

A Comparison of Approaches to Large-Scale Data Analysis

Magnus Kirø

Norwegian University of Science and Technology

October 25, 2012

Presentation Goals

The purpose of paper is to consider MapReduce and a regular Database Management Systems for large-scale data analysis.

- 1 **Parallel DBMS** and **MR**, two approaches to large-scale data analysis.
- 2 Thee **Architectural Elements** of DBMS and MR.
- 3 Outlinei of **Benchmark Results** and **Best Practice** for large-scale data analysis.

1 Two Approaches

- DBMS
- Map Reduce

2 Architectural Elements

- Schema Support
- Indexing
- Programming Model
- Data Distribution
- Execution Strategy
- Flexibility
- Fault Tolerance

3 Results

- Benchmark Results
- Best Practice

Database Management System

Title: **Tracking the flu pandemic by monitoring the Social Web**
Authors: V. Lampos and N. Cristianini
Submitted to: IAPR Cognitive Information Processing 2010 (accepted)

- Twitter and Health Protection Agency data for weeks 26-49, 2009 (on average 160,000 tweets collected per day geolocated in 54 urban centres in the UK)
- Frequency of **41 flu related words** (markers) in Twitter corpus had a correlation of **>80%** with the HPA flu rates in all UK regions
- Learn a better list of weighted markers **automatically**:
 - Generate a list of candidate markers (1560 words taken from flu related web pages)
 - Use **LASSO** for feature selection

DBMS -

Validation schemes:

- 1 Train on one region, validate regularisation parameter on another, test on the remaining regions (for all possible combinations)

Train/Validate (regions)	A	B	C	D	E
A	-	0.9594	0.9375	0.9348	0.9297
B	0.9455	-	0.9476	0.9267	0.9003
C	0.9154	0.9513	-	0.8188	0.908
D	0.9463	0.9459	0.9424	-	0.9337
E	0.8798	0.9506	0.9455	0.8935	-
Total Avg.					0.9256

97 selected words: lung, unwell, temperature, like, headache, season, unusual, chronic, child, dai, appetite, stay, symptom, spread, diarrhoea, start, muscle, weaken, immune, feel, liver, plenty, antiviral, follow, sore, people, nation, small, pandemic, pregnant, thermometer, bed, loss, heart, mention, condition, ...

- 2 Aggregate data from all regions, test on weeks 28 and 41 (2009) and train using the rest of the data set

DBMS - hei

- 1 Inferred vs Official flu rate in North England

figures/Lasso_Inference_regionC_1

- 2 Inferred vs Official rates in all regions (aggregated data set)

figures/Lasso_Inference_Aggregated

MR -

Title: **Flu detector - Tracking epidemics on Twitter**
Authors: V. Lampos, T. De Bie, and N. Cristianini
Submitted to: ECML PKDD 2010 Demos (under review)

- Extending and making more robust the methodology of P1
- Larger data sets (bigger time series) and more (2675) candidate features
- Select a list of features (markers) using BoLASSO (bootstrap version of LASSO)
- Then learn weights of those markers via linear least squares regression
- Stricter evaluation of the methodology - **Available online**
- Put all this into practice and come up with the **Flu detector**

Schema Support

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

Indexing

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Programming Model

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Data Distribution

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Execution Strategy

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

Flexibility

- ① ... (content omitted)
- ② ... (content omitted)
- ③ ... (content omitted)
- ④ ... (content omitted)

Fault Tolerance

- ❶ ... (content omitted)
- ❷ ... (content omitted)
- ❸ ... (content omitted)
- ❹ ... (content omitted)

A more time specific tentative plan

A more time specific tentative plan

This is the last slide.

Any questions?