DOWN ENG	12	0	0	
VISIBILITY FOG	0	13	0	
CEILING VISIBILITY	0	15	0	
DOWNWIND RWY	0	0	12	
SAW OTHER ACFT	0	0	10	
CLRED FOR RWY	0	0	44	
TOOK EVASIVE ACTION	0	0	44	
SUPPLEMENTAL FROM	17	10	31	
CALLBACK WITH REVEALED FOLLOWING	37	13	24	
CALLBACK WITH REVEALED FOLLOWING HAT	13	0	0	



## 1. Topic/Text Cube for comparative analysis

- 2. Cell ranking and summarization
- 3. Causal factor annotation

## 2. Cell Ranking and Summarization



#### Ranking

• Find critical anomaly and flight phases related to birds by searching flight reports with keyword query: "bird", "geese", "duck"

Light	Anomaly	Phase	Report
Night	equipment : critical	intermediate altitude	we hit what appeared to be a flock of large white <b>birds</b> . Some of them seem to be <b>duck</b> s
Daylight	equipment : critical	intermediate altitude	our flight encountered a small flock of ducks
Daylight	equipment : less severe	ground: parked	there were no <b>bird</b> strikes noticed while taxiing
			Cell (Light=Daylight

Cell (Anomaly=equipment : critical, Phase=intermediate altitude)

- Traditional IR: Ranking all documents
- **Our Approach**: Ranking all cells in the order of relevance given keyword query

## **Summarization of Cell Content**



- A fundamental challenge in Online Analytical Processing (OLAP) of multidimensional text database is to <u>summarize</u> the content in its text cells.
  - 1. <u>Neutral Summarization</u>

Give the most representative documents within a text cell

2. Topic-biased Summarization

Give the most relevant documents to a query within a text cell that also cover the content of the text cell well

ACN	Time	Airport	•••	Light	Narrative
101285	199901	MSP	•••	Daylight	Document 1
101286	199901	CKB		Night	Document 2
101291	199902	LAX	• • •	Dawn	Document 3

Table 1: An example of text database in ASRS

#### • For example:

- 1. What are the main points made in those reports about anomalies during night in Jan. 1999? (neutral)
- 2. What have the pilots said about landing at LAX in 1999? (topic-biased: topic="landing")



Micro Clustering



#### Neutral Summarization

...so that if we saw the ARPT, we could land... ...due to stronger than forecasted winds and weather going...

...resulted in RWY excursion during engine fail...

#### Topic-biased Summarization

...so that if we saw the ARPT, we could <u>land</u>...

...after an hour, the weather had not much improved which forced us to <u>land</u>...

...SMA engine failure, forced landing at LGB by instructor...

#### Micro-Text-Cluster is More Efficient than Direct Summarization





2009 Aviation Safety Program Technical Conference



- Topic/Text Cube for comparative analysis
  Cell ranking and summarization
- 3. Causal factor annotation



- Goal
  - Identify the cause(s) of an incident described in an ASRS report narrative from a set of 14 predefined causes, or shaping factors.
- Challenges
  - Reports are informally written (acronyms, abbreviations, and typos)
  - Reports can have multiple labels (shaping factors)
  - Labeled data is scarce
  - Skewed class distributions: 10 of the 14 classes account for 25% of labels

#### Approach 1: Semantic Lexicon Construction



- Motivation: insufficient annotated data implies that it is difficult to identify relevant words for each shaping factor
- Idea: learn words and phrases related to each shaping factor in a bootstrapping manner.
- Experimental Results:
  - Both rule-based and learning-based approaches outperform the baseline.
  - The learned words are useful features.
  - The best result of 52.7% illustrates the difficulty of the task.



- Motivation: insufficient annotated data
- Idea: develop a bootstrapping algorithm to augment the labeled data in an iterative fashion.
- Repeat the following in each iteration:
  - augment either the positively labeled data or the negatively labeled data, depending on which set is smaller
  - find the word w that best discriminates the positives and the negatives based on the currently labeled data
  - label all unlabeled documents that contain w with the corresponding label
- Experimental Results:
  - This approach is specially effective in improving the F-score for the 10 minority (i.e., under-represented) shaping factors, where positive instances are scarce.



