

# Coding open answers

Peter Lugtig      [p.lugtig@uu.nl](mailto:p.lugtig@uu.nl)

# Text data

- Ubiquitous in surveys
  - But rarely analyzed
- Even more so in big data
  - Social Media, speech (=text)
- What to do?
  - Historically: manually code them
  - Now/future: machine learning or combination

# Coding open answers

Important to keep in mind how you will use the data.

- Are detailed codes necessary?
  - Population size, time, importance of question
- What was goal of the survey?
  - e.g. if it was to evaluate quality of teaching, focus on this.
- Are there existing coding schemes?
  - e.g. occupational, educational coding

# How to code open answers manually

- Open
  - Break down, compare, and categorize data.
- Axial
  - Make connections between categories after open coding.
- Selective
  - Select the core category, relate it to other categories and confirm and explain those relationships.
  - See Boeiye (2009) for more details

# Example of coding of open answers

- “What is the most important problem facing our country?”
  1. “abortion”
  2. “economy, welfare should be checked into better, choices of how tax money is spent.”
  3. “honey, i really couldn’t tell you.”
  4. “moral problem – general moral decay; drugs; crime; poor education; divorce; alcohol.”
  5. “our national debt”
  6. “our physical responsible and the rest will fall into place”
  7. “taxes, eats up your income”
  8. “the federal bonus they keep promising to give to wwii veterans”
  9. “the life amendment”

Example from ANES 1984-2000; Goodrich 2008

# Example of coding of open answers (2)

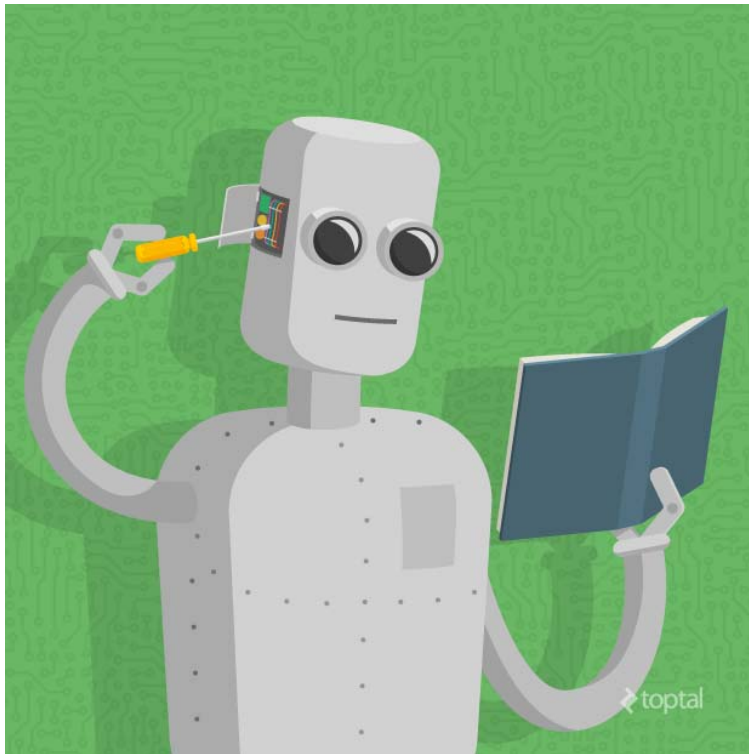
- |                           |                            |                        |                           |                    |                                 |
|---------------------------|----------------------------|------------------------|---------------------------|--------------------|---------------------------------|
| • Abortion                | • Children (not schools)   | • Education            | • Healthcare              | • Middle East      | • Seniors (not Social Security) |
| • AIDS/STDS               | • Communism                | • Employment           | • Housing                 | • Morality         | • Sexual Preference             |
| • <b>Arms/Weapons</b>     | • Cost of Living           | • Energy               | • <b>Hunger</b>           | • Native Americans | • <b>Social Security</b>        |
| • Big Government          | • <b>Crime (not drugs)</b> | • Environment          | • Immigration             | • Nicaragua        | • Somalia                       |
| • <b>Budget</b>           | • Defense/Military         | • Equality             | • Impeachment             | • <b>Nuclear</b>   | • South America                 |
| • Business                | • Discrimination/Race      | • Families (not value) | • Imports/Exports         | • Oil/Gas          | • Space Program                 |
| • Campaign Finance Reform | • <b>Drugs</b>             | • Far East             | • Insurance               | • Overpopulation   | • Star Wars                     |
| • Central America         | • Drunk Driving            | • Farming              | • International Relations | • Persian Gulf     | • <b>Taxes</b>                  |
|                           | • <b>Economy</b>           | • Guns                 | • Iran                    | • Pornography      | • Teen Pregnancy                |
|                           |                            |                        | • Iraq                    | • <b>Poverty</b>   | • Terrorism                     |
|                           |                            |                        | • Japan                   | • President        | • Unions                        |
|                           |                            |                        | • Justice System          | • Puerto Rico      | • Wages                         |
|                           |                            |                        | • Kuwait                  | • Religion         | • <b>War/Peace</b>              |
|                           |                            |                        | • Media                   | • <b>Russia</b>    | • Wealth Inequality             |
|                           |                            |                        | • Medicare/Medicade       | • Saudi Arabia     | • <b>Welfare</b>                |

