

This class is great, take it.

DS-GA 1015, Text as Data  
Arthur Spirling

Feb 2, 2021

# Professor and Lectures



# Professor and Lectures



Prof Arthur Spirling

# Professor and Lectures



Prof Arthur Spirling

`arthur.spirling@nyu.edu`

# Professor and Lectures



Prof Arthur Spirling

`arthur.spirling@nyu.edu`

`https://nyu.zoom.us/j/6678696568`

# Professor and Lectures



Prof Arthur Spirling

`arthur.spirling@nyu.edu`

`https://nyu.zoom.us/j/6678696568`

OH Wednesday, 1030-1130AM.

# Professor and Lectures



Prof Arthur Spirling

[arthur.spirling@nyu.edu](mailto:arthur.spirling@nyu.edu)

<https://nyu.zoom.us/j/6678696568>

OH Wednesday, 1030-1130AM.

Flip Tuesdays 11AM-1240PM, 101, 19W4.

# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material.



# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material. See < Lecture > tab.

# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material. See < Lecture > tab.  
And < Resources > tab for slides.

# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material. See < Lecture > tab.  
And < Resources > tab for slides.

~ 30 minutes of “in-person” flipped material in room 101 of 19W4th Street. This will be recorded and livecast.

# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material. See < Lecture > tab.  
And < Resources > tab for slides.

~ 30 minutes of “in-person” flipped material in room 101 of 19W4th Street. This will be recorded and livecast. Use professor’s personal zoom link.

# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material. See < Lecture > tab.  
And < Resources > tab for slides.

~ 30 minutes of “in-person” flipped material in room 101 of 19W4th Street. This will be recorded and livecast. Use professor’s personal zoom link.

Note your **cohort**: Cohort *A* next week (Feb 9); Cohort *B* week after (Feb 16), then alternate.

# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material. See < Lecture > tab.  
And < Resources > tab for slides.

~ 30 minutes of “in-person” flipped material in room 101 of 19W4th Street. This will be recorded and livecast. Use professor’s personal zoom link.

Note your **cohort**: Cohort *A* next week (Feb 9); Cohort *B* week after (Feb 16), then alternate. Bring your phone and record your seat the first time you come.

# Structure of “Lecture”

~ 70 mins of pre-recorded lecture material. See < Lecture > tab.  
And < Resources > tab for slides.

~ 30 minutes of “in-person” flipped material in room 101 of 19W4th Street. This will be recorded and livecast. Use professor’s personal zoom link.

Note your **cohort**: Cohort *A* next week (Feb 9); Cohort *B* week after (Feb 16), then alternate. Bring your phone and record your seat the first time you come.

flipped structure is **subject to demand**: if turnout is low, we will revert to ‘live’ but recorded lectures every week.

# TA and Sections







Ms Lucia Motolinia

# TA and Sections



Ms Lucia Motolinia

lucia.motolinia@nyu.edu

# TA and Sections



Ms Lucia Motolinia

lucia.motolinia@nyu.edu

OH Friday, 10–11AM

# TA and Sections



Ms Lucia Motolinia

lucia.motolinia@nyu.edu

OH Friday, 10–11AM

Sec Thursday, 2–250PM (remote)

→ (start this week!)

# What this class is about...

What this class is about...



new  
no  
people  
need  
research

## Text as the new frontier of...

new  
no  
people  
need  
research

Text as the new frontier of...  
data:



# What this class is about...



## Text as the new frontier of...

[data](#): lots of it (literally petabytes) on the web... not to mention archives.

# What this class is about...



## Text as the new frontier of...

**data:** lots of it (literally petabytes) on the web... not to mention archives.

**methods:**

# What this class is about...



## Text as the new frontier of...

**data:** lots of it (literally petabytes) on the web... not to mention archives.

**methods:** unstructured data needs to be harvested and modeled.

# What this class is about...

## Text as the new frontier of...

**data:** lots of it (literally petabytes) on the web...not to mention archives.

**methods:** unstructured data needs to be harvested and modeled.

social science:



# What this class is about...



## Text as the new frontier of...

**data:** lots of it (literally petabytes) on the web... not to mention archives.

**methods:** unstructured data needs to be harvested and modeled.

**social science:** politicians give speeches, thinkers write articles,

new  
ent  
ne  
peo  
ke  
need  
ext  
put  
form  
research

## Text as the new frontier of...

**data:** lots of it (literally petabytes) on the web...not to mention archives.

**methods:** unstructured data needs to be harvested and modeled.

**social science:** politicians give speeches, thinkers write articles, nations sign treaties,

# What this class is about...



## Text as the new frontier of...

**data:** lots of it (literally petabytes) on the web... not to mention archives.

**methods:** unstructured data needs to be harvested and modeled.

**social science:** politicians give speeches, thinkers write articles, nations sign treaties, users connect on Facebook etc.



