

# Digital humanities: modeling semi-structured data from traditional scholarship

---

**Tom Lippincott**

IntroHLT Fall 2019

Human Language Technology Center of Excellence

Center for Language and Speech Processing

Intro: A few thoughts on “Digital humanities”

Motivating study: Post-Atlantic Slave Trade

Model: Graph-Entity Autoencoders

Bonus study: Authorship attribution of ancient documents

Ongoing work

## **Intro: A few thoughts on “Digital humanities”**

---

# What is “digital humanities”?

## Some responses:

- “an idea that will increasingly become invisible” -Stanford
- “a term of tactical convenience” -UMD
- “I don’t: I’m sick of trying to define it” -GMU
- “a convenient label, but fundamentally I don’t believe in it”  
-NYU
- “an unfortunate neologism” -Library of Congress

# What is “digital humanities”?

## Themes at DH2019

- Visualization
- Geographic information systems
- Social and ethical issues
- Education
- VR, maker spaces
- OCR
- Machine learning

# Working definitions

**Digital humanities**

**Traditional researcher**

**(Traditional) scholarly dataset**

**Computational researcher**

# Working definitions

## **Digital humanities**

Traditional inquiries enabled by computational intelligence

## **Traditional researcher**

**(Traditional) scholarly dataset**

## **Computational researcher**

# Working definitions

## **Digital humanities**

Traditional inquiries enabled by computational intelligence

## **Traditional researcher**

Academic from field that doesn't typically employ quantitative methods (History, Literary Criticism, . . .

## **(Traditional) scholarly dataset**

## **Computational researcher**



# Working definitions

## **Digital humanities**

Traditional inquiries enabled by computational intelligence

## **Traditional researcher**

Academic from field that doesn't typically employ quantitative methods (History, Literary Criticism, . . .

## **(Traditional) scholarly dataset**

Data assembled by a traditional researcher in the field

## **Computational researcher**

# Working definitions

## **Digital humanities**

Traditional inquiries enabled by computational intelligence

## **Traditional researcher**

Academic from field that doesn't typically employ quantitative methods (History, Literary Criticism, . . .

## **(Traditional) scholarly dataset**

Data assembled by a traditional researcher in the field

## **Computational researcher**

Design and bring machine learning models to bear on datasets

## Why is collaboration rare?

Traditional researchers have insight into the *data*

Machine learning researchers can pair *data* with appropriate models

# Why is collaboration rare?

## Traditional researchers have insight into the *data*

- Data is painstakingly gathered and coveted
- Hypotheses are *subtle* but not numerically evaluated
- May publish one or two papers during PhD, but *dissertation* is primary focus

## Machine learning researchers can pair *data* with appropriate models

# Why is collaboration rare?

## Traditional researchers have insight into the *data*

- Data is painstakingly gathered and coveted
- Hypotheses are *subtle* but not numerically evaluated
- May publish one or two papers during PhD, but *dissertation* is primary focus

## Machine learning researchers can pair *data* with appropriate models

- Data is aggressively shared to encourage rigorous evaluation
- Tasks are often *shallow* and *prespecified*
- Publish multiple papers per year

## Topic models: the rare success story

## Widely used

- Low barrier to entry: everyone has “documents”
- Little expertise required
- Output easy to visualize and interpret

## Widely used

- Low barrier to entry: everyone has “documents”
- Little expertise required
- Output easy to visualize and interpret

## Widely abused

- Deceptively easy to use: it will give you *something*
- You can always find “patterns”: confirmation bias abounds
- Older than some undergrads: LDA from early 2000s



## A guiding challenge:

Can we leverage sophisticated modeling techniques without losing the advantages that popularize topic models and recreating some of the same bad community practices?

### Financial analysts, investigative reporters ...

- Concerned with specific domains
- Need to gather, build, and understand datasets
- Wide range of technical abilities
- The DH story is relevant to industry, government, etc

## **Motivating study: Post-Atlantic Slave Trade**

---

# Shipping manifests

# Shipping manifests

SLAVE CLEARANCE.—Printed and Sold by A. E. MILES.

**Manifest** of Slaves, Passengers on board the Schooner *Wilcox* Capt. *J. W. Martin*.  
Tons, bound from Charleston, S. C. for New Orleans

*Master, Luther Higgins*

NAMES.	SEX.	AGE.	HEIGHT.		CLASS.	OWNERS or SHIPPERS.	RESIDENCE.
			FEET.	INCHES.			
<i>Ullis</i>	<i>Male</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>Black</i>	<i>Amid' Gardinier</i>	<i>New Orleans</i>
<i>Jack</i>	<i>do</i>	<i>25</i>	<i>5</i>	<i>—</i>	<i>do</i>	<i>do</i>	<i>do</i>
<i>Hector</i>	<i>do</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>do</i>	<i>do</i>	<i>do</i>
<i>Adam</i>	<i>do</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>do</i>	<i>do</i>	<i>do</i>
<i>Maria</i>	<i>Female</i>	<i>19</i>	<i>5</i>	<i>4</i>	<i>do</i>	<i>do</i>	<i>do</i>
<i>Mary</i>	<i>do</i>	<i>7</i>	<i>3</i>	<i>6</i>	<i>Mulatto</i>	<i>do</i>	<i>do</i>
<i>See Tons</i>					<i>Christen K. K. 32</i>	<i>J. W. Martin</i>	
<i>Examined</i>	<i>Walter</i>	<i>Sept. 27</i>	<i>1852</i>		<i>Examined and found to agree</i>	<i>William Vandergriff</i>	
<i>Compt.</i>	<i>Wm B. &amp; Baylon</i>					<i>Wm B. &amp; Baylon</i>	
						<i>New Orleans, Sept 24<sup>th</sup> 1852</i>	

# Shipping manifests

Manifest of Slaves, Passengers on board the Schooner *Wilcox* Capt. *J. W. Martin*.  
SLAVE CLEARANCE.—Printed and Sold by A. E. MILES.  
 Master, *Walter H. Taylor*. Tons, bound from *Charleston, S. C.* for *New Orleans*.

