Computational Intelligence for the Humanities

1

Tom Lippincott (tom@cs.jhu.edu)

Background: Working with humanists

Starcoder: A neural ensemble for humanities research

Example: The post-Atlantic slave trade

Ongoing work

Background: Working with humanists



Finance



Medicine



Military



Journalism



Finance



Medicine



Military



Journalism

- Concerned with specific domains, esoteric questions
- · Need to gather and structure targeted data
- · Explore and reason over those structures
- Find and present a coherent, interesting narrative
- · Wide range of technical abilities



Finance



Medicine



Military



Journalism



Humanist scholars

Why might computational intelligence help humanists?

Why might computational intelligence help humanists?

Scholars are capable of deep, focused reasoning, but ...

- · have baggage from centuries of received wisdom
- are often deeply invested in supporting/refuting a conclusion
- · not possible to attend evenly to data and hypotheses

Why might computational intelligence help humanists?

Scholars are capable of deep, focused reasoning, but ...

- · have baggage from centuries of received wisdom
- are often deeply invested in supporting/refuting a conclusion
- · not possible to attend evenly to data and hypotheses

Al has limited capacity for abstract reasoning, but ...

- · Easier to avoid bias towards a particular conclusion
- · Can uniformly process data and explore hypothesis-spaces
- May produce well-formed probabilities, well-specified representations

Supervised, focused tasks ("I can't label/transcribe/etc all this data!")

Supervised, focused tasks ("I can't label/transcribe/etc all this data!")

- Starts with careful task definition
- Requires (maybe lots of) example pairs (datum, label/transcription)
- OCR, translation, any sort of predefined categorization, annotation

Supervised, focused tasks ("I can't label/transcribe/etc all this data!")

- Starts with careful task definition
- Requires (maybe lots of) example pairs (datum, label/transcription)
- OCR, translation, any sort of predefined categorization, annotation

- · Starts with careful, unbiased representations
- · No specific task, minimal bias and annotation cost
- Anomaly detection, topic models, clustering, intrinsic graph metrics

Supervised, focused tasks ("I can't label/transcribe/etc all this data!")

- Starts with careful task definition
- Requires (maybe lots of) example pairs (datum, label/transcription)
- OCR, translation, any sort of predefined categorization, annotation

- · Starts with careful, unbiased representations
- · No specific task, minimal bias and annotation cost
- Anomaly detection, topic models, clustering, intrinsic graph metrics

Minimal disruption to existing practice

"Sustainability"

Broad goals for collaboration

Minimal disruption to existing practice

- Humanists can keep working directly with primary sources as they always have
- Computational scholars can develop and swap in new architectures and techniques

"Sustainability"

Broad goals for collaboration

Minimal disruption to existing practice

- Humanists can keep working directly with primary sources as they always have
- Computational scholars can develop and swap in new architectures and techniques

"Sustainability"

- · Can't manage dozens of bespoke collaborations
- · Need to focus on the boundaries: data in, hypotheses out
- · Start with simple model that's easily expanded

Starcoder: A neural ensemble for humanities research

owner_name	owner₋job	owner_age	car₋make	car_price	dealer_name
John	Lawyer	23	Honda	25000	Crazy Ray's
Jane	Doctor	31	Ford	35000	Crazy Ray's

owner_name	owner₋job	owner_age	car₋make	car_price	dealer_name
John	Lawyer	23	Honda	25000	Crazy Ray's
Jane	Doctor	31	Ford	35000	Crazy Ray's

Entities with properties and relationships

- Property-type (e.g. text, category, number, image, location...)
- An **Entity-type** is a coherent *bundle of potential properties* (e.g. an *owner* has a name, job, and age)
- A Relationship-type is a predicate with a specific interpretation that might link entities of the appropriate entity-types (e.g. owner_of, sold_by)

Graph view of the same data



Humanist describes their domain in a schema

owner_name	owner_job	owner_age	car_make	car_price	dealer_name
John	Lawyer	23	Honda	25000	Crazy Ray's
Jane	Doctor	31	Ford	35000	Crazy Ray's

Humanist describes their domain in a schema

owner_name	owner_job	owner_age	car_make	car_price	dealer₋name
John	Lawyer	23	Honda	25000	Crazy Ray's
Jane	Doctor	31	Ford	35000	Crazy Ray's

```
[ "properties" : {
    "owner_age" : {"type" : "scalar", "meaning" : "human_age"},
    ...
},
    "entity_types" : {
    "owner" : ["owner_name", "owner_job", ...]
    ...
},
    "relationships" : {
        "owned_by" : {
            "source_entity_type" : "car",
            "target_entity_type" : "owner"
        },
        ...
}
```

