Data Catalog

Office of Chief Data Officer





Onintze Contreras Cornell University Information Science

What's the Problem?



	Location code for the asylum office with jurisdiction over		the asylum office with jurisdiction over the credible
APRSAFE-AP-CCO	the credible fear case	APSS.csv	fear case APRSAFEAPLASTN
APRSAFE-AP-LAST-NAME	last name.	APSS.csv	AME last name. APRSAFEAPFIRSTN
APRSAFE-AP-FIRST-NAME	first name.	APSS.csv	AME first name. APRSAFEAPMIDDL ENAME middle

Last name (EFTS field 2.907).

First name (EFTS field 2.908).

middle name.

APRSAFE-AP-MIDDLE-NAME

LAST_NAME

FIRST_NAME

APRSAFEAPCCO Location code for

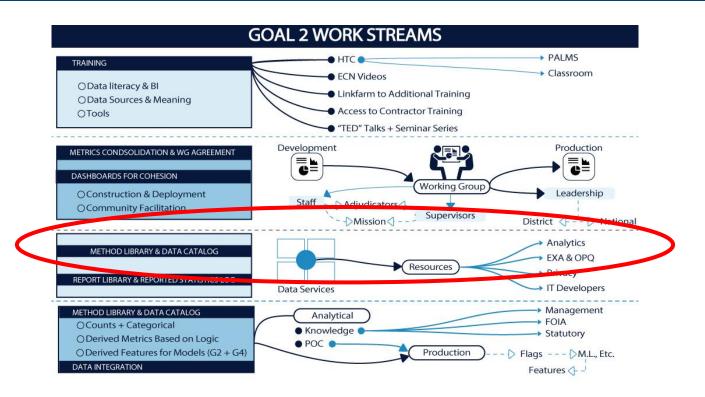
name.

BBSS.csv

BBSS.csv

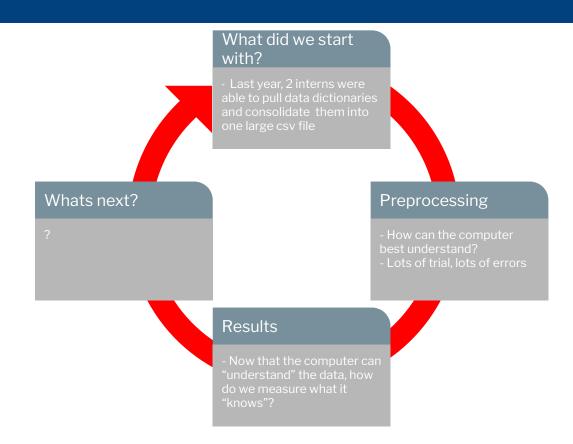
APSS.csv

What is it? Why?





The Process





What did we Start With?

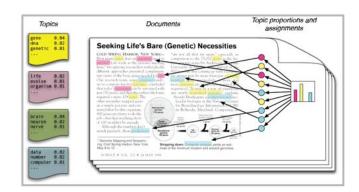
	A	В	C	D
1	Field	Field_Desc	Source	Combined
2	APRSAFE-AP-APPL-A-NUM	Applicant's A-number	APSS.csv	APRSAFEAPAPPLANUM Applicants Anumber
3	APRSAFE-AP-PRIN-A-NUM	Principal applicant-is A-number	APSS.csv	APRSAFEAPPRINANUM Principal applicants Anumber
4	APRSAFE-AP-CCO	Location code for the asylum office with jurisdiction over the credible fear case	APSS.csv	APRSAFEAPCCO Location code for the asylum office with jurisdiction over the credible fear case
5	APRSAFE-AP-LAST-NAME	last name.	APSS.csv	APRSAFEAPLASTNAME last name.
6	APRSAFE-AP-FIRST-NAME	first name.	APSS.csv	APRSAFEAPFIRSTNAME first name.
7	APRSAFE-AP-MIDDLE-NAME	middle name.	APSS.csv	APRSAFEAPMIDDLENAME middle name.
8	APRSAFE-AP-AKA-LAST-NAME	Alias for last name	APSS.csv	APRSAFEAPAKALASTNAME Alias for last name
9	APRSAFE-AP-AKA-FIRST-NAME	Alias for first name	APSS.csv	APRSAFEAPAKAFIRSTNAME Alias for first name
10	APRSAFE-AP-BIRTH-DATE	Birthdate	APSS.csv	APRSAFEAPBIRTHDATE Birthdate
11	APRSAFE-AP-SDX-LNAME-CODE		APSS.csv	APRSAFEAPSDXLNAMECODE
12	APRSAFE-AP-SDX-FNAME-CODE		APSS.csv	APRSAFEAPSDXFNAMECODE
13	APRSAFE-AP-SEX-CODE	Sex.	APSS.csv	APRSAFEAPSEXCODE Sex.
14	APRSAFE-AP-COB-CODE	Country of Birth	APSS.csv	APRSAFEAPCOBCODE Country of Birth
15	APRSAFE-AP-CITIZEN1-CODE	Applicant-ís 1st citizenship code	APSS.csv	APRSAFEAPCITIZEN1CODE Applicants 1st citizenship code
16	APRSAFE-AP-CITIZEN2-CODE	Applicant-is 2nd citizenship code	APSS.csv	APRSAFEAPCITIZEN2CODE Applicants 2nd citizenship code
17	APRSAFE-AP-APPL-LANG1-CODE	Applicant's 1st language code	APSS.csv	APRSAFEAPAPPLLANG1CODE Applicant's 1st language code
18	APRSAFE-AP-PORT-OF-ENTRY-CODE	Port of Entry Code	APSS.csv	APRSAFEAPPORTOFENTRYCODE Port of Entry Code
19	APRSAFE-AP-US-ARRIVAL-DATE	US Arrival Date	APSS.csv	APRSAFEAPUSARRIVALDATE US Arrival Date





