# Extraction of Quantitative information from economic news

Utkarsh Upadhyay
**Mentor:**
Prof. R.M.K. Sinha

Indian Institute of Technology, Kanpur

April 14, 2009

# Outline

Introduction

The Solution

Meanwhile . . .

The Problem

Current state

Further work

# What was the problem, again?

Determine in real time from economic news:

1. Turn Over
2. Profit/Loss
3. Revenues

Can be seen as a very restricted case of a very general problem.

# The exact Problem

Utkarsh Upadhyay
**Mentor:**
Prof. R.M.K. Sinha

# Outline

Introduction

## The Solution

Meanwhile . . .

The Problem

Current state

Further work
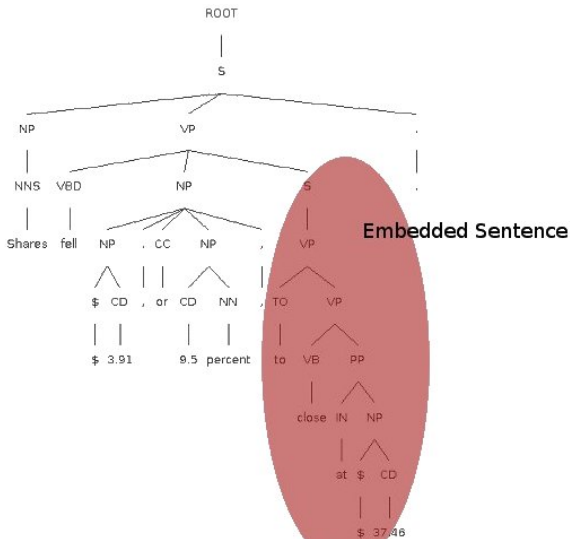
# Where's the news?

Utkarsh Upadhyay
**Mentor:**
Prof. R.M.K. Sinha

The news is fetched from the web:

1. *When?*: Updates real time RSS sources
2. Yahoo Finance News
3. *What form?*: Human readable form, HTML
   - Manually designed text filter
   - However, (17%) articles have tables

# Keeping it small

Breaking a sentence into simplest form.

# The right path

| Frequency | Path | Usual Description |
|-----------|------|-------------------|
| 14.2% | $VB \uparrow VP \downarrow PP$ | PP argument/adjunct |
| 11.8% | $VB \uparrow VP \uparrow S \downarrow NP$ | Subject (ARG0/ARG1) |
| 10.1% | $VB \uparrow VP \uparrow NP$ | Object (ARG0/ARG1) |
| . . . | | |

Table: Path from the predicate to the frame

# Read/Tag/Repeat

Stanford Manual Annotator?

Figure: How to do it right.

# What we get

```
<useful>
        <keyword>Shares</keyword>
        <direction>
                <predicate>fell </predicate>
        </direction>
        <relative_amount>
                <absolute_amount>$ 3.91</absolute_amount>
                , or
                <percent>9.5 percent</percent> ,
        </relative_amount>
</useful>
```

# Putting them together

- Tregex
- Tsurgeon

Figure: The whole sentence

# The path

Figure: A path from the *predicate* to the tag (*rel amt*)

# Outline

Introduction

The Solution

Meanwhile . . .

The Problem

Current state

Further work

# IB trading Olympiad

1. Direct Connection to the Internet
2. Ptolemy
3. Sectoral Arbitrage strategy

Paper in Advanced Data Analysis, Business Analytics and Intelligence conference in IIM-A.

# Outline

Utkarsh Upadhyay
Mentor:
Prof. R.M.K. Sinha

# The Sparsity

Very sparse data:

## Example

For the relation:be -> absolute_amount
 VBD(up)VP(down)NP(down)NP: 2
 VBD(up)VP(down)PP(down)NP(down)PP(down)NP: 1

Adding more features?

# Changing stance

- ▶ Some tags discarded
- ▶ Dictionary based name extractor
- ▶ RegEx based time extractor

# Outline

Introduction

The Solution

Meanwhile . . .

The Problem

Current state

Further work

# Useless news

Figure: No useless news

# Unparsed news

Figure: No unparsed news

# Absolute amount

Figure: Absolute Amount

# Relative amount

Figure: Relative Amount

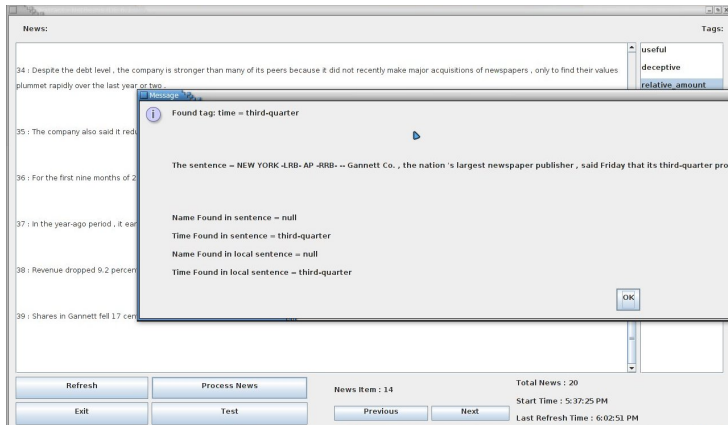# Not as helpless

Utkarsh Upadhyay
**Mentor:**
Prof. R.M.K. Sinha

Figure: Same time references

# Maintainers are psychopathic killers

- NetBeans
- GUI
- JavaDoc
- JUnit testing

Figure: Dependency Tree of the project

# Outline

Utkarsh Upadhyay
**Mentor:**
Prof. R.M.K. Sinha

# Next rung. . .

1. More data
2. Fixing temporal references, Grammatical dependency
3. HTML text extraction
4. Knowledgebase
5. Finer classification, discarded information

# Thank you

Utkarsh Upadhyay
utkarshu@iitk.ac.in