# UNSUPERVISED EMAIL LABELLING

by:

Dibya Ranjan (10243)

Vishal Kumavat(10818)

Advisor:
**Prof. Amit Mukerjee**

**CS365- Artificial Intelligence**

**Dept. of Computer Science & Engineering**

**IIT Kanpur**

# IMPORTANCE OF PROBLEM

- Emails have become the basic part of one's life.

- One practically receives 50-60 emails per week.

- Most of the times it become difficult to manage these emails and one may sometimes cannot find

# PAST WORKS & APPLICATIONS

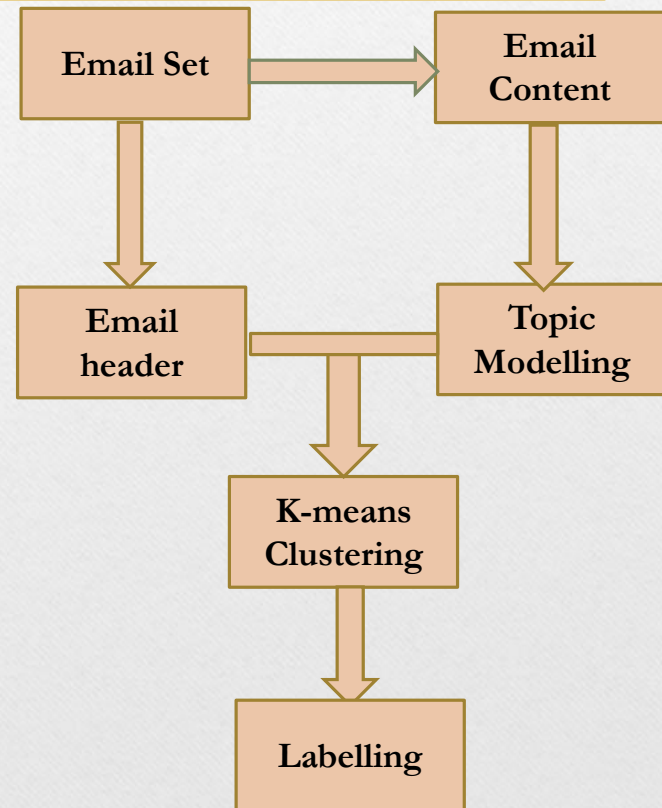- **Real-Time Topic Modelling of Microblogs**