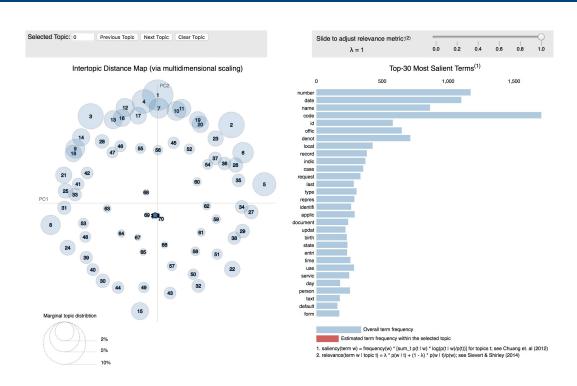
Preprocessing

- Load it into Python (Applicant's A-number)
- Stripped punctucation (applicantx96s a-number)
- Removed Stop Words (applicants a-number)
- Stemmed (applicant a-number)
- Tokenized ("applicant", 'a-number")

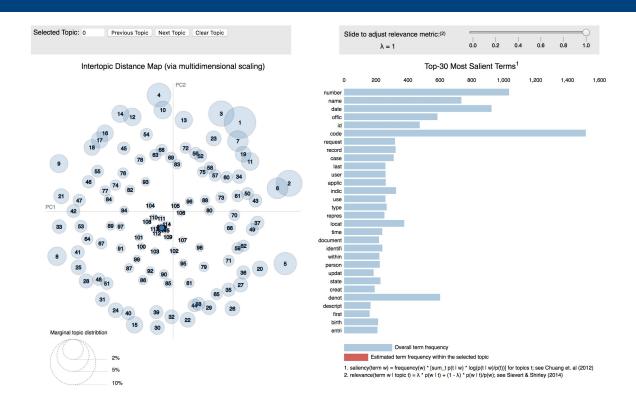




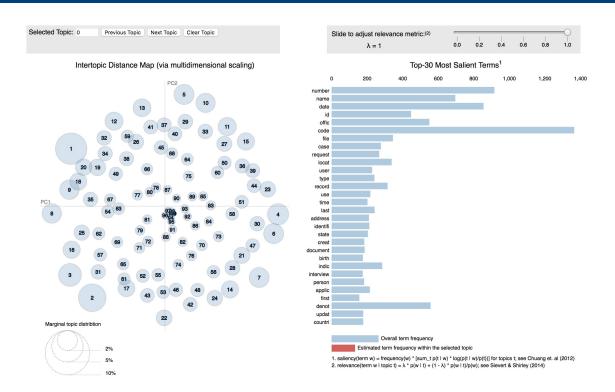




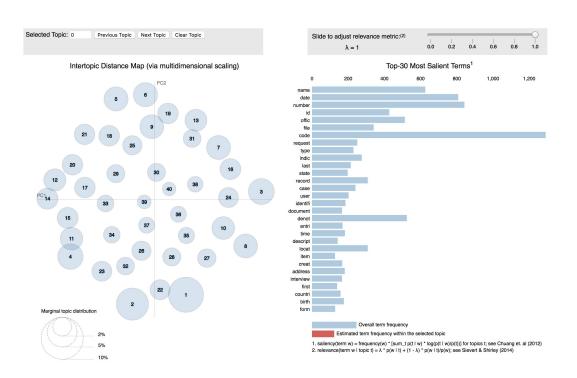




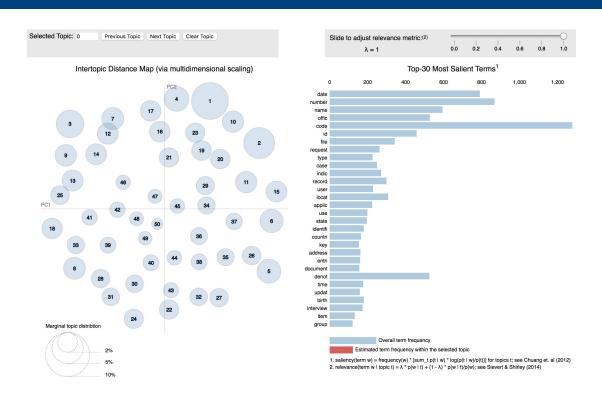




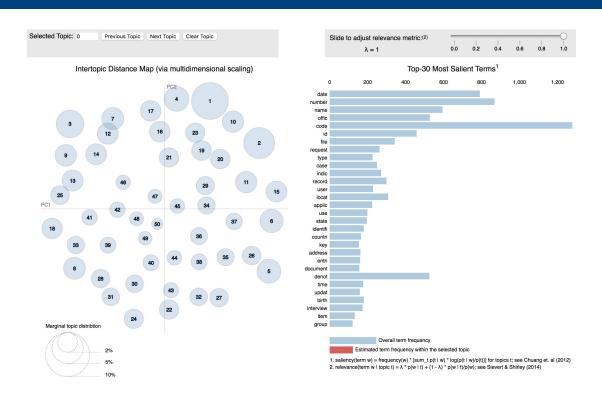






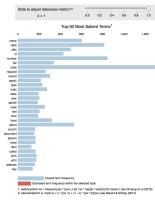


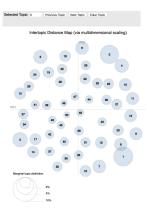