Codes in bold: among 9 most frequently codes in any election between 1984 and 2000

Example from ANES 1984-2000; Goodrich 2008

# Coding open answers automatically

- Large texts are impossible to code by hand
- Answer: Machine learning



# Machine learning in coding text

- Large subfield in IT
  - Do counts of words
  - Sentiment analysis

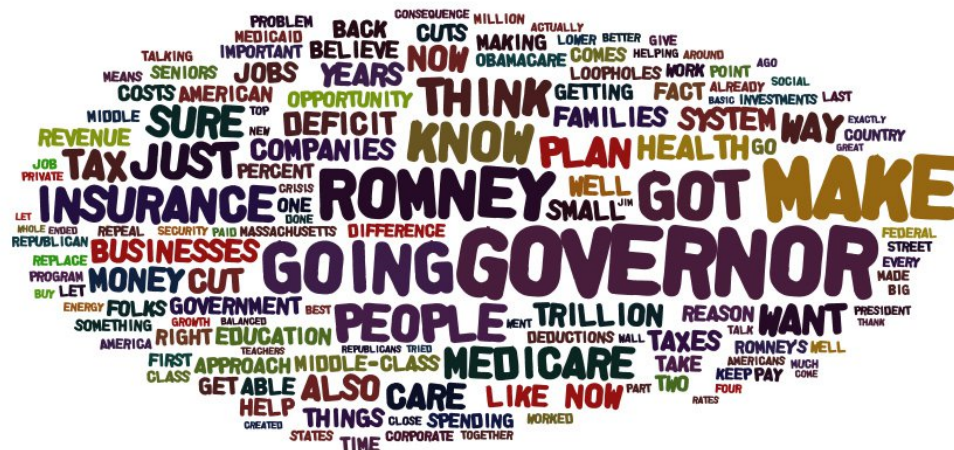




# Machine learning in coding text

- Sentiment analysis
  - Associations with topics
  - Using the twitter 'hose'
  - U.S. election presidential Debate 2012

# Barack Obama



- Sentiment analysis
  - Associations with topics
  - Using the twitter 'hose'
  - U.S. election presidential Debate 2012

[illegible][illegible]

# Machine learning in coding text

- Large subfield in IT
  - Do counts of words
  - Frequencies of text fragments -> group these
    - Sentiment analysis

## Twitter Political Index: The Swing States

Net difference in @BarackObama vs. @MittRomney #Twindex scores in @USAToday- and @Gallup-defined "swing states."  
October 16, 2012 to November 6, 2012 (partial day)

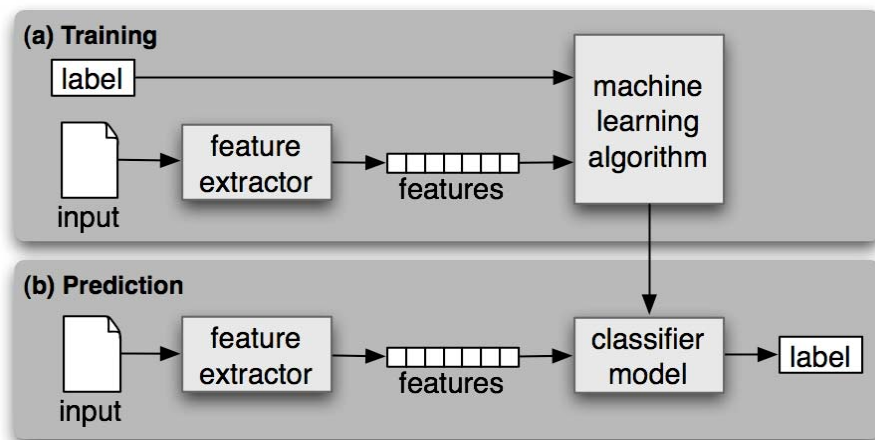
Follow @gov for more about government & politics on Twitter.

