

Week of:	Description of Activities:
8/24	<ul style="list-style-type: none"> <li>• TOTAL 2-3 hours</li> <li>• Met with Brandon, Viv, Mike. Introductions. Will meet weekly on Tuesdays 7:30am PT <ul style="list-style-type: none"> <li>◦ <a href="#">Teams meeting link for Tuesdays</a></li> <li>◦ Reviewed Student Project FY21 Notes (mission areas, project description, etc.)</li> <li>◦ Don't have access to <a href="https://code.usgs.gov/bserna/usgs-dash-project-generator">https://code.usgs.gov/bserna/usgs-dash-project-generator</a> (need login)</li> </ul> </li> </ul>
8/31	<ul style="list-style-type: none"> <li>• TOTAL 5-6 hours</li> <li>• Tuesday 1 hour <ul style="list-style-type: none"> <li>◦ Met with Brandon. Will connect on Slack, install CookieCutter, Goals: open dashboard code template and play around with sample data, share a GitHub project with Brandon by the weekend.</li> </ul> </li> <li>• Sunday 1 hour <ul style="list-style-type: none"> <li>◦ Read about CookieCutter, tried to unzip and open the usgs-dash-project-generator-master.zip.cpgz file without success.</li> <li>◦ Set up Slack #usgs_science_data_internship</li> <li>◦ Set up GitHub <a href="https://github.com/PortfolioSQA/USGS_Catalog_Dash">https://github.com/PortfolioSQA/USGS_Catalog_Dash</a></li> </ul> </li> <li>• Monday 3-4 hours <ul style="list-style-type: none"> <li>◦ Set-up environment and became more familiar with dash</li> <li>◦ Created a simple dashboard using sample data using my own code. Will request new cookiecutter template file again Tuesday morning.</li> </ul> </li> </ul>
9/7	<ul style="list-style-type: none"> <li>• TOTAL 7 hours</li> <li>• Tuesday 1 hour <ul style="list-style-type: none"> <li>◦ Met with Brandon. Created Slack Channel. Resent CookieCutter file and opened together. Goals for the week: check out the dash template, look at xml from Science Base Catalog (project -&gt; more details) and think about summary graphics/tables.</li> <li>◦ Notes: 8 USGS regions, 6-7 mission areas, many science centers! Connect with elastic search index to fill contents (makes aggregation harder), new elastic search package works with pandas dataframe.</li> </ul> </li> <li>• Friday/Sunday 2 hours <ul style="list-style-type: none"> <li>◦ Downloaded the cookiecutter template</li> <li>◦ Played with the dashboard app</li> <li>◦ Looked through Science Data Catalog xml files to determine which data is available</li> <li>◦ Proposed a new idea for dashboard to Brandon- Query/Summarize by map boundaries: <ul style="list-style-type: none"> <li>▪ User first filters by Mission Area, Science Topic, or none (dropdown menus)</li> <li>▪ User then demarcates geographic boundary requirement on interactive US Map (default is entire US)</li> <li>▪ Return record count and all records using geographic location and other filters</li> <li>▪ User may filter results further by Science Center, date, or sort by most frequently accessed records</li> </ul> </li> </ul> </li> <li>• Monday 4 hours <ul style="list-style-type: none"> <li>◦ Read about Dash, Mapbox, &amp; Leafly</li> <li>◦ Installed Shapely, Created Mapbox account</li> </ul> </li> </ul>

	<ul style="list-style-type: none"> <li>○ Comments: Polygons are not easy to work with for the user, could choose by selecting geographic boundary by county, state, USGS region, center, or select point and radial distance.</li> <li>○ What to include in sample data: name, keyword, center, lat/long, dates</li> </ul>
9/14	<ul style="list-style-type: none"> <li>• TOTAL 11 hours</li> <li>• Tuesday 3 hours <ul style="list-style-type: none"> <li>○ morning meeting canceled – will meet next Tuesday</li> <li>○ Installed faker and created sample dataset to use in dashboard</li> <li>○ Spent time understanding the bounding coordinates in the xml files <pre>&lt;bounding&gt; &lt;westbc&gt;-123.60443114739&lt;/westbc&gt; &lt;eastbc&gt;-97.764587398423&lt;/eastbc&gt; &lt;northbc&gt;45.032773794668&lt;/northbc&gt; &lt;southbc&gt;38.897445871348&lt;/southbc&gt; &lt;/bounding&gt;</pre> </li> </ul> </li> <li>• Sunday/Monday 8 hours <ul style="list-style-type: none"> <li>○ Worked on sample data &amp; dashboard</li> </ul> </li> <li>• Questions for Brandon/to figure out: <ul style="list-style-type: none"> <li>○ How to use a styled map from MapBox</li> <li>○ Pricing on MapBox – Feasible for USGS?</li> <li>○ Other options</li> <li>○ Discuss best way to map boundaries of data collection (&amp; data collected for entire US)</li> <li>○ Which summaries are most important for USGS? Why? (record count by mission area/science center/popular/keywords)</li> <li>○ Layer USGS regions on map?</li> <li>○ Use rows from data table without displaying all need columns (join?)</li> </ul> </li> </ul>