Shaper	<b>Positive Expanders</b>	Negative Expanders
Familiarity	unfamiliar, layout, unfamiliarity, rely	
<b>Physical Environment</b>	cloud, snow, ice, wind	
<b>Physical Factors</b>	fatigue, tire, night, rest, hotel, awake, sleep, sick	declare, emergency, advisory, separation
Preoccupation	distract, preoccupied, awareness, situational, task, interrupt, focus, eye	declare, ice, snow, crash, fire, rescue, anti, smoke
Pressure	bad, decision, extend, fuel, calculate, reserve, diversion, alternate	

## Conclusions



- EventCube: A framework for mining and analyzing ASRS database to help achieve three IHVM milestones: <u>Detection</u>, <u>Diagnosis</u>, and <u>Prognosis</u>
  - Unification of relational data mining and text mining
  - Integrative analysis of ASRS data
- We proposed algorithms for constructing specific EventCube instances:
  - Text Cube and Topic Cube support OLAP on text dimension and comparative analysis of anomalies
  - Sequential pattern mining enhances analysis with pattern units
- We developed algorithms for interactive analysis:
  - Multidimensional keyword search: TopCells
  - Multidimensional cell summarization: MiTexCluster
- We developed algorithms for automatic analysis
  - Causal factor analysis
  - Data stream classification
- We also developed a prototype system



Different word distributions reflect different content in the two cells

Cell 1: (Weather="fog", \*) Cell 2: (Weather="Turbulence", \*)

Term	Probability	Term	Probability
runway	0.022331	feet	0.015784
approach	0.018184	aircraft	0.015541
feet	0.017086	flight	0.012476
aircraft	0.012505	flight level	0.009766
fiabt	0.012000	turbulence	0.009297
night	0.007746	runway	0.009082
tower	0.007044	approach	0.008561
weather	0.006380	altitude	0.007906
time	0.006355	time	0.006232
mile	0.006087	air traffic control	0.006059
airport	0.005806	weather	0.005340
visibility	0.005640		

Tania 424



#### Three sample topics in a cell

Topic #33

Term	Probability
flight	0.014851
turbulence	0.013191
aircraft	0.012837
feet	0.011509
approach	0.009760
runway	0.008078
kts	0.007171
air traffic control	0.007104
weather	0.007016
flight level	0.006197
captain	0.006064
deas	0 005644

1 opic #34						
Term	Probability					
feet	0.023412					
visual flight rules	0.017479					
flight	0.015342					
weather	0.014002					
approach	0.013109					
instrument flight rules	0.011738					
clouds	0.009824					
mile	0.009601					
ZZZ	0.009250					
conditions	0.008867					
airport	0.008325					

Горіс	#35
-------	-----

Term	Probability
wake	0.021715
aircraft	0.019422
runway	0.015918
feet	0.013929
turbulence	0.013453
approach	0.009560
flight	0.009387
degs	0.008435
air traffic control	0.008089
turbulence	0.006618
time	0.006186
kts	0.006056

## **System Demo: Cell Ranking**



Enter key words (space separated) below:					— Top Cells for keywo				word		
Key V	Key Words: excursion				Search		auerv="excursion"				
The top ranked cells are:						query excursion				$\leq$	
Rank	Year	State	Person	Weather	Light	Make/Model	Flight Phase	Primary Area	Event Anomaly	Resolutory Action	Score
1	2000	*	*	Rain	Night	*	landing : roll	*	aircraft equipment problem : critical	*	34.3084299
2	*	*	*	*	Daylight	Boeing	cruise : enroute altitude change	Flight Crew Human Performance	altitude deviation : overshoot	controller : issued advisory	33.3752199
3	*	*	*	*	*	McDonnell Douglas	*	Airport	excursion : taxiway	none taken : anomaly accepted	33.0382225
4	*	ОН	*	*	*	*	*	*	non adherence : far	flight crew : returned to original clearance	32.5993384



## Please visit me! http://infonetcube.cs.uiuc.edu/nasa/



- Further enhance EventCube through
  - Multidimensional cause analysis/classification/prediction
  - Multidimensional data stream management/analysis
  - Construction and analysis of information networks (entities, relations)
  - Natural language processing and text mining for causal analysis
- Closer collaboration with aviation experts to evaluate our prototype system and obtain feedback