ANALYTICS BY TOPSY



# Machine learning in coding text

- Sentiment analysis
  - How to do it
  - Supervised learning using a training set
  - Support Vector Machines
    - See Bird, Klein and Loper (2014)
    - Accuracy with human coding: ~90%



# Machine learning

- Sentiment analysis not very sophisticated
- Other approaches relevant for survey research
  - Structural Topic Models (Roberts et al 2012;2014)

# Structural Topic Model

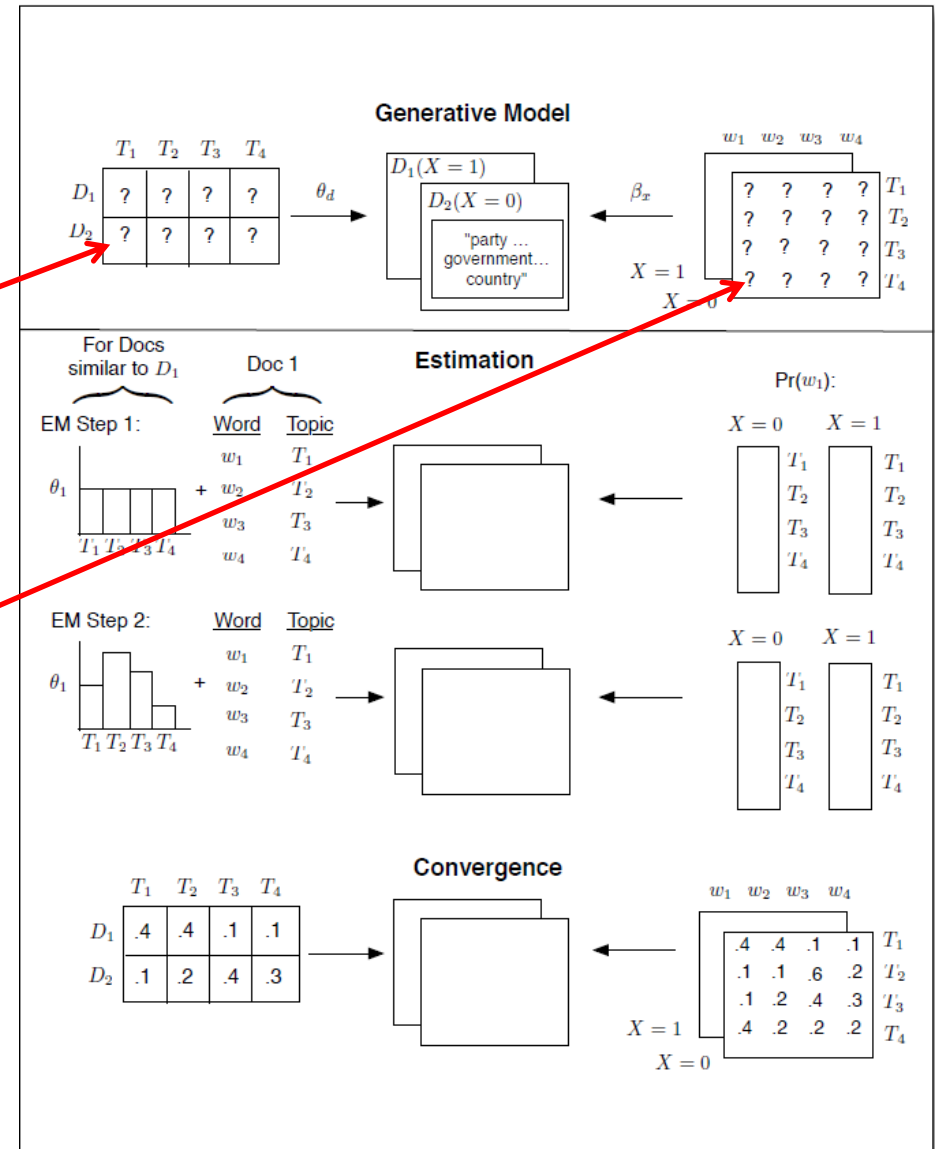
- Mixture models (as in statistics)
  - LDA (Latent Dirichlet Analysis)
- 3-level mixture model
  - 1. Document (case) is mixture of topics
  - 2. Topic is mixture of words
  - 3. words
- Between the levels a log-normal link function (multinomial)