# What this class is about...



## Text as the new frontier of...

**data:** lots of it (literally petabytes) on the web... not to mention archives.

**methods:** unstructured data needs to be harvested and modeled.

**social science:** politicians give speeches, thinkers write articles, nations sign treaties, users connect on Facebook etc.

Introduction to quantitative 'text-as-data' approaches as strategies to learn more about social scientific phenomena of interest.

# Overview

# Overview

new  
ent  
ne  
peo  
ke  
need  
ext  
put  
form  
research  
law together  
matter republic  
south  
us  
sc

since  
stand  
responsibility  
scho  
parents t  
live SUCC  
research

# Overview



- Descriptive inference:

# Overview



- Descriptive inference: how to characterize text,

# Overview



- **Descriptive inference:** how to characterize text, vector space model,



# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words,



# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures,

# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity,

# Overview

- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity,

# Overview

- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity, style.

new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext  
ne  
feet  
us  
sc  
s  
live  
succ  
law  
matter  
parents  
stand  
responsibility  
since  
sma  
race  
g  
scho

- A set of small navigation icons typically found in Beamer presentations, including symbols for back, forward, search, and other slide controls.

new  
no  
people  
need  
research

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext  
ne  
feet  
us  
sc  
s  
live  
succ  
law  
matter  
parents  
stand  
responsibility  
since  
scho

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext  
ne  
feet  
us  
sc  
s  
live  
succ  
law  
matter  
parents  
stand  
responsibility  
since  
sma  
race  
sch

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡



new  
no  
people  
need  
research

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

new  
no  
people  
need  
research

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ≡ ▶ ◀ ≡ ≡ ▶ ≡ ≡ ≡ ↺ 🔍 ↻

# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity, style.
- **Supervised techniques:** dictionaries, classification, scaling, machine learning approaches.
- **Unsupervised techniques:** clustering,

new  
no  
people  
need  
research

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

new  
no  
people  
need  
research

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity, style.
- **Supervised techniques:** dictionaries, classification, scaling, machine learning approaches.
- **Unsupervised techniques:** clustering, principal components, scaling, topic models,

# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity, style.
- **Supervised techniques:** dictionaries, classification, scaling, machine learning approaches.
- **Unsupervised techniques:** clustering, principal components, scaling, topic models, embeddings.

new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻



# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity, style.
- **Supervised techniques:** dictionaries, classification, scaling, machine learning approaches.
- **Unsupervised techniques:** clustering, principal components, scaling, topic models, embeddings.
- **Special topics:** modeling debate and questions,

# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity, style.
- **Supervised techniques:** dictionaries, classification, scaling, machine learning approaches.
- **Unsupervised techniques:** clustering, principal components, scaling, topic models, embeddings.
- **Special topics:** modeling debate and questions, bursts,

new  
no  
people  
need  
research

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

# Overview



- **Descriptive inference:** how to characterize text, vector space model, collocations, bag-of-words, dissimilarity measures, diversity, complexity, style.
- **Supervised techniques:** dictionaries, classification, scaling, machine learning approaches.
- **Unsupervised techniques:** clustering, principal components, scaling, topic models, embeddings.
- **Special topics:** modeling debate and questions, bursts, networks, text reuse.

# Quantitative vs Qualitative

# Quantitative vs Qualitative



# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.

# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.
- Partly **still** is:



# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.
- Partly **still** is: need to make qualitative judgements about what the text reveals,

# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.
- Partly **still** is: need to make qualitative judgements about what the text reveals, and **validation** requires substantive knowledge

# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.
- Partly **still** is: need to make qualitative judgements about what the text reveals, and **validation** requires substantive knowledge
- 'Distant reading' instead of '**close reading**'

# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.
- Partly **still** is: need to make qualitative judgements about what the text reveals, and **validation** requires substantive knowledge
- 'Distant reading' instead of '**close reading**'—not focussed on **interpretation** in light of norms or belief systems.

# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.
- Partly **still** is: need to make qualitative judgements about what the text reveals, and **validation** requires substantive knowledge
- 'Distant reading' instead of '**close reading**'—not focussed on **interpretation** in light of norms or belief systems.
- Important: **quantitative** work is **reliable** and **replicable** (easily)

# Quantitative vs Qualitative



- For most of its history text analysis was **qualitative**.
- Partly **still** is: need to make qualitative judgements about what the text reveals, and **validation** requires substantive knowledge
- 'Distant reading' instead of '**close reading**'—not focussed on **interpretation** in light of norms or belief systems.
- Important: **quantitative** work is **reliable** and **replicable** (easily) and can cope with **large volume** of material.

# What this class is not about...

What this class is not about...

new  
ent  
ne  
peo  
ke  
need  
ext  
put  
form  
research  
law together  
matter republic  
south  
us  
sc



new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext  
ne  
feet  
us  
sc

- Data acquisition:

# What this class is not about...

- Data acquisition: many sources of text,



## What this class is not about...



- Data acquisition: many sources of text, but **web-scraping** better taught elsewhere.

# What this class is not about...



- Data acquisition: many sources of text, but [web-scraping](#) better taught elsewhere.

e.g. <http://www.crummy.com/software/BeautifulSoup/>

# What this class is not about...



- Data acquisition: many sources of text, but [web-scraping](#) better taught elsewhere.

e.g. <http://www.crummy.com/software/BeautifulSoup/>

We have many CDS people who can probably help with this...

# What this class is not about...



- Data acquisition: many sources of text, but **web-scraping** better taught elsewhere.

e.g. <http://www.crummy.com/software/BeautifulSoup/>

We have many CDS people who can probably help with this...

- Regular expressions and **basic text manipulation**:

# What this class is not about...



- Data acquisition: many sources of text, but **web-scraping** better taught elsewhere.

e.g. <http://www.crummy.com/software/BeautifulSoup/>

We have many CDS people who can probably help with this...

- Regular expressions and **basic text manipulation**: won't generally be required,

# What this class is not about...



- Data acquisition: many sources of text, but **web-scraping** better taught elsewhere.

e.g. <http://www.crummy.com/software/BeautifulSoup/>

We have many CDS people who can probably help with this...

- Regular expressions and **basic text manipulation**: won't generally be required, though helpful if known.



# What this class is not about...



- Data acquisition: many sources of text, but **web-scraping** better taught elsewhere.

e.g. <http://www.crummy.com/software/BeautifulSoup/>

We have many CDS people who can probably help with this...

- Regular expressions and **basic text manipulation**: won't generally be required, though helpful if known.
- **CS** 'stuff' like machine translation, OCR, algorithm design etc.

# What this class is not about...



- Data acquisition: many sources of text, but **web-scraping** better taught elsewhere.

e.g. <http://www.crummy.com/software/BeautifulSoup/>

We have many CDS people who can probably help with this...

- Regular expressions and **basic text manipulation**: won't generally be required, though helpful if known.

- **CS** 'stuff' like machine translation, OCR, algorithm design etc.

→ excellent options elsewhere.

# Requirements

# Requirements

new  
ent  
ne  
peo  
ke  
need  
ext  
pub  
form  
research  
law together  
matter republic  
south  
us  
sc  
live SUCC  
parents t  
responsibility  
stand  
since  
race  
sma  
scho

# Requirements



- A first class in statistics and/or inference

# Requirements



- A first class in statistics and/or inference
- Basic knowledge of calculus, probability, densities, distributions, statistical tests, hypothesis testing, the linear model, maximum likelihood and generalized linear models is assumed.

# Requirements



- A first class in statistics and/or inference
- Basic knowledge of calculus, probability, densities, distributions, statistical tests, hypothesis testing, the linear model, maximum likelihood and generalized linear models is assumed.
- Familiarity with core language and software environment of this course: R.

# Requirements

- A first class in statistics and/or inference
  - Basic knowledge of calculus, probability, densities, distributions, statistical tests, hypothesis testing, the linear model, maximum likelihood and generalized linear models is assumed.
  - Familiarity with core language and software environment of this course: R.
- check in with me if unsure.



# Terminology

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ .

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete**,

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete**, **uniform**

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
- **MLE** is obvious,

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
- **MLE** is obvious,  $\max(x) = m = \hat{N} = 68$ ,

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
- **MLE** is obvious,  $\max(x) = m = \hat{N} = 68$ , but **biased** (though **consistent**).



# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
- **MLE** is obvious,  $\max(x) = m = \hat{N} = 68$ , but **biased** (though **consistent**). Note that  $\max$  is **sufficient statistic** in this case.