NAMES.	SEX.	AGE.	HEIGHT.		CLASS.	OWNERS or SHIPPERS.	RESIDENCE.
			FEET.	INCHES.			
<i>Ullis</i>	<i>Male</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>Black</i>	<i>Amid' Gardinier</i>	<i>New Orleans</i>
<i>Jack</i>	<i>d°</i>	<i>25</i>	<i>5</i>	<i>—</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Hector</i>	<i>d°</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Adam</i>	<i>d°</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Maria</i>	<i>Female</i>	<i>19</i>	<i>5</i>	<i>4</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Mary</i>	<i>do</i>	<i>7</i>	<i>3</i>	<i>6</i>	<i>Mulatto</i>	<i>— do —</i>	<i>d°</i>

*New Orleans*      *Charleston, S. C.*      *J. W. Martin*

*Examined Walter H. Taylor, Sept. 27, 1832*      *Examined and found to agree*  
*Somerset*      *Wm. H. Taylor*      *William Vandergift*  
*New Orleans, Sept 24, 1832*

slave name	slave sex	slave age	owner name	journey date	vessel type

# Shipping manifests

**Manifest** of Slaves, Passengers on board the Schooner *Willow* Capt. J. W. Martin. SLAVE CLEARANCE.—Printed and Sold by A. E. Miller.  
Master, Luther Higgins, Jr. Tons, loaded from Charleston, S. C. for New Orleans.

NAMES.	SEX.	AGE.	HEIGHT.		CLASS.	OWNERS or SHIPPERS.	RESIDENCE.
			FEET.	INCHES.			
Willis	Male	20	5	8	Black	Amidu Gardunphy	New Orleans
Jack	d <sup>o</sup>	20	5	—	d <sup>o</sup>	d <sup>o</sup>	d <sup>o</sup>
Hector	d <sup>o</sup>	20	5	8	d <sup>o</sup>	d <sup>o</sup>	d <sup>o</sup>
Adam	d <sup>o</sup>	20	5	8	d <sup>o</sup>	d <sup>o</sup>	d <sup>o</sup>
Maria	Female	19	5	4	d <sup>o</sup>	d <sup>o</sup>	d <sup>o</sup>
Mary	do	17	3	6	Mulatto	— do —	d <sup>o</sup>

New Orleans  
Charleston, S. C. Sept. 27, 1832  
Examined and found to agree  
William Vandergriff  
New Orleans, Sept 24, 1832

slave name	slave sex	slave age	owner name	journey date	vessel type
Willis	m	20	Amidu	1832/9/24	Schooner

# Shipping manifests

**Manifest** of Slaves, Passengers on board the Schooner *Willow* Capt. *J. W. Martin*.  
SLAVE CLEARANCE.—Printed and Sold by A. E. Miller.  
 Master, *Isaiah H. H. H. H.* Tons, loaded from Charleston, S. C. for New Orleans.

NAMES.	SEX.	AGE.	HEIGHT.		CLASS.	OWNERS or SHIPPERS.	RESIDENCE.
			FEET.	INCH.			
<i>Willis</i>	<i>Male</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>Black</i>	<i>Amidu Gardunphy</i>	<i>New Orleans</i>
<i>Jack</i>	<i>d°</i>	<i>25</i>	<i>5</i>	<i>—</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Hector</i>	<i>d°</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Adam</i>	<i>d°</i>	<i>20</i>	<i>5</i>	<i>8</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Maria</i>	<i>Female</i>	<i>19</i>	<i>5</i>	<i>4</i>	<i>d°</i>	<i>d°</i>	<i>d°</i>
<i>Mary</i>	<i>d°</i>	<i>17</i>	<i>5</i>	<i>6</i>	<i>Mulatto</i>	<i>d°</i>	<i>d°</i>

*New Orleans* *Charleston S. C. Sept. 27 1832* *J. W. Martin*  
*Examined and found to agree*  
*William Vandergriff*  
*New Orleans Sept 24<sup>th</sup> 1832*

slave name	slave sex	slave age	owner name	journey date	vessel type
Willis	m	20	Amidu	1832/9/24	Schooner
Maria	f	19	Amidu	1832/09/24	Schooner



# Fugitive notices

February 21.

1770

## Three Pistoles Reward.

RAN AWAY on the 15th day of October last, a black fellow, named Davy, stout made, about five feet three inches and a quarter high; has a scar on his left cheek, which is very apparent. It is suspected that he has made towards Pennsylvania, to which state he has twice attempted to make his escape. Who ever secures the above slave in any part (if taken out of the county) so as I get him, (if taken in the county, and delivered to me) shall receive the above reward.

Raphael Bourne.

Near Port Tobacco, Feb. 21.

1770

# Fugitive notices

February 21. 1770

## Three Pistoles Reward.

RAN AWAY on the 15th day of October last, a black fellow, named Davy, stout made, about five feet three inches and a quarter high; has a scar on his left cheek, which is very apparent. It is suspected that he has made towards Pennsylvania, to which state he has twice attempted to make his escape. Who ever secures the above slave in any part (if taken out of the county) so as I get him, (if taken in the county, and delivered to me) shall receive the above reward.