9/21	<ul style="list-style-type: none"> <li>• TOTAL 16 hours</li> <li>• Tuesday 1 hour <ul style="list-style-type: none"> <li>○ Met with Brandon</li> <li>○ To do: Get pricing information for Mapbox, get map working with filtering, fix table, complete aggregate graphs, &amp; place code in template</li> <li>○ Next week: Get SDC data to play with mapping, use styled maps</li> </ul> </li> <li>• Friday 3 hours <ul style="list-style-type: none"> <li>○ 'Learning App' Got map filtering working with sample data, fixed table, aggregated by visits</li> <li>○ Posted mapbox pricing</li> <li>○ Explored SDC data</li> </ul> </li> <li>• Saturday/Sunday 10 hours <ul style="list-style-type: none"> <li>○ Graphed SDC data, but need to clean it up to be able to filter by date/keyword etc.</li> <li>○ ETL SDC data <ul style="list-style-type: none"> <li>▪ Dropped NA for missing spatial (8 rows)</li> <li>▪ Replaced NA for dates, kept year (see questions)</li> <li>▪ Combined keyword columns</li> </ul> </li> <li>○ 'Sketch' of where to place map/table/graphs in template</li> </ul> </li> <li>• Monday 2 hours <ul style="list-style-type: none"> <li>○ 'learning_app' committed to Github – has SDC data graphed in map by date/science center, ETL python file in folder</li> <li>○ 'SDC_Map_Dash' template committed to Github - wireframe: <ul style="list-style-type: none"> <li>▪ Will start working on this version once we decide on layout, filtering</li> <li>▪ Filter by date, map, filter by keyword (not sci center)</li> </ul> </li> </ul> </li> </ul>

	<ul style="list-style-type: none"> <li> <ul style="list-style-type: none"> <li>▪ In each tab – Datatable + Summary Stats <ul style="list-style-type: none"> <li>• mappable US data (not nationwide)</li> <li>• Nationwide datasets – mapped?</li> <li>• world/earth/international datasets – mapped?</li> </ul> </li> </ul> </li> <li>• Questions: <ul style="list-style-type: none"> <li>○ If nan data is filtered by date, we lose about 10% (600+ rows). Lose the rows? Save them? Include Map/graph? I replaced nan with 1900-present</li> <li>○ Best way to determine global, continental, local data (using PlaceKeyword). Agree with three tabs?</li> <li>○ What to do with keywords?</li> <li>○ Could color markers by ...</li> <li>○ Moved away from selecting map area to fill table – table fills map (zooming allows user to determine if data is available in a particular area). Do we want map selection ability? Seems messier.</li> <li>○ Using integer for year – to_datetime doesn't like 1500s</li> <li>○ Next steps</li> </ul> </li> </ul>
9/28	<ul style="list-style-type: none"> <li>• TOTAL 7-8 hours</li> <li>• Tuesday 1 hour <ul style="list-style-type: none"> <li>○ Met with Brandon</li> <li>○ Notes: <ul style="list-style-type: none"> <li>▪ Viv may send invitation to bi-weekly meetings Mondays 9am PT – have not received invite</li> <li>▪ Filter by Science Center &amp; USGS Thesaurus Keyword (not by date) - done</li> <li>▪ Figure out a 'pretty' way to select all key words/all science centers – note: dash doesn't seem to have one, we can select all, but shows all science centers...</li> <li>▪ For now, separate global, US, mappable data by place keyword in three tabs - problematic</li> <li>▪ Create data table and aggregate statistics for each tab</li> <li>▪ Color markers by science center</li> </ul> </li> <li>○ Other: Fix map so when table returns nothing, no markers are mapped, when global or US data tab is selected – map shows what? Need to ignore rows in data table without lat/lon for the map graph data? Nan Science Center?</li> </ul> </li> <li>• Tuesday 2 hours <ul style="list-style-type: none"> <li>○ Updated journal, separated keywords, created unique USGS thesaurus keyword and science center list for dropdowns</li> </ul> </li> <li>• Monday 4-5 hours <ul style="list-style-type: none"> <li>○ Broke the app trying to add keyword search and spent a couple hours trying to fix it. Started from scratch.</li> <li>○ Have not successfully separated international data (see graph) so the tabs don't make sense right now. I could work on separating by lon/lat if we want to do this.</li> <li>○ If we still want to separate by area in three tabs, I need to place the map and science center choices inside each tab to be able to change the map/selection</li> <li>○ So many Science Centers... can use a select all/deselect etc. need to change font size. How should we display this?</li> </ul> </li> </ul>