# Structural Topic Model

- Heuristic model

- from Roberts et al 2016

- Generation:
- (Gibbs sampler)
- Multiple Topics (T) within documents(D)
- Multiple Words (w) within topics (T)

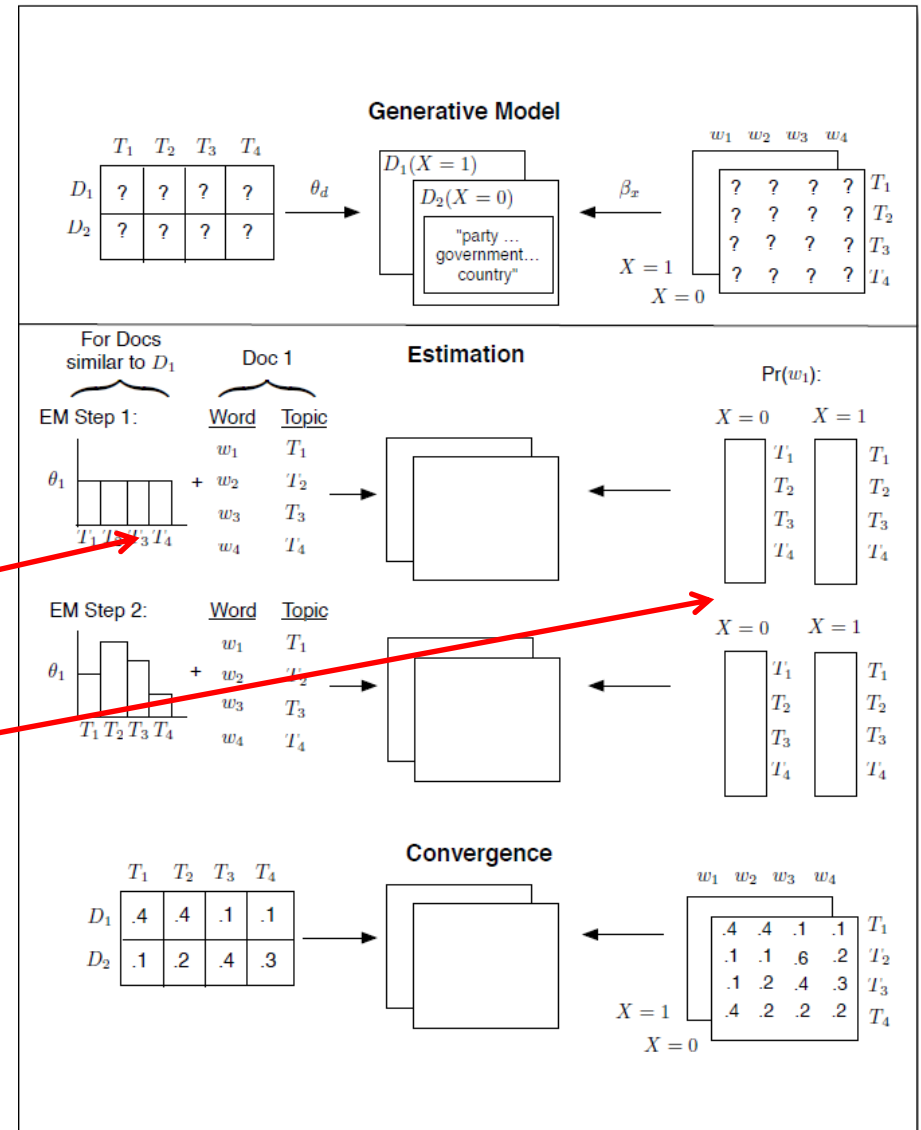


# Structural Topic Model (2)

- Heuristic model

- from Roberts et al 2016

- Estimation:
- Per document (case):
- Words and topics form document
- How likely is every word related to topic





# Structural Topic Model (3)

- Mixture models (as in statistics)
  - LDA (Latent Dirichlet Analysis)
- Can be linked to covariates
  - Do proper statistical analysis
  - E.g. do males use specific words or word-clusters more often?