Raphael Bourne an.  
Near Port Tobacco, Feb. 21. 1770

slave name	slave sex	escape date	escape location	owner name	notice reward	notice date

# Fugitive notices

February 21. 1796

## Three Pistoles Reward.

RAN AWAY on the 15th day of October last, a black fellow, named Davy, stout made, about five feet three inches and a quarter high; has a scar on his left cheek, which is very apparent. It is suspected that he has made towards Pennsylvania, to which state he has twice attempted to make his escape. Who ever secures the above slave in any part (if taken out of the county) so as to get him, (if taken in the county, and delivered to me) shall receive the above reward.

Raphael Bourman.  
Near Port Tobacco, Feb. 21. 1796

slave name	slave sex	escape date	escape location	owner name	notice reward	notice date
Davy	m	1795/10/15	Port Tobacco	Bourman	3 Pistoles	1796/02/21

## Some numbers

- 45k manifest entries spanning five cities
- 11k fugitive notices from 70 gazettes
- 28k unique slave names
- 7k unique owner names
- Not big data, but thousands of studies like this at a research university!

## Difficulties with data in the wild

## Difficulties with data in the wild

- Unnormalized
  - People/places/things recorded *many* times
  - “What’s the age/height/sex distribution of escapees?”

## Difficulties with data in the wild

- Unnormalized
  - People/places/things recorded *many* times
  - “What’s the age/height/sex distribution of escapees?”
- Noisy
  - Vessel type: Bark, Barke, BAque, Barque, Barques
  - Slave name: “Nelly’?, Nelly’s child”, “not visible”
  - Owner sex: 3k missing



## Difficulties with data in the wild

- Unnormalized
  - People/places/things recorded *many* times
  - “What’s the age/height/sex distribution of escapees?”
- Noisy
  - Vessel type: Bark, Barke, BAque, Barque, Barques
  - Slave name: “Nelly’?, Nelly’s child”, “not visible”
  - Owner sex: 3k missing
- Underspecified entities
  - *Majority* of slaves have no last name
  - Can’t tell if two “Johns” are the same person

## What might a historian want to do with this data?

- Follow one slave throughout their life
- Group owners according to the nature of their workforce
- Determine what drove valuation in transactions and rewards
- Reconstruct slave families when there are no last names
- Map out trade “ecosystems” of sellers, shippers, owners, etc

# Fundamental observation

## There is an implicit *database schema* here

- **Field:** a recorded *value* with a clear interpretation (age, name, manufacturer ...)
- **Entity-type:** a coherent *bundle of fields* (person, location, object ...)
- Entity-types and fields have been determined by traditional scholars and common sense
- Relations *between* entities are also (conservatively) implied by the tabular format

## Fundamental observation

There is an implicit *database schema* here

- **Field**: a recorded *value* with a clear interpretation (age, name, manufacturer ...)
- **Entity-type**: a coherent *bundle of fields* (person, location, object ...)
- Entity-types and fields have been determined by traditional scholars and common sense
- Relations *between* entities are also (conservatively) implied by the tabular format

**This sets things up so we (ML researchers) can tackle the *general* problem**

# Entities, field types, and relations

## Traditional scholarly data

<i>slave_name</i>	<i>Jim</i>
<i>slave_age</i>	20
<i>owner_name</i>	<i>Jane</i>
<i>owner_sex</i>	<i>F</i>
<i>vessel_name</i>	<i>Uncas</i>
<i>vessel_type</i>	<i>Brig</i>
<i>voyage_date</i>	6/2/1823
<i>voyage_dest</i>	29.9, 90.0
...	...

## Numbers

<i>slave_name</i>	<i>Jim</i>
<i>slave_age</i>	<i>20</i>
<i>owner_name</i>	<i>Jane</i>
<i>owner_sex</i>	<i>F</i>
<i>vessel_name</i>	<i>Uncas</i>
<i>vessel_type</i>	<i>Brig</i>
<i>voyage_date</i>	<i>6/2/1823</i>
<i>voyage_dest</i>	<i>29.9, 90.0</i>
...	...

# Entities, field types, and relations

## Categories

<i>slave_name</i>	<i>Jim</i>
<i>slave_age</i>	20
<i>owner_name</i>	<i>Jane</i>
<i>owner_sex</i>	<i>F</i>
<i>vessel_name</i>	<i>Uncas</i>
<i>vessel_type</i>	<i>Brig</i>
<i>voyage_date</i>	6/2/1823
<i>voyage_dest</i>	29.9, 90.0
...	...

## Strings

*slave\_name* Jim

*slave\_age* 20

*owner\_name* Jane

*owner\_sex* F

*vessel\_name* Uncas

*vessel\_type* Brig

*voyage\_date* 6/2/1823

*voyage\_dest* 29.9, 90.0

... ..



# Entities, field types, and relations

## More complex fields

<i>slave_name</i>	<i>Jim</i>
<i>slave_age</i>	20
<i>owner_name</i>	<i>Jane</i>
<i>owner_sex</i>	<i>F</i>
<i>vessel_name</i>	<i>Uncas</i>
<i>vessel_type</i>	<i>Brig</i>
<i>voyage_date</i>	6/2/1823
<i>voyage_dest</i>	29.9, 90.0
...	...

## Entities

<i>slave_name</i>	<i>Jim</i>
<i>slave_age</i>	<i>20</i>
<i>owner_name</i>	<i>Jane</i>
<i>owner_sex</i>	<i>F</i>
<i>vessel_name</i>	<i>Uncas</i>
<i>vessel_type</i>	<i>Brig</i>
<i>voyage_date</i>	<i>6/2/1823</i>
<i>voyage_dest</i>	<i>29.9, 90.0</i>
...	...