10/5	<ul style="list-style-type: none"> <li>• TOTAL 10 hours</li> <li>• Tuesday 1 hour <ul style="list-style-type: none"> <li>○ Dropdown with select all – return all data</li> <li>○ Remove tabs – or use for summaries</li> <li>○ Zoom into populated marker area</li> </ul> </li> </ul>

	<ul style="list-style-type: none"> <li>○ Color markers</li> <li>○ Fill null Science Data Center with 'Undetermined'</li> <li>○ Counts/Graph</li> <li>○ Keywords</li> <li>○ <a href="https://data.usgs.gov/modelcatalog/">https://data.usgs.gov/modelcatalog/</a></li> <li>○ <a href="https://sciencebase.gov/datarelease/summary/">https://sciencebase.gov/datarelease/summary/</a></li> <li>○ Word cloud for the datatable keywords (% top 50 terms)</li> <li>• Friday 4 hours <ul style="list-style-type: none"> <li>○ Combined all data (instead of mappable, US, global)</li> <li>○ Fixed dropdown &amp; added All Science Centers</li> <li>○ Added keyword search</li> <li>○ Removed tabs</li> <li>○ Fixed null SC -&gt; Undetermined</li> <li>○ Return Count</li> </ul> </li> <li>• Saturday 4 hours <ul style="list-style-type: none"> <li>○ Fixed &amp; formatted datatable</li> <li>○ WordCloud</li> <li>○ Next Steps <ul style="list-style-type: none"> <li>▪ Zoom to populated markers</li> <li>▪ Color markers by science center</li> <li>▪ Make sure key words are searched lower case</li> <li>▪ Speed up wordcloud</li> <li>▪ Heroku</li> </ul> </li> </ul> </li> <li>• Sunday 1 hour <ul style="list-style-type: none"> <li>○ Researched zoom to markers and marker color</li> </ul> </li> </ul>
10/12	<ul style="list-style-type: none"> <li>• TOTAL 9 hours</li> <li>• Tuesday 5 hours <ul style="list-style-type: none"> <li>○ 1-hour meeting with Brandon. Notes: <ul style="list-style-type: none"> <li>▪ 10/19 check-in biweekly informal 8-8:30 am PT</li> <li>▪ 10/26 SDM Biweekly 9-10:30 PT</li> <li>▪ <del>Text changes (see github)</del></li> <li>▪ <del>place filters in same search box</del></li> <li>▪ <del>update datatable headers</del></li> <li>▪ <del>remove browse these data</del></li> <li>▪ <del>cache wordcloud for 'all'</del></li> <li>▪ <del>keywords - lowercase, remove punctuation</del></li> <li>▪ <del>controlled terms - USGS, etc.</del></li> <li>▪ <del>bug - CA/iso, count 2, returns 1??</del></li> <li>▪ Color by center/zoom into markers</li> <li>▪ <del>update button to download CSV</del></li> <li>▪ <del>Place 'downloading state' into dt and wc</del></li> <li>▪ <del>Heroku</del></li> </ul> </li> <li>○ 4 hours - completing above list</li> </ul> </li> <li>• Friday 4 hours <ul style="list-style-type: none"> <li>○ 2 hours – completing above list</li> <li>○ 2 hours - Connect with git SSH, fix files and attempted to deploy with Heroku</li> </ul> </li> </ul>
10/19	<ul style="list-style-type: none"> <li>• TOTAL 2.5 hours</li> <li>• Monday 1.5 hours <ul style="list-style-type: none"> <li>○ Meeting 0.5 hours <ul style="list-style-type: none"> <li>▪ 8am <a href="#">Join Microsoft Teams Meeting</a></li> <li>▪ CDI: <a href="https://my.usgs.gov/confluence/display/cdi/Home">https://my.usgs.gov/confluence/display/cdi/Home</a></li> <li>▪ Model catalog: <a href="https://data.usgs.gov/modelcatalog/">https://data.usgs.gov/modelcatalog/</a></li> </ul> </li> </ul> </li> </ul>

	<ul style="list-style-type: none"> <li> <ul style="list-style-type: none"> <li>Sciencebase data releases: <a href="https://www.sciencebase.gov/catalog/items?q=&amp;filter=systemType%3DData+Release&amp;filter=browseCategory%21%3DData+Release+-+In+Progress">https://www.sciencebase.gov/catalog/items?q=&amp;filter=systemType%3DData+Release&amp;filter=browseCategory%21%3DData+Release+-+In+Progress</a> </li> <li>1 hour – tried to get the app to deploy on Heroku <ul style="list-style-type: none"> <li>Manage.py file? heroku ps:scale web=1 ???</li> </ul> </li> </ul> </li> <li>Tuesday – 1 hour <ul style="list-style-type: none"> <li>Met with Brandon, deployed on Heroku</li> </ul> </li> </ul>