Schema combined with spreadsheet/XML/etc generates JSON entities

```
{"type" : "owner", "name" : "John", "age" : 23,
     "job" : "Lawyer", "id" : 0},
{"type" : "owner", "name" : "Jane", "age" : 31,
     "job" : "Doctor", "id" : 1},
{"type" : "dealer", "dealer_name" : "Crazy Ray's",
     "id" : 2},
{"type" : "car", "make" : "Honda", "price" : 25000,
     "owned_by" : 0, "sold_by" : 2},
{"type" : "car", "make" : "Ford", "price" : 35000,
     "owned_by" : 1, "sold_by" : 2}
```

1

Design a model-generator that matches the schema

Encoder, decoder, and autoencoder mechanisms Capture the *entities* and *fields*

Graph convolutional mechanism Capture the *relationships*

Encoder, decoder, and autoencoder mechanisms Capture the *entities* and *fields*

Graph convolutional mechanism Capture the *relationships*

Starcoder, from the Kleene-closure (asterisk/wildcard)

Intro to encoders/decoders/autoencoders

Feed-forward network



Feed-forward network



Feed-forward network



Encoder



Decoder



Encoders and decoders are often paired



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



Summary of coder mechanisms

An **encoder** transforms data into a fixed-length representation

• A **decoder** takes a fixed-length representation and generates data

• An **autoencoder** is an encoder and decoder working together to preserve data through a **bottleneck**
On to graph convolutions...



Grid (image, text ...)



Grid (image, text ...)



Each position incorporates its "receptive field"

Grid (image, text ...)



Repeat process, expand field Graph nodes (e.g. entities) Graph nodes (e.g. entities)



Adjacent nodes (related entities)

Graph convolutional network (GCN)

Graph nodes (e.g. entities)



Each node incorporates its neighbors

Graph convolutional network (GCN)

Graph nodes (e.g. entities)



Info spreads according to graph • Extends CNNs from grids to graphs

• Information can pass along edges

· Each layer allows nodes to see one further "hop"

· Encoders, decoders, autoencoders

Graph convolutional mechanism

Combine these to match the data being modeled



















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

• . . .

Bottleneck (embedding) similarity

- · Compute distance between two entities
- · Find flat or hierarchical clusters of entities

Field generation

- · Generate likely value of missing field
- · Detect an improbable value of a present field
- · Observe response of one field to another

Example: The post-Atlantic slave trade

Manifest of Slaves, Paramore on load the Schooner Willow Cal. J. W. W. Mantrue Burge Licher Storme Ver Master, baithen fiftyme Stys CLASS. OWNERS on SHIPPERS. AGE. NAMES. SEX. FEET. INCHES Black Amide Sardnue fil Mew Orleans 8. d Willis Male 20 25 do acto do d. N 20 do do 8 N do 20. do timah 19. 10 Mary. do to agree hamined Balize Septer Tamined and 150 Correct. m. B. S. Baylon New: Orleans. Dept 24th

Manifest of Slaves, Gaussiger in load the Scheoner Willow Cat. J. W. Manifest on Slaves, Jonanger in load the Scheoner Childow St. S. C. for Meur Orleans Master, baithen fiftyme Stys AGE. CLASS. OWNERS on SHIPPERS. NAMES. SEX. FEET. INCHES Black Amide Sardnue fil Mew Orleans 8. Willis Male 20 8 20 do lack d. N 20 do do 8 ata N do 20. do temah 19. late Mary. do The Cloves hamined Balize Dep Tamined ana Correct. mg3 & Baston New. Orleans. Jo

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

N	MES.	SEX.	AGE.	HEIGHT.	CLASS.	OWNERS on SI	IIPPERS.	RESIDENCE.
Mall	is m	ali	20. 1	8.	Black,	Amidu So	admie fit	. Mew Orles
Jack	- a	10	20- 4	7 -	X:	1 do		do
Adam	- 19	t.º	20.	5- 8.	do	de.	·	do_
Maria	- Pen	rah	19. 0	5 40	nulatte	X - a	10	do
C					Letter.	BR IM	11. 2	, a
Tix @ la	-15		- 195.27.4	sen	aler	2	Maran	k
								1 1 - 1
1	amined ;	Poalez.	- Septer	22.20	rad	6 ramouro	6 and 1	and for agree
14	Corner?	Baliz	- Depter	32 2. al	isd	6 xamauri	band further	andigitty
1	Camered : Corret?	malez.	Mepler M. B. S. Sa.	32 2. "	ind	New: Ort	land f	and gilly
1	Carmined : Correct?	Baliz.	- Ocher 12 3. 8. Sa.	32 2. al	rid	New: Or	land f	and to asses
lave	Compete Compete Slave	Bala Sla	Ve 0	2 2. Al	jou	Nuw: On	land f	ound to arrie anderitty pt 24th 183.
lave ame	save sex	lia (2 	Ve o e n	wner ame	jou da	Nuw, Ort	land f	ound to arrie and giffy of 24th 183. Sessel ype
lave ame Villis	slave Sex m	Isala sla age 20	Ve 0 e n	wner ame	jou da 18	Viur, Ort Irney te 32/9/24	tand the and t	essel ype Schooner
lave ame Villis	save sex m	slar slar age 20	ve o e n A	wner ame midu	jou da 18	Vrur, Ort Irrney te 32/9/24	teans the	here to accel and i gift and i gift at 24 to 180. At 24 to 180.