# Entities, field types, and relations

## Entities

*slave\_name* Jim

*slave\_age* 20

*owner\_name* Jane

*owner\_sex* F

*vessel\_name* Uncas

*vessel\_type* Brig

*voyage\_date* 6/2/1823

*voyage\_dest* 29.9, 90.0


...

...

# Entities, field types, and relations

## Slave-to-owner


<i>slave_name</i>	<i>Jim</i>
<i>slave_age</i>	<i>20</i>
<i>owner_name</i>	<i>Jane</i>
<i>owner_sex</i>	<i>F</i>
<i>vessel_name</i>	<i>Uncas</i>
<i>vessel_type</i>	<i>Brig</i>
<i>voyage_date</i>	<i>6/2/1823</i>
<i>voyage_dest</i>	<i>29.9, 90.0</i>
...	...



# Entities, field types, and relations

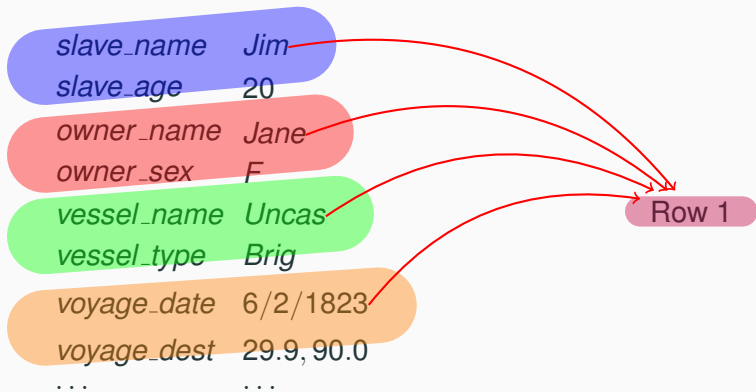
## Vessel-to-voyage, slave-to-voyage

<i>slave_name</i>	<i>Jim</i>
<i>slave_age</i>	20
<i>owner_name</i>	<i>Jane</i>
<i>owner_sex</i>	<i>F</i>
<i>vessel_name</i>	<i>Uncas</i>
<i>vessel_type</i>	<i>Brig</i>
<i>voyage_date</i>	6/2/1823
<i>voyage_dest</i>	29.9, 90.0
...	...

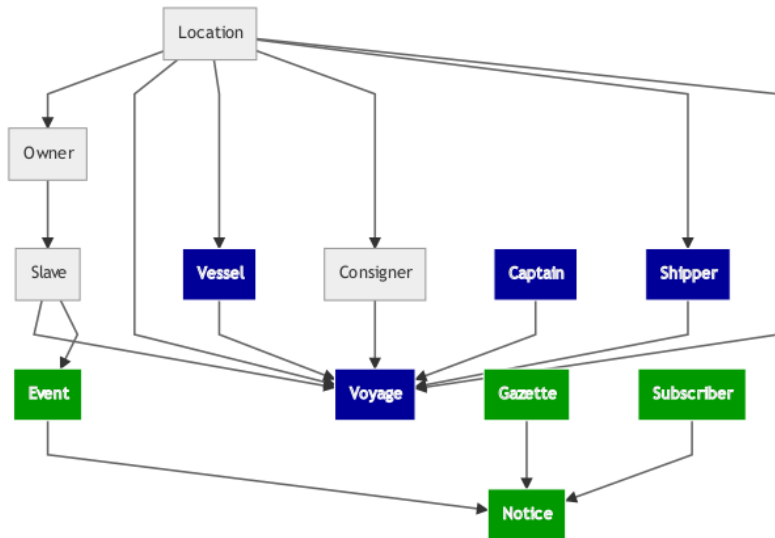


# Entities, field types, and relations

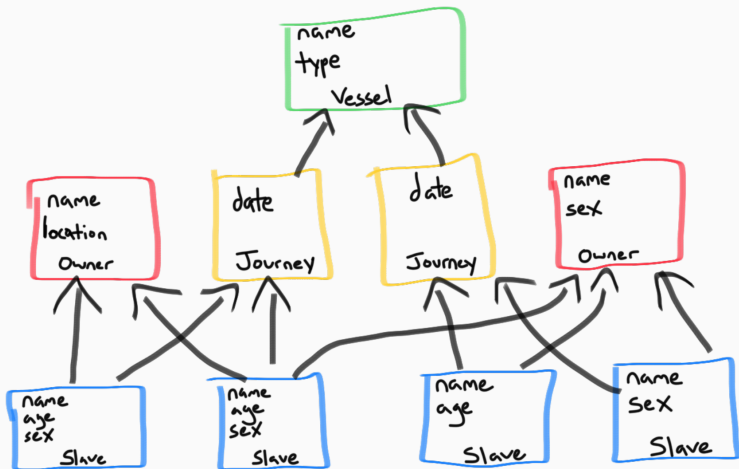
## Fewer assumptions



# Full schema of possible entity relationships

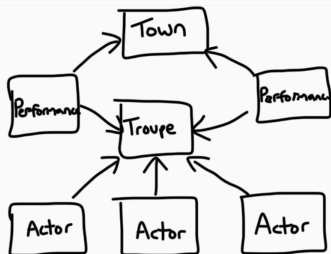


## Example data point: one graph component

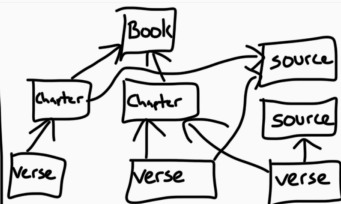




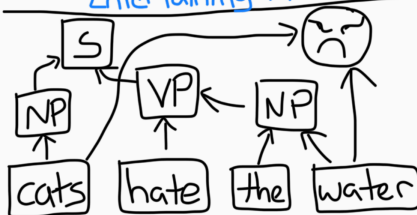
# Subsumes studies from a wide range of domains