10/26	<ul style="list-style-type: none"> <li>TOTAL 4 hours</li> <li>Monday 3 hours</li> <li><a href="#">Join Microsoft Teams Meeting 9am</a> <ul style="list-style-type: none"> <li>Missed the meeting this morning. I'll be there next time.</li> </ul> </li> <li>3 hours <ul style="list-style-type: none"> <li>Fixed loading state for count</li> <li>If I place a loading state on map, it doesn't map correctly!? Can you help me figure out why that would be?</li> <li>sci_center colors – works with static data!!! Talk about how to fix this, if we want it.</li> </ul> </li> <li>Tuesday 1 hour <ul style="list-style-type: none"> <li>Meet with Brandon and Lisa Zolly</li> <li>Data needed: <ul style="list-style-type: none"> <li>Science center</li> <li>Latest harvest</li> <li>Status</li> <li>Doi citation info</li> <li>DOI</li> </ul> </li> <li>Filter by: <ul style="list-style-type: none"> <li>Science Center</li> <li>Dates (Beg, End, Updated, Latest Harvest)</li> <li>Status – Active/Inactive</li> </ul> </li> </ul> </li> </ul>
11/2	<ul style="list-style-type: none"> <li>TOTAL 11 hours</li> <li>Monday 1.5 hours <ul style="list-style-type: none"> <li>8am <a href="#">Join Microsoft Teams Meeting</a></li> </ul> </li> <li>Tuesday Meeting (0.5 hours)</li> <li>Looked at data (1 hour)</li> <li>To Do: <ul style="list-style-type: none"> <li>Filter by Science Center (datasource)</li> <li>Filter by Status - Active/Inactive (where is this?)</li> <li>Filter by Dates (avail: first_harvest_date, last_harvest_date, last_mdate_check_date, mdate) which?</li> <li>Filter: doi: None Do we use this?</li> <li>Return Table: SciCenter, # citations, doi, status, date?</li> <li>Return Count: dataset count with filters</li> <li>Return # of citations in the last month (don't have by date)</li> <li>Return: Pie Chart of Active/Inactive (if not filtered by this)</li> <li>Flag certain datasets? inactive with citation?</li> <li>Other Questions: <ul style="list-style-type: none"> <li>How do I know if data is active/inactive? Is this if the data has been harvested? What would the row look like if it weren't?</li> <li>For the ORCID availability for retrospective DOI assignment - Do I just check if doi is None?</li> </ul> </li> </ul> </li> </ul>

	<ul style="list-style-type: none"> <li>▪ Lisa talked about rate of growth of the collection over time (monthly growth? Which dates?) Bar chart</li> <li>▪ Do we want to show the related primary publications?</li> <li>▪ Include non-primary pubs for citation count?</li> <li>▪ Layout similar to dash template? Use tabs?</li> </ul> <ul style="list-style-type: none"> <li>• Friday 3.5 hours <ul style="list-style-type: none"> <li>○ etl</li> <li>○ New file from Cookiecutter template using sample data</li> <li>○ Filter by Science Center (datasource)</li> <li>○ Filter by Status - Active/Inactive</li> <li>○ Filter by Dates (beg, end, last updated (mdate), last harvest)</li> <li>○ Return Table: SciCenter, # citations, doi, status</li> </ul> </li> <li>• Saturday/Sunday 6 hours <ul style="list-style-type: none"> <li>○ Filter by Dates</li> <li>○ Return Count: dataset count with filters</li> <li>○ Return: Pie Chart of Active/Inactive (if not filtered by this)</li> </ul> </li> <li>• Questions/Issues: <ul style="list-style-type: none"> <li>○ When I align the CSV download link to the right, it breaks! Why?</li> <li>○ Using mdate instead of last_harvest_date doesn't work... Worked for a while trying to figure out problem with date type, couldn't find it</li> <li>○ Fix tab1 so it gives the datatable for selection (doesn't return to all). Why do the other tabs retain filter info, but not datatable?</li> </ul> </li> </ul>
11/9	<ul style="list-style-type: none"> <li>• TOTAL 8 hours</li> <li>• Monday 2.5 hours <ul style="list-style-type: none"> <li>○ 9am <a href="#">Join Microsoft Teams Meeting 9am</a></li> <li>○ Fix count and commit to new repository on GitHub: <a href="https://github.com/PortfolioSQA/SDC_Manager_Dashboard">https://github.com/PortfolioSQA/SDC_Manager_Dashboard</a></li> </ul> </li> <li>• Tuesday Meeting 1 hour <ul style="list-style-type: none"> <li>○ TODO: Fix table (so it doesn't change to default when tabs are changed), <del>Fix</del> <del>mdate</del></li> <li>○ Talked about NLP project - creating a custom NER to identify models from epubs warehouse?? for the Model catalog. If time permits start researching this.</li> <li>○ Also talked about volunteering for usability testing for Sophie</li> </ul> </li> <li>• Saturday 3 hours <ul style="list-style-type: none"> <li>○ Fixed mdate</li> <li>○ Worked on datatable persistence for a couple hours... I can't figure out what to do to keep the table from reverting to default in tab.</li> <li>○ Changed 0/1 =&gt; inactive/active</li> </ul> </li> <li>• Sunday 1.5 hours <ul style="list-style-type: none"> <li>○ Research on how to create custom NER model</li> <li>○ Tried for another 1-2 hours to debug the table persistence. Used dcc.store, persistence = True, and filtering in another callback. Didn't get it to work ☹</li> </ul> </li> </ul>
