# NLP Project Tasks

1. Read in scripts + scrape
   - try PDFMiner
   - PyPDF2

   -make subdata frames

2. Explore each script -pre-load
   - # of words
   - # of main characters ⎫ store these somewher
   - # of words per character ⎬ -Pos → then is upper?
   - # of settings
   - # unsaid cues

   → unsaid cues?
   ↳ put unsaid cues in a col
   ↳ put character names in a col
   ↳ put counts of char speaks
   talking is indented

(2.A) Add genre field to table
   → import table

3. Subsets
   - Text per character
   - locations

   — or extract based on regex → then do PoS

4. Clean text — clean, then PoS → more clean
   - find ways to remove what we don't want ⎫ use NLTK
   - Tokenize → Lemmatize
   - Remove Stopwords

   - Remove "noise"
   - extract key parts? → identify what is considered key?

Sentiment Analysis — look @ HW 4
    → use ANEW to score words
        → any customization needed?

    → calculate scores                    → split by actor/role?
   → split by genre                              main
   → split by winner loser

6. Topic Modeling        ☐ TF-IDF?                        words
    → put words in BoW + use sklearn    ⇒ share ___ in
                                                    each topic
   → which method do we want to use?
            ↳ NMF?
   → split by genre?
   → split by winner loser?
7. Semantic Similarity ___ —Lable

   ↓ ↓                → winners vs losers
   ↓ ↓ ?, these similar    the → genres?
— — — — — — — — — — — — — — — —

8. K-Means Clustering
        → KNN


Visuals — Screen shob of diff formats of PDFs
    — word cloud for — genres
              — winners vs loser
        graphs → — sentiment across genre grap
                — sentiment winner vs losers

— — — — — — — — — — — — — —

Other — If time Classification model
        → use 4 unused scripts + maybe current films that
                        could be nominated for
                        2020 oscars?