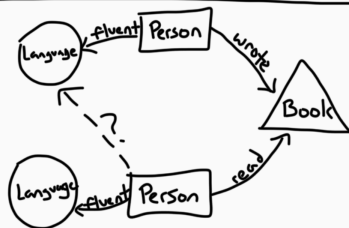
Entertaining America



Source criticism



Targeted sentiment



Network analysis

# **Model: Graph-Entity Autoencoders**

---

## General scholarly questions

- What is the overall picture of a *particular entity*?
- In what ways can we *group* entities?
- How are fields *correlated*?
- What *missing fields and relationships* can be recovered?

## General scholarly questions

- What is the overall picture of a *particular entity*?
- In what ways can we *group* entities?
- How are fields *correlated*?
- What *missing fields and relationships* can be recovered?

### Three basic operations we'd like:

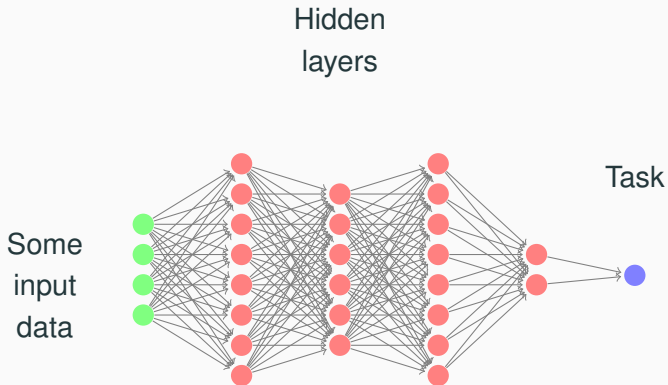
- Measure similarity of two entities
- Generate plausible field-values
- Score a proposed relationship between two entities

**Encoders, decoders, and autoencoders**

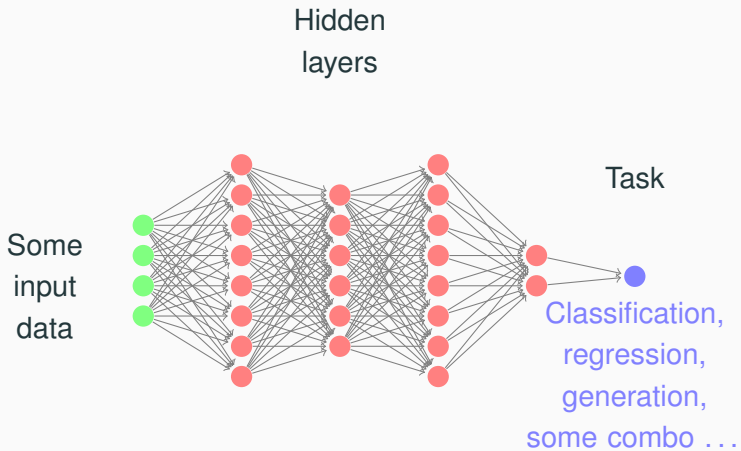
**Graph convolutional networks**



# Basic feed-forward neural model

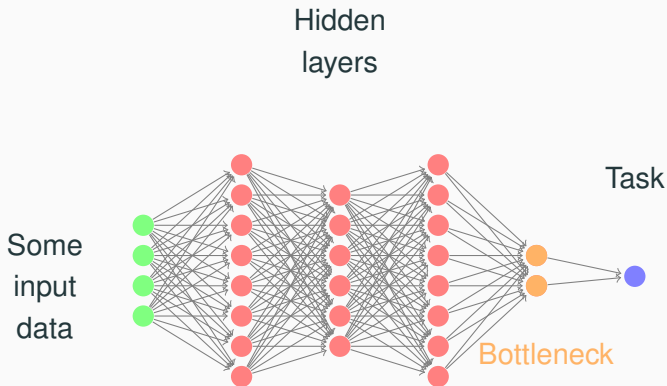


# Basic feed-forward neural model

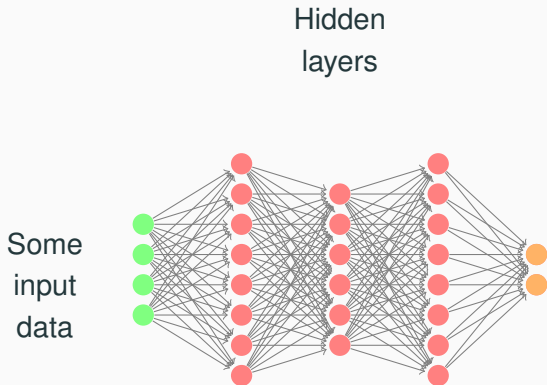




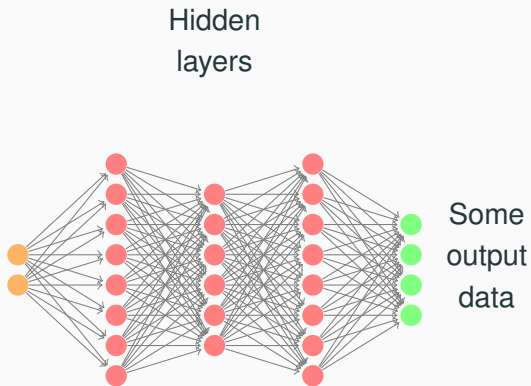
# Basic feed-forward neural model



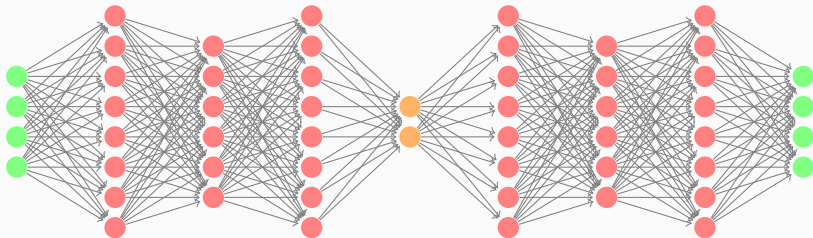
## Basic feed-forward neural model (an “encoder”)