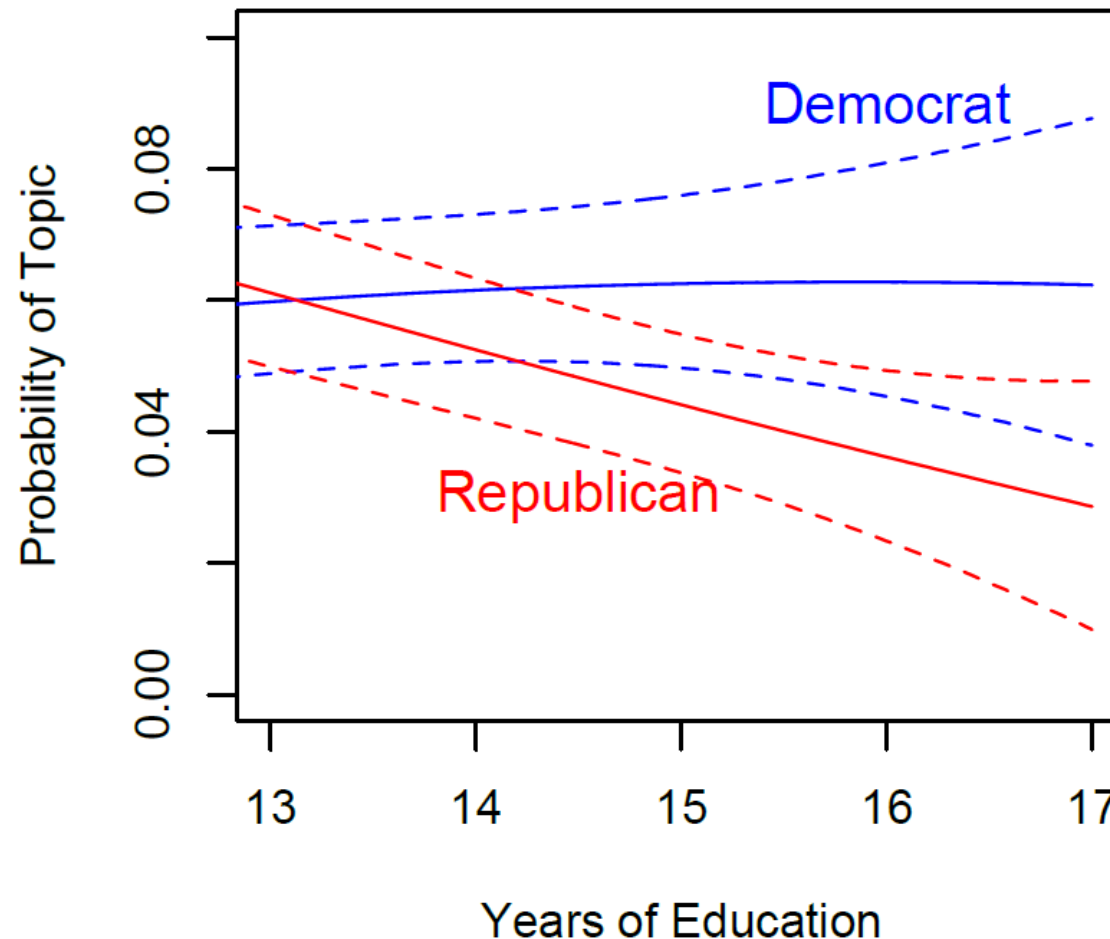
# The STM package in R - output

iraq , war  
war  
donotknow  
job  
budget , peopl  
monei  
terror  
govern  
deficit  
none  
unemploy  
crisi  
countri  
affair , foreign  
think , econom  
view , world , around  
financi  
bipartisan , person , problem  
financ , keep , nationalsecur  
get  
debt , parti  
foreign , polici  
bail , out , divis  
go  
need  
togeth , work  
relat , nation , defens  
anyth  
honesti , safeti  
tax , thing  
militari , politician  
immigr , sure  
right , now  
noth , done  
everyon , bodi , organ  
healthcar  
terrorist  
republican , democrat , troop  
market , stock  
busi , liber , mean  
probabl , give , futur  
obama , barack , control  
good , along , get  
lose , peopl , home  
guess , new , polit  
everyth  
global , environ , imag  
american , unit , back  
corrupt