    — *by Yogesh Tewari and Rajesh Kawad*

- **Topic Modelling the Sarah Palin Emails**

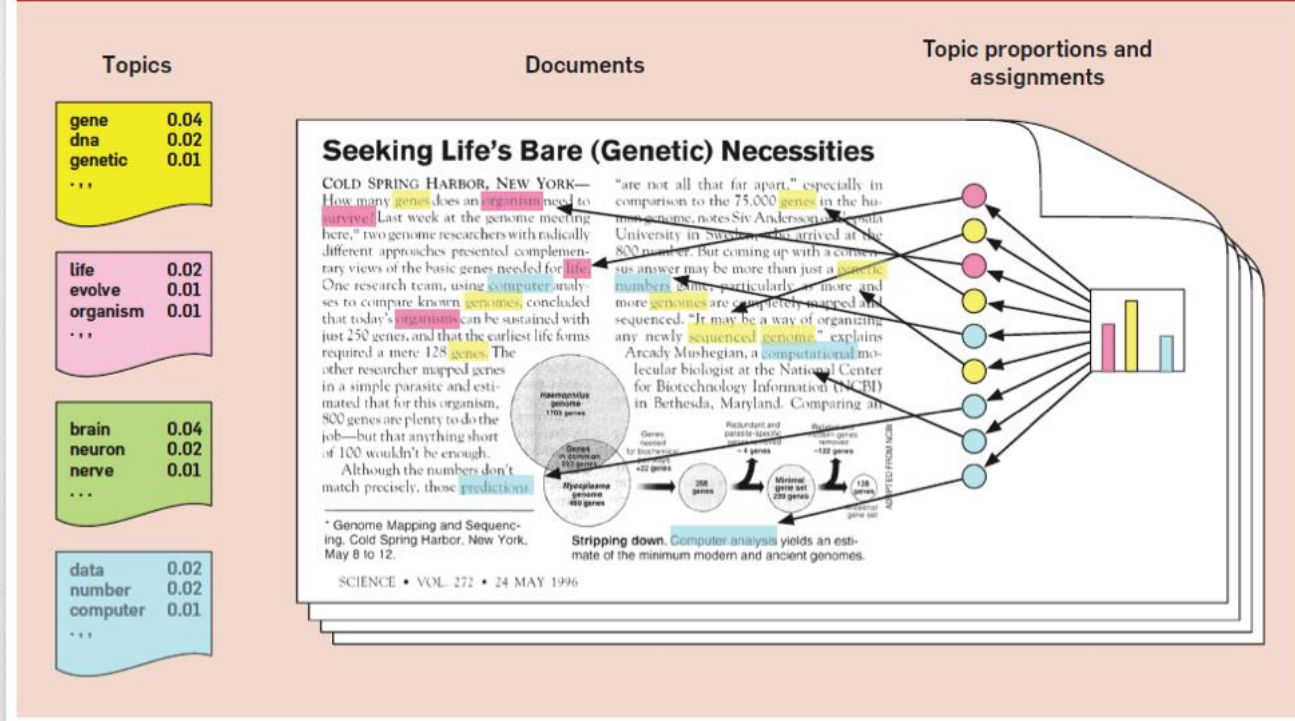    — *from Edwin Chen's blog*

# OUR APPROACH

1. Email is divided in two parts
   - Email header
   - Email content

2. Then topic modelling on email content is applied

3. Combining the features of header and from topic modelling we apply K mean clustering
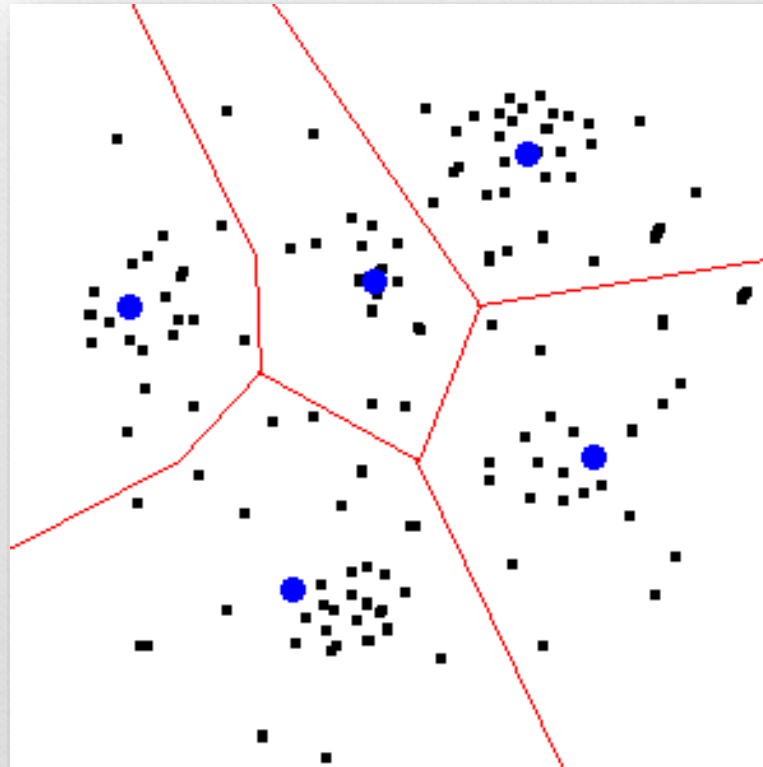
# Algorithms
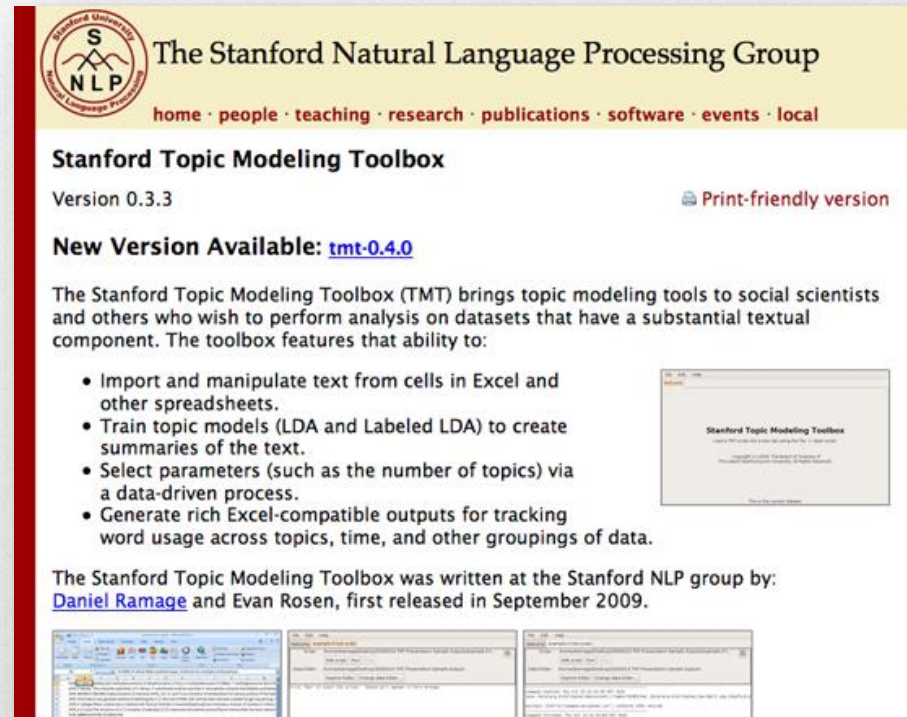
## Topic Modelling