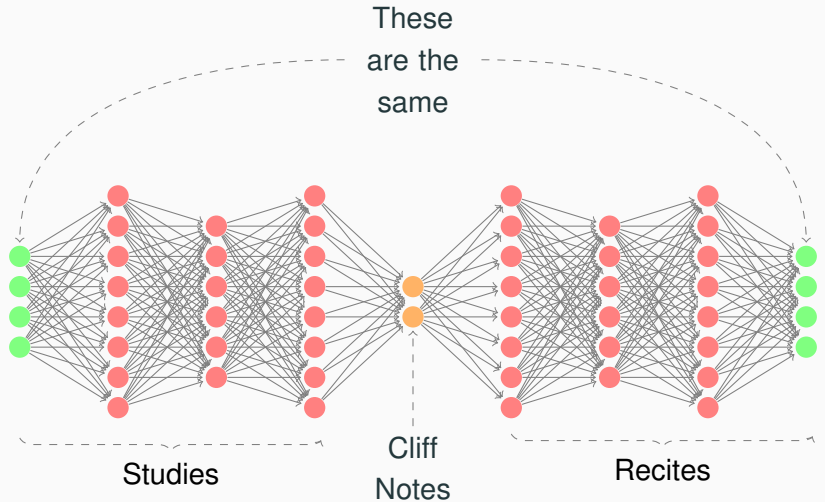
# A “decoder” goes in the opposite direction



# Encoders and decoders are often paired



# If the goal is to reconstruct the input, it's an "autoencoder"



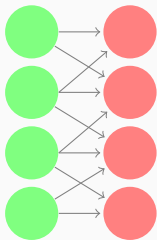
# Convolutional networks (CNNs)

Grid  
(image,  
text ...)



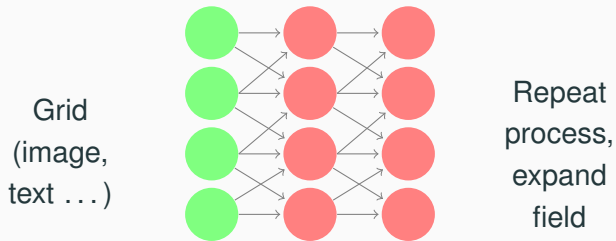
# Convolutional networks (CNNs)

Grid  
(image,  
text ...)



Each  
position  
incor-  
porates  
its “re-  
ceptive  
field”

# Convolutional networks (CNNs)



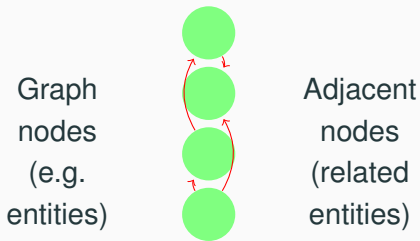


# Graph-convolutional networks (GCNs)

Graph  
nodes  
(e.g.  
entities)

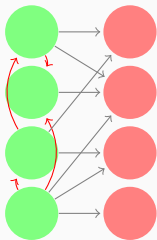


# Graph-convolutional networks (GCNs)



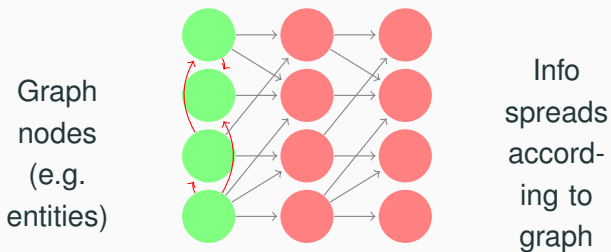
# Graph-convolutional networks (GCNs)

Graph  
nodes  
(e.g.  
entities)



Each  
node  
incorporates its  
neighbors

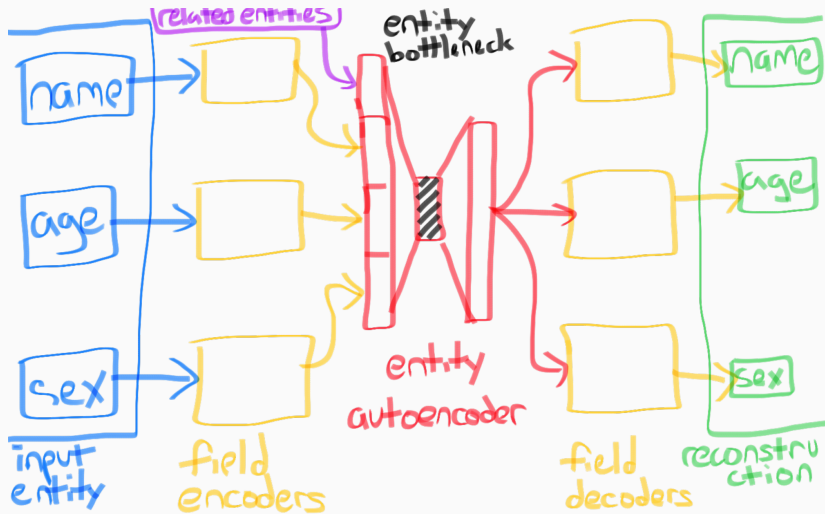
# Graph-convolutional networks (GCNs)



## Full model summary

- Read data, determine:
  - Fields and field-types
  - Entity-types
  - Relationships
- Each field allocated appropriate encoder-decoder pair
- Each entity-type allocated autoencoder
- Autoencoders use GCN-like mechanism to incorporate adjacent bottlenecks

# Model sketch



# Training is a complex process

- Random field dropout
- Graph component subselection
- Ways to combine loss functions
- ...