11/16	<ul style="list-style-type: none"> <li>• TOTAL 4.5 hours</li> <li>• Monday 0.5 hours <ul style="list-style-type: none"> <li>○ Monday 8am <a href="#">Join Microsoft Teams Meeting</a></li> </ul> </li> <li>• Tuesday – meeting cancelled due to power outage</li> <li>• Wednesday 3 hours <ul style="list-style-type: none"> <li>○ Deploy to Heroku</li> <li>○ Tab update works in Heroku. Weird!</li> <li>○ Fixed datatable sort</li> <li>○ Need to fix table so we can sort by status and right align download link</li> </ul> </li> </ul>

	<ul style="list-style-type: none"> <li>• Sunday 1 hour <ul style="list-style-type: none"> <li>◦ Created random dates in Excel for sample data</li> </ul> </li> </ul>
11/23	Thanksgiving Week – WEEK OFF
11/30	Final Exam – WEEK OFF
12/7	<ul style="list-style-type: none"> <li>• TOTAL 12.5 hours</li> <li>• 1.5 hours Monday <a href="#">Join Microsoft Teams Meeting 9am</a></li> <li>• 1 hour Tuesday meeting <ul style="list-style-type: none"> <li>◦ Methods: CNN, SVM, etc...</li> <li>Tasks: Image classification, image segmentation, etc...</li> <li>Applications (use cases): precipitation-induced landslide warning, tracking rainfall thresholds,</li> </ul> </li> <li>• 3 hours Tuesday <ul style="list-style-type: none"> <li>◦ Create sample data (dates for beg, end, update, harvest), change code to incorporate new 'dates'</li> </ul> </li> <li>• 3 hours Wednesday <ul style="list-style-type: none"> <li>◦ Debug code changes for new sample dates, deploy on Heroku, update journal, e-mail Sophie with updates</li> </ul> </li> <li>• 2 hours Thursday <ul style="list-style-type: none"> <li>◦ Create notes document for model scraping</li> <li>◦ Review text</li> </ul> </li> <li>• 2 hours Sunday <ul style="list-style-type: none"> <li>◦ extract abstract, identify method, task and application</li> <li>◦ Create txt files with abstracts</li> </ul> </li> </ul>
12/14	<ul style="list-style-type: none"> <li>• TOTAL 10 hours</li> <li>• 0.5 hours <ul style="list-style-type: none"> <li>◦ Monday 8am <a href="#">Join Microsoft Teams Meeting</a></li> </ul> </li> <li>• 1 hour <ul style="list-style-type: none"> <li>◦ 2pm - Meet with Lisa, Brandon, and Sophie</li> </ul> </li> <li>• TO DO LIST: <ul style="list-style-type: none"> <li>◦ <del>Rename dashboard</del></li> <li>◦ <del>Link – learn how metrics are calculated</del></li> <li>◦ Table – font choice?</li> <li>◦ <del>Org of interest –&gt; Choose your science center/program</del></li> <li>◦ <del>Select dataset status</del></li> <li>◦ Explanation of active/inactive</li> <li>◦ <del>Put date selection in same box</del></li> <li>◦ Justify radio buttons</li> <li>◦ <del>Colon after headers</del></li> <li>◦ <del>Remove beg/end dates</del></li> <li>◦ <del>Last harvested, last updated –&gt; start date?</del></li> <li>◦ Flexibility in search by date? Calendar year?</li> <li>◦ <del>Dataset Count –&gt; Active, Inactive, Total, pie chart</del></li> <li>◦ <del>Make stacked bar graph for active inactive by date</del></li> <li>◦ <del>Remove pie chart tab</del></li> <li>◦ <del>doi – links?</del></li> <li>◦ <del>Alternate IDs for datasets?</del></li> <li>◦ <del>Delete widgets in graph area</del></li> <li>◦ <del>Format date label</del></li> <li>◦ <del>Download –&gt; more usable file name</del></li> <li>◦ Help &amp; documentation (Tool tips)</li> <li>◦ Errors: If it's not working what happens?</li> </ul> </li> <li>• 8.5 hours</li> </ul>

	<ul style="list-style-type: none"> <li>○ Complete TODO list, redesign site, insert links for DOI, alternate identifiers, stacked bar chart for counts by date, delete widgets in graphs</li> </ul>
12/21	<b>WEEK OFF - HOLIDAY</b>