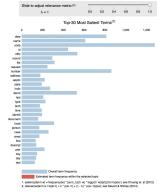


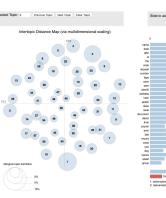


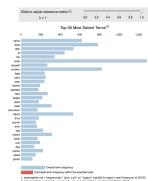




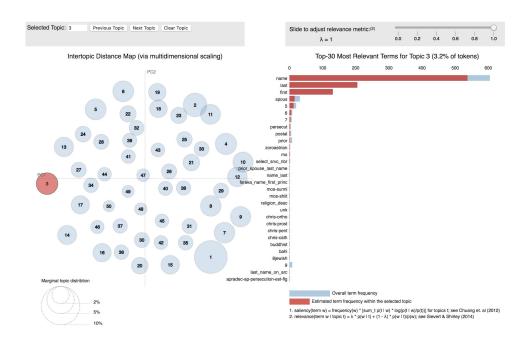












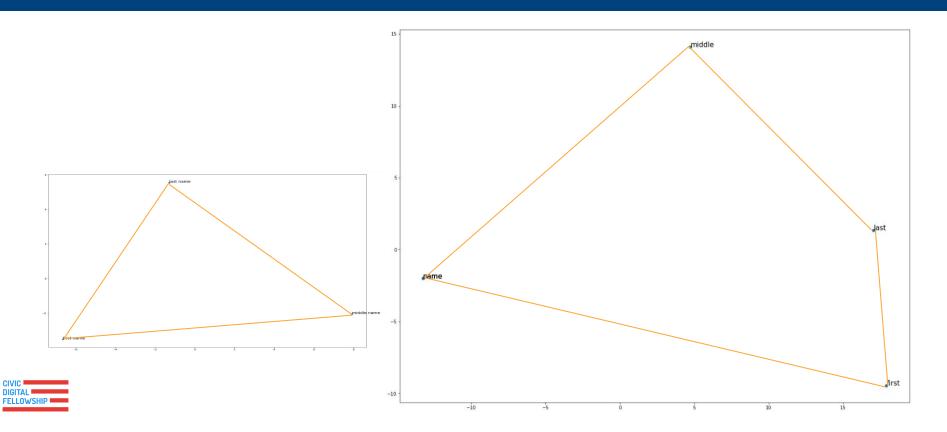


Takeaways

- Work with Data Standards Team (Data Governance)
- Try other word vectorizations:
 - Word2vec, GloVec
 - Different groupings
- Utilize named entity recognition customize (Stanford Named Entity Recognizer (NER))
- Semantic Parsing (NLTK)



Takeaways (continued)



Thanks!

- OCDO (Damian, Justin, everyone)
- RAIO(Logan, Michael, Nevin, Christine, everyone)
- All of you!