## How can we use a trained model?

- Compute distance between two entities
- Find flat or hierarchical clusters of entities
- Generate likely value of missing field
- Detect an improbable value of a present field
- Observe response of one field to another



# Example insights looking at most-similar entities

## Mistranscriptions

Baltiomre, Austin Woolfolk  $\Leftrightarrow$  Baltimore, Austin Woolfolk  
New Orleans, William Wiliams  $\Leftrightarrow$  New Orleans, William Williams

## Semantically-equivalent variants

Baltimore, George Y. Kelso  $\Leftrightarrow$  Baltimore, Kelso & Ferguson  
New Orleans, Leon Chabert  $\Leftrightarrow$  Louisiana, Leon Chabert

## Same slave transported multiple times

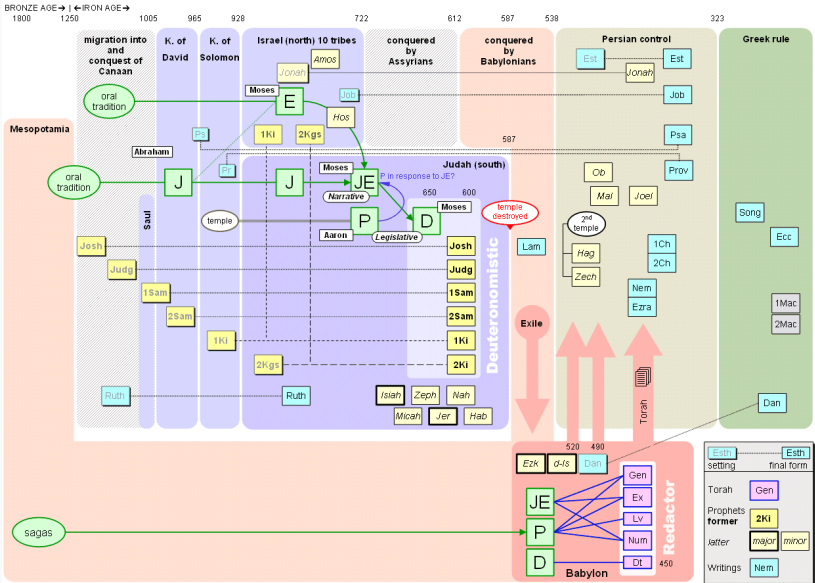
Louisa, F, 16yo  $\Leftrightarrow$  Louisa, F, 17yo  
Waters, F, 14yo  $\Leftrightarrow$  Waters, F, 15yo  
Kesiah, F, 20yo  $\Leftrightarrow$  Kesiah, F, 22yo  
Taylor, F, 15yo  $\Leftrightarrow$  Taylor, F, 16yo

**Bonus study: Authorship  
attribution of ancient documents**

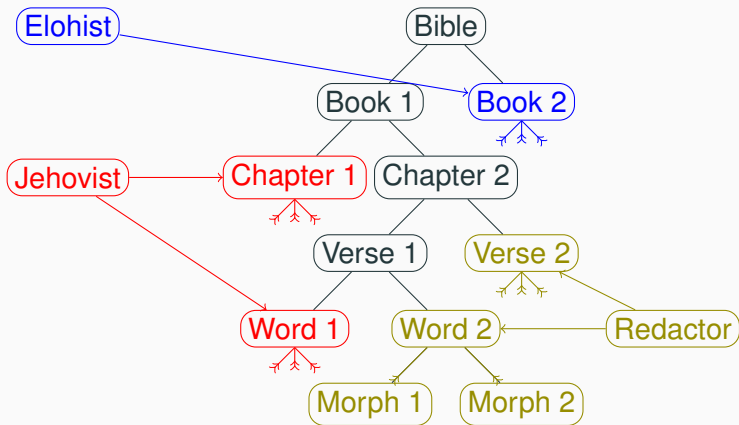
---

# Transmission of a text: the "Documentary Hypothesis"

Hebrew Bible sources timeline (Jewish Canon)



# Hypothesis as pointers into document structure



# Thomas Mendenhall: The Characteristic Curves of Composition



## SCIENCE.—SUPPLEMENT.

FRIDAY, MARCH 11, 1887.

### THE CHARACTERISTIC CURVES OF COMPOSITION.

AGUSTUS DeMORGAN somewhere remarks (I think it is in his 'Budget of paradoxes') that some time somebody will institute a comparison among writers in regard to the average length of

mean word-length suggested itself. The new method, while scarcely more laborious than that proposed by DeMorgan, promised to yield results more quickly and of a definitely higher order. It also had the advantage of including, in its application, all that was necessary to the determination of mean word-length; so that, in reality, it furnished two distinct tests.