12/28	<ul style="list-style-type: none"> <li>• TOTAL 7-8 hours</li> <li>• 30 min -Monday 8am <a href="#">Join Microsoft Teams Meeting</a> <ul style="list-style-type: none"> <li>○ <a href="https://data.usgs.gov/datacatalog/">https://data.usgs.gov/datacatalog/</a></li> </ul> </li> <li>• 1 hour – Tuesday Meeting with Brandon <ul style="list-style-type: none"> <li>○ Questions: title, <del>markdown</del> – open link in new tab/window for DOI, filter by active/inactive?, explanation for active inactive, Errors</li> <li>○ <b>TODO for dashboard:</b> table → dbc, tab (dash → SDC Dashboard), Sentence – <del>bold or large numbers,</del></li> <li>○ <del><a href="https://dash-bootstrap-components.opensource.faculty.ai/docs/components/table/">https://dash-bootstrap-components.opensource.faculty.ai/docs/components/table/</a></del> – dbc table</li> <li>○ <a href="https://data.usgs.gov/datacatalog/api/docs/v1">https://data.usgs.gov/datacatalog/api/docs/v1</a> , swagger API – look at documentation to get data except citation (leave as null for now)</li> <li>○ <b>TODO for NER:</b> <del>exploratory functional prototype: look at</del> <a href="https://paperswithcode.com/">https://paperswithcode.com/</a> , format text in jsonl, get prodigy (costs \$390),</li> </ul> </li> </ul> <pre>{ "text": "This is a text" } { "text": "This is another text" }</pre> <p><a href="https://data.usgs.gov/modelcatalog/search">https://data.usgs.gov/modelcatalog/search</a></p> <p>Training Data</p> <ol style="list-style-type: none"> <li>1. <a href="https://pubs.usgs.gov/tm/14/a2/tm14a2.pdf">https://pubs.usgs.gov/tm/14/a2/tm14a2.pdf</a></li> <li>2. <a href="https://pubs.usgs.gov/of/2008/1159/downloads/pdf/OF08-1159.pdf">https://pubs.usgs.gov/of/2008/1159/downloads/pdf/OF08-1159.pdf</a></li> <li>3. <a href="https://pubs.usgs.gov/of/2016/1136/ofr20161136.pdf">https://pubs.usgs.gov/of/2016/1136/ofr20161136.pdf</a></li> <li>4. <a href="https://pubs.usgs.gov/of/2007/1088/pdf/of07-1088_508.pdf">https://pubs.usgs.gov/of/2007/1088/pdf/of07-1088_508.pdf</a></li> <li>5. <a href="https://www.mdpi.com/2073-4441/8/1/17">https://www.mdpi.com/2073-4441/8/1/17</a></li> <li>6. <a href="https://pubs.usgs.gov/wri/1990/4130/report.pdf">https://pubs.usgs.gov/wri/1990/4130/report.pdf</a></li> <li>7. <a href="https://pubs.usgs.gov/tm/12b1/">https://pubs.usgs.gov/tm/12b1/</a></li> <li>8. <a href="https://pubs.usgs.gov/tm/2006/tm6b3/">https://pubs.usgs.gov/tm/2006/tm6b3/</a></li> <li>9. <a href="https://www.mdpi.com/1999-4893/1/2/52">https://www.mdpi.com/1999-4893/1/2/52</a></li> <li>10. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485782ce25b5135abf28">https://data.usgs.gov/modelcatalog/data/5eb4485782ce25b5135abf28</a></li> <li>11. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485082ce25b5135abee3">https://data.usgs.gov/modelcatalog/data/5eb4485082ce25b5135abee3</a></li> <li>12. <a href="https://data.usgs.gov/modelcatalog/data/5f6240eb82ce38aaa2361498">https://data.usgs.gov/modelcatalog/data/5f6240eb82ce38aaa2361498</a></li> <li>13. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485e82ce25b5135abf70">https://data.usgs.gov/modelcatalog/data/5eb4485e82ce25b5135abf70</a></li> <li>14. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485f82ce25b5135abf86">https://data.usgs.gov/modelcatalog/data/5eb4485f82ce25b5135abf86</a></li> <li>15. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485f82ce25b5135abf7c">https://data.usgs.gov/modelcatalog/data/5eb4485f82ce25b5135abf7c</a></li> <li>16. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485582ce25b5135abf16">https://data.usgs.gov/modelcatalog/data/5eb4485582ce25b5135abf16</a></li> <li>17. <a href="https://data.usgs.gov/modelcatalog/data/5f036b3f82ce0afb2446e04a">https://data.usgs.gov/modelcatalog/data/5f036b3f82ce0afb2446e04a</a></li> <li>18. <a href="https://data.usgs.gov/modelcatalog/data/5eb4486182ce25b5135abfb0">https://data.usgs.gov/modelcatalog/data/5eb4486182ce25b5135abfb0</a></li> <li>19. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485282ce25b5135abef4">https://data.usgs.gov/modelcatalog/data/5eb4485282ce25b5135abef4</a></li> <li>20. <a href="https://data.usgs.gov/modelcatalog/data/5eb4486182ce25b5135abfaa">https://data.usgs.gov/modelcatalog/data/5eb4486182ce25b5135abfaa</a></li> </ol> <ul style="list-style-type: none"> <li>• 3 hours – Wednesday <ul style="list-style-type: none"> <li>○ Update journal, look at prodigy (\$390 license), get 10 more text files, create jsonl file (see above links)</li> </ul> </li> <li>• 3 hours --Thursday <ul style="list-style-type: none"> <li>○ Dashboard updates, dbc tables are SUPER SLOW - decided to format dash data table (I have another file with the dbc table I can show if you'd like)</li> </ul> </li> </ul>