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
- **MLE** is obvious,  $\max(x) = m = \hat{N} = 68$ , but **biased** (though **consistent**). Note that  $\max$  is **sufficient statistic** in this case.
- **Bayesian estimator** with uniform **prior** has **closed form** for the mean:  
$$E(N|x) = \frac{k-1}{k-2}(m-1).$$

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
- **MLE** is obvious,  $\max(x) = m = \hat{N} = 68$ , but **biased** (though **consistent**). Note that  $\max$  is **sufficient statistic** in this case.
- **Bayesian estimator** with uniform **prior** has **closed form** for the mean:  $E(N|x) = \frac{k-1}{k-2}(m-1)$ . Thus, (average)  $\hat{N} = 89$ .

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
  - **MLE** is obvious,  $\max(x) = m = \hat{N} = 68$ , but **biased** (though **consistent**). Note that  $\max$  is **sufficient statistic** in this case.
  - **Bayesian estimator** with uniform **prior** has **closed form** for the mean:  $E(N|x) = \frac{k-1}{k-2}(m-1)$ . Thus, (average)  $\hat{N} = 89$ .
- no need for **simulation**.

# Terminology

- **sample** is  $x = \{3, 4, 20, 50, 68\}$ . Problem requires calculation of maximum of **discrete, uniform distribution,  $\hat{N}$** .
  - **MLE** is obvious,  $\max(x) = m = \hat{N} = 68$ , but **biased** (though **consistent**). Note that  $\max$  is **sufficient statistic** in this case.
  - **Bayesian estimator** with uniform **prior** has **closed form** for the mean:  
 $E(N|x) = \frac{k-1}{k-2}(m-1)$ . Thus, (average)  $\hat{N} = 89$ .
- no need for **simulation**.
- Straightforward to implement via **function writing** in **R**.



sample, discrete, uniform, distribution,  $\hat{N}$ , MLE,  
biased, consistent, sufficient statistic, Bayesian,  
estimator, prior, closed form, simulation, function  
writing, R

sample, discrete, uniform, distribution,  $\hat{N}$ , MLE,  
biased, consistent, sufficient statistic, Bayesian,  
estimator, prior, closed form, simulation, function  
writing, R

If the terms in blue were familiar to you (even if you can't recall their exact meaning),



sample, discrete, uniform, distribution,  $\hat{N}$ , MLE,  
biased, consistent, sufficient statistic, Bayesian,  
estimator, prior, closed form, simulation, function  
writing, R

If the terms in blue were familiar to you (even if you can't recall their exact meaning), you are probably in good shape for this class.







<https://www.r-project.org/>



Contrary to (un)popular opinion,

<https://www.r-project.org/>



Contrary to (un)popular opinion, R has **excellent** text handling/modeling capabilities.

<https://www.r-project.org/>



<https://www.r-project.org/>

Contrary to (un)popular opinion, R has **excellent** text handling/modeling capabilities.

Free, and massive online community writing packages and extending modeling abilities.



<https://www.r-project.org/>

Contrary to (un)popular opinion, R has **excellent** text handling/modeling capabilities.

Free, and massive online community writing packages and extending modeling abilities.

We will use **quanteda** and other packages.  
Need R version 4.0.3



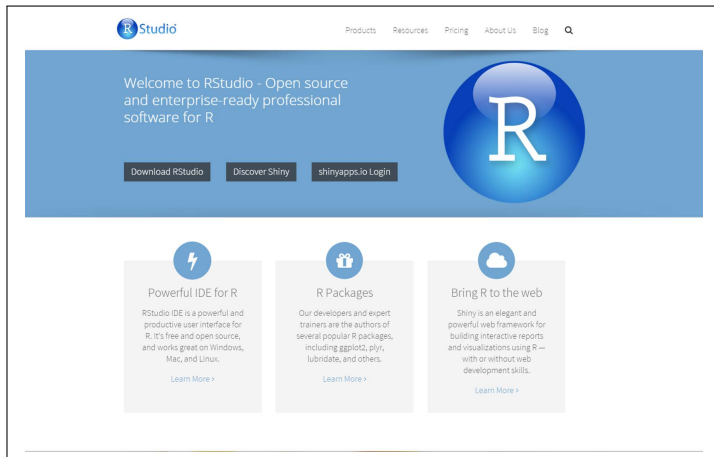
# Writing R: RStudio

# Writing R: RStudio

`https://www.rstudio.com/`

# Writing R: RStudio

`https://www.rstudio.com/`



# Readings

# Readings

new  
ent  
ne  
peo  
ke  
need  
ext  
pub  
form  
research  
law together  
matter republic  
south  
us  
sc  
live SUCC  
parents t  
since  
stand  
responsibility  
race  
sma  
scho

# Readings



- No required text book(s),

new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext  
ne  
feet  
us  
sc

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext  
ne  
feet  
us  
sc

- ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻



# Readings



- No required text book(s), though syllabus has some you might find interesting/helpful if you **really** want to part with \$s.
- Generally, we'll tell you where to get the readings, or **provide directly** if they are hard to obtain.