# Algorithms

## K Means Clustering

# TOOLS USED

- Topic Modelling Tool
  - ✓ It divides documents under different topics
  - ✓ It uses LDA(Latent Dirichlet Allocaton)



*From http://nlp.stanford.edu/software/tmt/tmt-0.4/*

# RESULTS

**List of Topics**

1. system average equipartition theorem law energy number kinetic needham water
2. south hindi film acting sullivan edward back time naa award
3. years yard national wilderness war parks park modern survived grossing
4. sunderland echo zinta role paper world earned debut films independent
5. rings ring dust uranus thespis moons narrow uranian addition dark
6. confederate london indian century ho female filmfare service thylacinus gods
7. battle union hawes kentucky army grant gen tennessee united confederates
8. gunnhild australian norway numerous england career death king life particles
9. thylacine tasmanian tiger general mother acted male devil species related
10. test including cricket hill actress gilbert record top movement actors

# RESULTS

**TOPIC** : system average equipartition theorem law energy number kinetic needham water ...

top-ranked docs in this topic (#words in doc assigned to this topic)

2.  (75) equipartition_theorem.txt
3.  (13) uranus.txt
4.  (7)  thylacine.txt
5.  (7)  elizabeth_needham.txt
6.  (6)  sunderland_echo.txt
7.  (5)  zinta.txt
8.  (5)  thespis.txt
9.  (5)  shiloh.txt
10. (5)  gunnhild.txt
11. (4)  hill.txt
12. (4)  hawes.txt
13. (2)  yard.txt

# RESULTS

**DOC** :uranus.txt

```
The rings of Uranus were discovered on March 10,
1977, by James L. Elliot, Edward W. Dunham, and
Douglas J. Mink. Two additional rings were discovered
in 1986 by the Voyager 2 spacecraft, and two outer
rings were found in 2003â€"2005 by the Hubble Space
Telescope. A number of faint dust bands and
incomplete arcs may exist between the main rings. The
rings are extremely darkâ€"the Bond albedo of the
rings' particles does not exceed 2%. They are likely
composed of water ice with the addition of so
```

Top topics in this doc (% words in doc assigned to this topic)

(50%) rings ring dust uranus thespis moons narrow uranian addition dark ...

(12%) system average equipartition theorem law energy number kinetic needham water ...

(10%) years yard national wilderness war parks park modern survived grossing ...

(7%) sunderland echo zinta role paper world earned debut films independent ...

(6%) gunnhild australian norway numerous england career death king life particles ...

# DATASET

We have used the topic modelling on online available dataset on around 2000 emails

Created the data set of around 400-500 emails from personal inbox and used them for the unsupervised labelling

.

# IMPROVEMENTS

We can also use more header information like subject and whether the email was sent to single or multiple receiver.

# CONCLUSION

- Classified the emails based on balance between email header information and email semantics information.
- A Parameter was found which give priority to email header or email semantics

# REFERENCES

- Mtech thesis on Email Classification Ozcaglar, Cagri. (2008)

- Topic Modelling theory-

  - http://clc.yale.edu/2011/10/07/how-to-do-your-own-topic-modeling/

  - http://www.fredgibbs.net/clio3workspace/blog/topic-modeling/

  - http://miriamposner.com/blog/?p=1335

  - http://blog.echen.me/2011/06/27/topic-modeling-the-sarah-palin-emails/


- Topic Modelling Tool-

  - http://nlp.stanford.edu/software/tmt/tmt-0.4/

- Dataset-

  - http://dhhumanist.org/Archives/Current/