# ECS 193AB Winter/Spring2017

## Template-based Data Extraction



| color | otherDesignation | brandName | vintage.year |
|---|---|---|---|
| Red | BORDEAUX ROUGE | Boyer Frercs | 1955 |
| Red | CHIANTI | Cimamori | |
| Red | | Boyer I'rem | 1955 |
| Red | | Puisseguin | 1955 |
| Red | | Spanish Roja | 1955 |

The library is interested in extracting data from scanned data from their historical collections. These data can thought of as tabular data. However that is not how they are formatted on the scanned documents, nor do standard OCR reliably capture what structure there is. Current structured extraction methodologies, like tabula (http://tabula.technology/), do not work well with data that is only semi-structured. Additionally, most extraction tools are separated from the OCR step. This can have implications, for example; when you know you are looking only for prices, (eg /\$(\d|[.])+/)

We would like an application or library that would allow the specification of search templates, these templates could be used across multiple scanned documents. Closely matching templates could be identified for further refinement or curation. Templates would act on multiple lines, and might associate multiple locations on a document into a single datum.

As ideas, we have considered approximate regular expression searches like TRE (https://en.wikipedia.org/wiki/TRE_(computing)), and implementing the functions as a postgresql extension, using tesseract for OCR, possibly rescanning to better match parts like prices. We've also considered a javascript library

Quinn Hart <qjhart@ucdavis.edu (mailto:qjhart@ucdavis.edu)> – Digital Applications Manager – UCD Library