Manifest of Slaves, Paranger on load the Scherones Willow Cart I. W - Marters Bun late billow 25-Master, baithen fiftyme 5%5-AGE. CLASS. OWNERS OR SHIPPERS. NAMES. SEX. FEET. Plack Amide Sardmuchty Mew Orleans Male. di Willis 20 8 Bes's 25-Jack d. do N 20 Horton do No 8 do 20. do Adam de Timah 19. 0 do Maria do Mary. 4 1.0 ened and hamen Poars 150 m. B. S. Saylon tot 24 to New, Orleans. slave slave slave owner journey vessel date name sex age name type Willis nh F 20 Amidu 1832/9/24 Schooner 19 1832/09/24 Maria Amidu Schooner

(9/29/0.1 revruary 21. Three Pifoles Reward. RAN AW AT on the 15th day of titlober 1 ft, a black fellow, named Davy, flout made, about for feet three inches and a quarter bigh ; bas a fear on bis left cheek, aub ch is very apparent. It is fufpitted that be bas made towards Pentifitounia, to which fate be bas stance attempted to make bis feape. Who ver fecures the above flows in any good (if taken out of the county) to as I got bins, (if taken in the county, and desvered to me) feall receive the above required. Rayhael Beurn an. Near Port Toharco, Fib: 21. 1070

name

sex

date

revruary 21. (0/29/0.1 Three Pifoles Reward. RAN AWAY on the 15th day of tittoher lift, a black fellow, named Davy, flout mude, about for feit three inches and a quarter bigb ; bas a fear on bis left cheek, aub ch is very apparent. It is fufpitted that be bas made towards Pentifituania, to which fate be has stanice attempted to make bis sfeape. Who ver fectures the above fires in any good (if taken out of the county) to as I g t bim, (if taken in the county, and despotered to me) feall receive the above reasourch. Rayhael Beurn an. Near Port Toharco, Fib: 21. 1 070 slave slave notice notice escape escape owner

name

location

date

reward



• 45k manifest entries spanning five cities

11k fugitive notices from 70 gazettes

• Not big data, but thousands of studies like this at a research university!

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, BArque, Barque, Barques
 - · Slave name: "Nelly'?, Nelly's child", "not visible"

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, BArque, Barque, Barques
 - · Slave name: "Nelly'?, Nelly's child", "not visible"
- · Missing and underspecified entities
 - Majority of slaves have no last name
 - Can't tell if two "John"s 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
- Map out trade "ecosystems" of sellers, shippers, owners, etc
- · Reconstruct slave families when there are no last names
- Determine what drove valuation in transactions and rewards

Data

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

Numeric

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

Categorical

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

Text

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

More complex fields

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

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



Slave-to-owner



Vessel-to-voyage, slave-to-voyage



Fewer assumptions



One graph



Train Starcoder...

Looking at the most-similar pairs by entity-type, some trends emerge:

Mistranscriptions

Semantically-equivalent variants

George Y. Kelso	\Leftrightarrow	Kelso & Fergusor
New Orleans	\Leftrightarrow	Louisiana

Same slave transported multiple times¹

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

¹Many more instances of e.g. "John" following this superficial pattern

Ongoing work

- · Scaling: easy to add new studies
- · Research: new architectures to implement and import
- · Visualization: auto-generate figures based on schemas
- Interaction: well-defined interface to explore ML output
- Remains grounded in the particular scholarly domain

Data

- · Derived from local newspaper ads
- · Actors, troupes, towns, performances
- · Geo-coded town maps

Questions/Goals

- · Location of performances (within a town)
- · How performances are chosen and described

Chaucer's (and Gower's) Metrical Voice

Data

- Poetry (e.g. Canterbury Tales) with (partial) stress annotations
- · Entities are lines, stanzas, chapters, tales

Questions/Goals

- Train Starcoder as words2stresses, lines2rhymes
- Stresses as features for authorship
- Does the traditional scansion obscure interesting properties of the text?

Studies

- Tax records of Medieval Paris
- Employment in early modern London
- Economy of the Caribbean Colonies

Questions/Goals

- Impetus to treat dates and coordinates as first-class field types
- · Differential look at economic policies

Cuneiform in the Ancient Near East



- Traditionally, inscriptions are grouped by physical object-type
- We also have geographic location, transcription, ruler ...
- · What is the relationship between these properties?
- · Is the traditional view a useful distinction, and how?

Compositional forces on foundational cultural texts

Hebrew Bible sources timeline (Jewish Canon)



39

Some familiar NLP tasks

Language ID

- Communication network
- Users, messages, and languages
- Capture how multi-lingual users tweet and follow

Sentiment analysis

- · Parse trees
- Words (leaves), nodes (constituents), polarities
- Capture how sentiment composes

Also, compare with GraphSAGE, KB embedding techniques, etc