Preliminary trials of the method have furnished

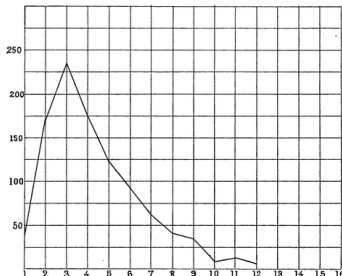


FIG. 1.—FIRST ONE THOUSAND WORDS IN 'OLIVER TWIST.'

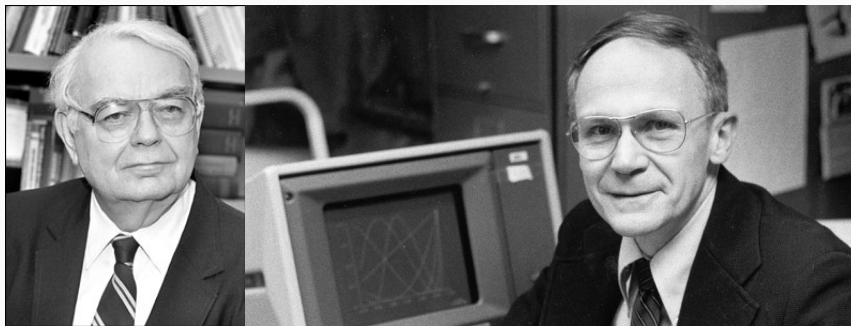
words used in composition, and that it may be found possible to identify the author of a book, a poem, or a play, in this way.

In reflecting upon this remark at various times within the past five or six years, always with the determination to test the value of the suggestion whenever time for the work seemed available, a more comprehensive and satisfactory method of analysis than that based simply upon

strong grounds for the belief that it may prove useful as a method of analysis leading to identification or discrimination of authorship, and it is therefore brought to the attention of the scientific and literary public in the hope that some one may be found who is at once able and willing to secure a satisfactory test of its validity.

The nature of the process is extremely simple, but it may be useful to point out its similarity to

# Mosteller and Wallace: Inference in an Authorship Problem



## The Federalist papers

- 85 articles written by Hamilton, Madison, and Jay
- 12 are unattributed
- Frequency analysis of *function words* determined Madison as author

# Back to the Documentary Hypothesis

## Problems

- The “authors” are also editors, redactors, synthesizers . . . they interact in context-dependent ways
- There is no predefined segmentation into “articles”
- We *know* more than function-words are important (e.g. name of God)

## Solutions

- Limit vocabulary to words that are used frequently by all authors
- Employ a GCN to exploit the document structure
- Take the DH for granted (for now)

## GEA predicts the author slightly better ...

Model	F-score
LR	41.39
MLP	47.45
<b>GEA</b>	<b>48.60</b>

Gold	Guess							
	J	E	P	1D	2D	nD	R	O
J	100	8	7	0	0	0	3	0
E	22	53	8	0	0	0	0	0
P	13	5	77	0	1	0	4	0
1D	2	0	2	7	1	0	0	0
2D	2	2	1	0	5	0	0	0
nD	0	0	0	1	0	0	0	0
R	3	3	11	0	0	0	33	0
O	2	0	1	0	0	0	1	0



## Sentiment and in-context word senses

- “wife” shows up as polygamous in older but monogamous in newer sources
- Redactor’s positive view of Aaron+Moses, violent story of rebellion

## Narrative continuity

- Abraham and Isaac story thought to *originally end with sacrifice*, changed by the Redactor
- “it was the season for grapes. ¶travel and geographic locations¿ They broke off some grapes.”

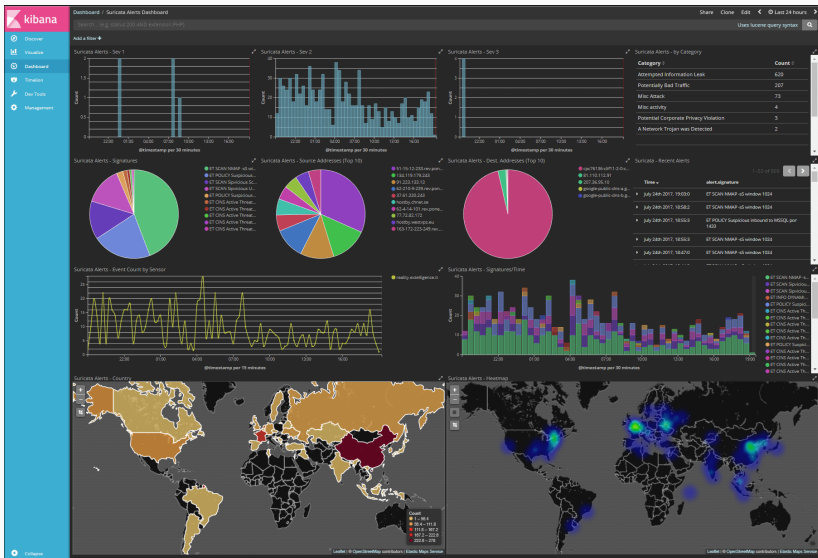
## Recipe for literary criticism

- Collect or construct useful resources from traditional scholarship
- Determine fit of potential “compositional actions” to observed document tree
- Choose the actions that are high-scoring *and* parsimonious
- Put the hypothesis in front of domain experts for verification/annotation

## Ongoing work

---

# Visualizing results



## Assembling other example studies

- JHU history department's "Entertaining America" (tabular)
- Northeastern U's Women Writers Collection (XML/TEI)
- Targeted sentiment analysis (JSON)
- Tennyson's poetic development (unconstrained text)

## Quick plug: come to David Mimno's talk!

- Nov. 15 at noon (Hackerman B17)
- CS Professor at Cornell
- Rare CS faculty working in DH (topic modeling)