eci

# The STM package in R - output

## STM War Topic and Education



# Comparison with manual coding

- |  |   |   |  |   |   |
|--|---|---|--|---|---|
| <ul style="list-style-type: none"><li>• Abortion</li><li>• AIDS/STDS</li><li>• <b>Arms/Weapons</b></li><li>• Big Government</li><li>• <b>Budget</b></li><li>• Business</li><li>• Campaign Finance Reform</li><li>• Central America</li></ul> | <ul style="list-style-type: none"><li>• <b>Children (not schools)</b></li><li>• Communism</li><li>• Cost of Living</li><li>• <b>Crime (not drugs)</b></li><li>• <b>Defense/Military</b></li><li>• Discrimination/Race</li><li>• <b>Drugs</b></li><li>• Drunk Driving</li><li>• <b>Economy</b></li></ul> | <ul style="list-style-type: none"><li>• <b>Education</b></li><li>• <b>Employment</b></li><li>• Energy</li><li>• <b>Environment</b></li><li>• Equality</li><li>• Families (not values)</li><li>• Far East</li><li>• Farming</li><li>• Guns</li></ul> | <ul style="list-style-type: none"><li>• <b>Healthcare</b></li><li>• <b>Housing</b></li><li>• <b>Hunger</b></li><li>• Immigration</li><li>• Impeachment</li><li>• Imports/Exports</li><li>• Insurance</li><li>• International Relations</li><li>• Iran</li><li>• Iraq</li><li>• Japan</li><li>• Justice System</li><li>• Kuwait</li><li>• Media</li><li>• Medicare/Medicade</li></ul> | <ul style="list-style-type: none"><li>• <b>Middle East</b></li><li>• <b>Morality</b></li><li>• Native Americans</li><li>• Nicaragua</li><li>• <b>Nuclear</b></li><li>• Oil/Gas</li><li>• Overpopulation</li><li>• Persian Gulf</li><li>• Pornography</li><li>• <b>Poverty</b></li><li>• President</li><li>• Puerto Rico</li><li>• Religion</li><li>• <b>Russia</b></li><li>• Saudi Arabia</li></ul> | <ul style="list-style-type: none"><li>• Seniors (not Social Security)</li><li>• Sexual Preference</li><li>• <b>Social Security</b></li><li>• Somalia</li><li>• South America</li><li>• Space Program</li><li>• Star Wars</li><li>• <b>Taxes</b></li><li>• Teen Pregnancy</li><li>• Terrorism</li><li>• Unions</li><li>• Wages</li><li>• <b>War/Peace</b></li><li>• Wealth Inequality</li><li>• <b>Welfare</b></li></ul> |
|--|---|---|--|---|---|

# STM analysis: example

- ▶ Conclusions:
  - ▶ Manual
  - ▶ Supervised (using training data – not worked out today)
  - ▶ Unsupervised (STM approach)
- ▶ Manual and Automatic coding yield similar data
  - ▶ 94% of raw responses yield the same code (Roberts et al 2016)
  - ▶ In a lot less time
  - ▶ Automatic coding yields other benefits
    - ▶ Statistics on classification quality etc.
    - ▶ Statistical models

# References

- Boeije (2009) Analysis in qualitative research. Sage: London.
- Bird, S., Klein, E. & Loper, E. (2014) the natural Language Processing with Python. <http://www.nltk.org/book/>
- Roberts, M.E., Stewart, B.D., Tingley, D., Lucas, C., Leder-Luis, J., Kushner gadarian, S., Albertson, B., Rand, D.G. (2014) Structural topic models for open ended survey responses. American Journal of Political Science, 58. 1064-1082. DOI:10.1111/ajps.12103
- Roberts, M.E., Stewart, B.M., & Tingley, D. (forthcoming) stm: R package for Structural Topic Models. Journal of Statistical Software
- Goodrich (2008) A coding methodology for open-ended survey questions. Paper for the 2008 new faces in political methodology conference.