1/4	<ul style="list-style-type: none"> <li>• TOTAL 9 hours</li> <li>• 2 hours - Monday: <ul style="list-style-type: none"> <li>○ Doctor's Appt- Couldn't join meeting <a href="#">Join Microsoft Teams Meeting 9am</a></li> <li>○ 2 hours Install prodigy, en_core_web_sm, etc. Label text, git commit, etc.</li> </ul> </li> <li>• 1 hour - Tuesday</li> </ul>



	<ul style="list-style-type: none"> <li>○ Meet with Brandon 7:30 – will meet with Mike &amp; Viv soon, Ruby Gem linguist shows which languages in code,</li> <li>○ TODO: <del>10-20 unlabeled texts for evaluation data (good and bad examples)</del>, run in prodigy (see slack), try again using shorter NEs, use en_core_web_lg, future: establish list of methods, concept of datasets on prodigy?, deploy dbc table on Heroku and see if its faster</li> </ul> <ul style="list-style-type: none"> <li>• 3 hours – Thursday <ul style="list-style-type: none"> <li>○ 1+ hours Meet with Sophie to talk about Usability Testing</li> <li>○ 2 hours Evaluation Data (in excel &amp; jsonl)</li> </ul> </li> <li>• 3 hours – Friday <ul style="list-style-type: none"> <li>○ Train/evaluate ner – 0% accuracy, relabeled training set 0% accuracy, read more documentation, tried with models only, frustrating - next combined to 40 for training set and re-labeled. Didn't help</li> </ul> </li> </ul> <ol style="list-style-type: none"> <li>1. <a href="https://doi.org/10.1126/science.aat4723">https://doi.org/10.1126/science.aat4723</a></li> <li>2. <a href="https://pubs.usgs.gov/of/2001/ofr-01-0002/">https://pubs.usgs.gov/of/2001/ofr-01-0002/</a></li> <li>3. <a href="https://data.usgs.gov/modelcatalog/data/5ff62dc1d34ea5387df035fa">https://data.usgs.gov/modelcatalog/data/5ff62dc1d34ea5387df035fa</a></li> <li>4. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485982ce25b5135abf3c">https://data.usgs.gov/modelcatalog/data/5eb4485982ce25b5135abf3c</a></li> <li>5. <a href="https://doi.org/10.1111/gwat.12397">https://doi.org/10.1111/gwat.12397</a></li> <li>6. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485482ce25b5135abf0e">https://data.usgs.gov/modelcatalog/data/5eb4485482ce25b5135abf0e</a></li> <li>7. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485682ce25b5135abf1c">https://data.usgs.gov/modelcatalog/data/5eb4485682ce25b5135abf1c</a></li> <li>8. <a href="https://data.usgs.gov/modelcatalog/data/5eb4485382ce25b5135abefc">https://data.usgs.gov/modelcatalog/data/5eb4485382ce25b5135abefc</a></li> <li>9. <a href="https://doi.org/10.3133/wri874163">https://doi.org/10.3133/wri874163</a></li> <li>10. <a href="https://doi.org/10.1002/2017JC013204">https://doi.org/10.1002/2017JC013204</a></li> <li>11. <a href="https://doi.org/10.1029/2011JB008968">https://doi.org/10.1029/2011JB008968</a></li> <li>12. <a href="https://doi.org/10.3133/wri974022">https://doi.org/10.3133/wri974022</a></li> <li>13. <a href="https://pubs.usgs.gov/tm/tm4f2/">https://pubs.usgs.gov/tm/tm4f2/</a></li> <li>14. <a href="https://doi.org/10.1016/j.ocemod.2010.07.010">https://doi.org/10.1016/j.ocemod.2010.07.010</a></li> <li>15. <a href="https://doi.org/10.3133/ofr20151009">https://doi.org/10.3133/ofr20151009</a></li> <li>16. <a href="https://data.usgs.gov/modelcatalog/data/5ef3952782ced62aaae3ef55">https://data.usgs.gov/modelcatalog/data/5ef3952782ced62aaae3ef55</a></li> <li>17. <a href="https://pubs.usgs.gov/tm/tm6a37/">https://pubs.usgs.gov/tm/tm6a37/</a></li> <li>18. <a href="https://doi.org/10.1002/joc.3625">https://doi.org/10.1002/joc.3625</a></li> <li>19. <a href="https://doi.org/10.3133/tm14A1">https://doi.org/10.3133/tm14A1</a></li> <li>20. <a href="https://doi.org/10.3133/tm6A43">https://doi.org/10.3133/tm6A43</a></li> </ol>