## **Research Publications (Led by UIUC Group)**



- Bolin Ding, Bo Zhao, Cindy Xide Lin, Jiawei Han, and Chengxiang Zhai, "*TopCells: Keyword-Based Search of Top-k Aggregated Documents in Text Cube*", to appear in Proc. 2010 Int. Conf. on Data Engineering (ICDE'10).
- Duo Zhang, Chengxiang Zhai, Jiawei Han, Nikunj C. Oza, Ashok N. Srivastava, "*Topic Cube: Topic Modeling for OLAP on Multidimensional Text Databases*", Journal version invited to special issue of SDM'09 (to appear)
- Cindy Xide Lin, Bolin Ding, Jiawei Han, Nikunj C. Oza, Ashok N. Srivastava, Bo Zhao, Feida Zhu, "*TextCube: Computing IR Measures for Multidimensional Text Database Analysis*", submitted to IEEE Trans. on Knowledge and Data Engineering, 2009
- Duo Zhang, Chengxiang Zhai and Jiawei Han, "*Topic Cube: Topic Modeling for OLAP on Multidimensional Text Databases*", Proc. 2009 SIAM Int. Conf. on Data Mining (SDM'09), Sparks, NV, April 2009. (Best of SDM'09)
- Bolin Ding, David Lo, Jiawei Han, and Siau-Cheng Khoo, "Efficient Mining of Closed Repetitive Gapped Subsequences from a Sequence Database", Proc. 2009 Int. Conf. on Data Engineering (ICDE'09), Shanghai, China, Mar. 2009.
- Feida Zhu, Xifeng Yan, Jiawei Han and Philip S. Yu, "*Mining Frequent Approximate Sequential Patterns*", in H. Kargupta, et al., (eds.), Next Generation of Data Mining, Chapman & Hall, 2009.
- Cindy Xide Lin, Bolin Ding, Jiawei Han, Feida Zhu, and Bo Zhao, "*Text Cube: Computing IR Measures for Multidimensional Text Database Analysis*", Proc. 2008 Int. Conf. on Data Mining (ICDM'08), Pisa, Italy, Dec. 2008.
- Luiz Mendes, Bolin Ding, and Jiawei Han, "*Stream Sequential Pattern Mining with Precise Error Bounds*", Proc. 2008 Int. Conf. on Data Mining (ICDM'08), Pisa, Italy, Dec. 2008.
- Jing Gao, Bolin Ding, Wei Fan, Jiawei Han, and Philip S. Yu, "*Classifying Data Streams with Skewed Class Distribution and Concept Drifts*", IEEE Internet Computing (Special Issue on Data Stream Management), 12(6): 37-49, 2008
- Chen Chen, Cindy Xide Lin, Xifeng Yan, Jiawei Han, "On Effective Presentation of Graph Patterns: A Structural Representative Approach", Proc. 2008 ACM Conf. on Information and Knowledge Management (CIKM'08).

## **Research Publications (Led by UTD Group)**



- Muhammad Arshad Ul Abedin, Vincent Ng, and Latifur Rahman Khan, "*Weakly Supervised Cause Identification from Aviation Reports via Semantic Lexicon Construction*", submitted to Journal of Artificial Intelligence Research, Feb. 2009.
- Isaac Persing and Vincent Ng, "*Semi-Supervised Cause Identification from Aviation Safety Reports*", in Proc. of the Joint Conf. of the 47th Annual Meeting of the Association for Computational Linguistics and 4th Int. Joint Conf. on Natural Language Processing of the Asian Federation of Natural Language Processing (ACL-IJCNLP 2009), Singapore, Aug. 2009.
- Clay Woolam, Latifur Khan, "Multi-concept Document Classification Using a Perceptron-Like Algorithm," 2008 IEEE / WIC / ACM Int. Conf. on Web Intelligence (WI'08), Dec. 2008, Sydney, NSW, Australia, Page: 570-574.
- Clay Woolam, Latifur Khan, "*Multi-label large margin hierarchical perceptron*," Int. Journal of Data Mining, Modelling and Management (IJDMMM), 1(1): 5-22 (2008), Interscience Publisher (invited paper).
- Mohammad M Masud, Jing Gao, Latifur Khan, Jiawei Han, and Bhavani Thuraisingham, "*A Multi-Partition Multi-Chunk Ensemble Technique to Classify Concept-Drifting Data Streams*", Proc. 2009 Pacific-Asia Conf. on Knowledge Discovery and Data Mining (PAKDD'09), Bangkok, Thailand, Apr. 2009.
- Mohammad Masud, Jing Gao, Latifur Khan, Jiawei Han, and Bhavani Thuraisingham, "A *Practical Approach to Classify Evolving Data Streams: Training with Limited Amount of Labeled Data*", Proc. 2008 Int. Conf. on Data Mining (ICDM'08), Pisa, Italy, Dec. 2008.