# Readings



- No required text book(s), though syllabus has some you might find interesting/helpful if you really want to part with \$s.
- Generally, we'll tell you where to get the readings, or provide directly if they are hard to obtain.
- Substantive readings are especially important,

# Readings



- No required text book(s), though syllabus has some you might find interesting/helpful if you **really** want to part with \$s.
- Generally, we'll tell you where to get the readings, or **provide directly** if they are hard to obtain.
- Substantive readings are especially important, because they'll help you understand what an **interesting question** looks like (in social science).

# Assessment

# Assessment

new  
ent  
ne  
peo  
ke  
need  
ext  
pub  
form  
research  
law together  
matter republic  
south  
us  
sc

since  
stand  
responsibility  
scho  
parents t  
live  
SUCC  
research  
us  
sc

# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills.

# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills. While you may consult with others when working, what you hand in must **be your own work** and **no one else's work**.

# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills. While you may consult with others when working, what you hand in must **be your own work** and **no one else's work**. Must use RMarkdown, or no grade.



# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills. While you may consult with others when working, what you hand in must **be your own work** and **no one else's work**. Must use RMarkdown, or no grade.

If you copy someone's work,

new  
no  
people  
need  
research  
together  
republic  
south  
form  
put  
ext  
ne  
feet  
us  
sc  
s  
live  
succ  
law  
matter  
parents  
stand  
responsibility  
since  
scho

If you copy someone's work, or allow someone to copy yours,

# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills. While you may consult with others when working, what you hand in must **be your own work** and **no one else's work**. Must use RMarkdown, or no grade.

If you copy someone's work, or allow someone to copy yours, you are breaking the rules of the class (and of NYU),

# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills. While you may consult with others when working, what you hand in must **be your own work** and **no one else's work**. Must use RMarkdown, or no grade.

If you copy someone's work, or allow someone to copy yours, you are breaking the rules of the class (and of NYU), and **bad** things happen: we will report you for **academic dishonesty**.

# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills. While you may consult with others when working, what you hand in must **be your own work** and **no one else's work**. Must use RMarkdown, or no grade.

If you copy someone's work, or allow someone to copy yours, you are breaking the rules of the class (and of NYU), and **bad** things happen: we will report you for **academic dishonesty**.

50%. **Final paper.**

# Assessment



50%. There will be **three homeworks** that will enable you to practice your skills. While you may consult with others when working, what you hand in must **be your own work** and **no one else's work**. Must use RMarkdown, or no grade.

If you copy someone's work, or allow someone to copy yours, you are breaking the rules of the class (and of NYU), and **bad** things happen: we will report you for **academic dishonesty**.

50%. **Final paper**. Must concern an original research question. Must involve **text-as-data** data and methods.

new  
no  
people  
need  
research

If you copy someone's work, or allow someone to copy yours, you are breaking the rules of the class (and of NYU), and **bad** things happen: we will report you for **academic dishonesty**.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

# Text as Data Speaker Series



# Text as Data Speaker Series

Researchers from all over the country,

# Text as Data Speaker Series

Researchers from all over the country, and industry (e.g. Facebook, Google) come to discuss their applied work and get feedback.

# Text as Data Speaker Series

Researchers from all over the country, and industry (e.g. Facebook, Google) come to discuss their applied work and get feedback. Mostly social science, NLP and computer science focussed.

# Text as Data Speaker Series

Researchers from all over the country, and industry (e.g. Facebook, Google) come to discuss their applied work and get feedback. Mostly social science, NLP and computer science focussed.

Sign up at: <http://cds.nyu.edu/text-data-speaker-series/>

# Text as Data Speaker Series

Researchers from all over the country, and industry (e.g. Facebook, Google) come to discuss their applied work and get feedback. Mostly social science, NLP and computer science focussed.

Sign up at: <http://cds.nyu.edu/text-data-speaker-series/>