1/11	<ul style="list-style-type: none"> <li>• TOTAL 9.5 hours</li> <li>• 2 hours - Monday: <ul style="list-style-type: none"> <li>○ 30 min. Monday 8am <a href="#">Join Microsoft Teams Meeting</a></li> <li>○ 1.5 hours – tried again with prodigy, read a bit on topic analysis, dbc table. No pagination, can't sort, slower</li> </ul> </li> <li>• 2 hours - Tuesday <ul style="list-style-type: none"> <li>○ Meet with Brandon 7:30</li> <li>○ TODO: Tutorials, analyze n-grams, constrain to better examples</li> <li>○ Peer review for a “resource review” that Sophie preparing for the CDI Usability Collaboration Area</li> </ul> </li> <li>• 2 hours - Wednesday <ul style="list-style-type: none"> <li>○ Prodigy tutorials <a href="https://prodi.gy/docs/named-entity-recognition">https://prodi.gy/docs/named-entity-recognition</a> (Food Ingredient entities)</li> </ul> </li> <li>• 1.5 hour - Thursday <ul style="list-style-type: none"> <li>○ Continued Prodigy tutorial</li> </ul> </li> <li>• 1 hour – Friday <ul style="list-style-type: none"> <li>○ Continued Prodigy tutorial</li> </ul> </li> <li>• 1 hour – Saturday <ul style="list-style-type: none"> <li>○ N-gram analysis of text</li> <li>○ <a href="https://github.com/PortfolioSQA/USGS_Catalog_Dash/blob/master/ngram_text_analysis.ipynb">https://github.com/PortfolioSQA/USGS_Catalog_Dash/blob/master/ngram_text_analysis.ipynb</a></li> </ul> </li> </ul>
1/18	<ul style="list-style-type: none"> <li>• TOTAL 10 hours</li> </ul>

	<p>Martin Luther King Jr Holiday Monday</p> <ul style="list-style-type: none"> <li>0.5 hour - Tuesday <ul style="list-style-type: none"> <li>Meet with Brandon 7:30 – 8:00</li> <li>TODO: methods NER: try water balance, transport model, etc. separate and then try DE, Lin reg, random forests, etc. (may have to use regex for that), Graphical interface - Keep track of articles that may not be models, New Github directory with files and commands</li> </ul> </li> <li>2.5 hours – Friday <ul style="list-style-type: none"> <li>New model for methods (1) Used all methods</li> </ul> </li> <li>4 hours - Saturday <ul style="list-style-type: none"> <li>New model for methods (2) Used geologic models (not stats and ML methods)</li> </ul> </li> <li>3 hours – Sunday <ul style="list-style-type: none"> <li>Model 3 + notes</li> </ul> </li> </ul>
1/25	<ul style="list-style-type: none"> <li>TOTAL 10.5 HOURS</li> <li>1.5 hours Monday <ul style="list-style-type: none"> <li>30 min Monday 8am <a href="#">Join Microsoft Teams Meeting</a></li> <li>Notes, github commit, ml model</li> </ul> </li> <li>0.5 hours Tuesday <ul style="list-style-type: none"> <li>Meet with Brandon 7:30-8:00</li> <li>TODO: <ul style="list-style-type: none"> <li><del>Get another 20 texts for testing (abstracts)</del></li> <li><del>Train again with 'bad examples' see if make gold improves</del></li> <li><del>Text (jsonl) files for entity seeds (separate models)</del></li> <li><del>Train for words before 'model' using verb 'model' as bad examples</del></li> <li><del>Research existing work for ML NER, articles and/or list for seed terms</del> <ul style="list-style-type: none"> <li>Goal: blog post (search for other blogs pertaining to the topic)</li> </ul> </li> </ul> </li> </ul> </li> <li>1.5 hours Friday <ul style="list-style-type: none"> <li>Gather 20 new texts from Model Catalog</li> </ul> </li> <li>4 hours Saturday <ul style="list-style-type: none"> <li>Train new model – geological terms with 'bad examples'</li> <li>Train new model – words that are prior to 'model'</li> </ul> </li> <li>2 hours Sunday <ul style="list-style-type: none"> <li>Anaconda environment broke after update. Spent an hour trying to fix it. (Didn't count this as hours). Took a break from Prodigy to research.</li> <li>Research available ML model detection algorithms, seed lists</li> </ul> </li> </ul>
2/1	<p>TOTAL __ HOURS</p> <ul style="list-style-type: none"> <li>1.5 hours Monday <a href="#">Join Microsoft Teams Meeting 9am</a></li> <li>__ Tuesday <ul style="list-style-type: none"> <li>Meet with Brandon</li> <li>Train new model – geological terms with 'bad examples' – REDO!! Try with and without patterns</li> </ul> </li> </ul>
2/8	Monday 8am <a href="#">Join Microsoft Teams Meeting</a>
2/15	Monday <a href="#">Join Microsoft Teams Meeting 9am</a>
2/22	Monday 8am <a href="#">Join Microsoft Teams Meeting</a>
3/1	Monday <a href="#">Join Microsoft Teams Meeting 9am</a>
3/8	Monday 8am <a href="#">Join Microsoft Teams Meeting</a>
3/15	Monday <a href="#">Join Microsoft Teams Meeting 9am</a>
3/22	Monday 8am <a href="#">Join Microsoft Teams Meeting</a>
3/29	Monday <a href="#">Join Microsoft Teams Meeting 9am</a>

4/5	Monday 8am <a href="#">Join Microsoft Teams Meeting</a>
4/12	Monday <a href="#">Join Microsoft Teams Meeting 9am</a>
4/19	Monday 8am <a href="#">Join Microsoft Teams Meeting</a>
4/26	Monday <a href="#">Join Microsoft Teams Meeting 9am</a>
5/3	