

Alan Perez  
AXP200075  
CS4395.001

## Portfolio: Finding or Building a Corpus

For the knowledge base, I extracted the top terms from each individual site and selected the most appropriate ones for the artist's album/music. I manually looked up the definition of each term and any relevant information regarding the artist/celebrity. I specifically looked at the top frequent terms and created a hardcoded dictionary for the knowledge base. My initial approach was to do it by automatically extracting the definition from the existing information within the URLs but ran into issues trying to implement it this way. Considering the small amount of data, I believe manual entry was the best option since it allows more control over the content.

I made the mistake of choosing a topic that didn't have enough data or rich text for the corpus, had I gone with a topic that was a lot broader then I'd have had sufficient data

```
{'http://50cent.com'}
{'http://www.thesmokinggun.com/archive/50cent1.html'}
{'https://web.archive.org/web/20070309174039/http://www.thesmokinggun.com/archive/50cent1.html'}
{'http://www.daveydc.com/interview50cent.html'}
{'https://web.archive.org/web/20070214131915/http://www.daveydc.com/interview50cent.html'}
{'http://www.50cent.com/'}
{'https://web.archive.org/web/19991124193831/http://www.50cent.com/'}
{'http://www.dubcnn.com/interviews/50cent/'}
{'https://web.archive.org/web/20070513011450/http://www.dubcnn.com/interviews/50cent/'}
{'https://www.rollingstone.com/artists/50cent/albums/album/301556/review/6067729/get_rich_or_die_tryin'}
{'https://web.archive.org/web/20090410035524/http://www.rollingstone.com/artists/50cent/albums/album/301556/review/6067729/get_rich_or_die_tryin'}
{'http://www.rockonthenet.com/artists-f/50cent.htm'}
{'https://web.archive.org/web/20070428234536/http://www.rockonthenet.com/artists-f/50cent.htm'}
{'https://web.archive.org/web/20070626212626/https://www.rollingstone.com/artists/50cent/albums/album/7072060/review/7045740/the_massacre'}
{'https://www.rollingstone.com/artists/50cent/albums/album/7072060/review/7045740/the_massacre'}
[['('cent', 0.2), ('upstoesign', 0.2), ('latest', 0.2), ('updates', 0.2), ('centcopyright', 0.2)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']
[['('stars', 0.016887816640562123), ('track', 0.014475271411338963), ('album', 0.013268998793727383), ('cent', 0.010856453558504222), ('massacre', 0.009650180940892641), ('get', 0.009650180940892641), ('shop', 0.009650180940892641), ('centcopyright', 0.009650180940892641), ('tryin', 0.009650180940892641)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']
[['('cent', 0.2000482509047045), ('upstoesign', 0.2000482509047045), ('latest', 0.2000482509047045), ('updates', 0.2000482509047045), ('centcopyright', 0.2000482509047045), ('stars', 0.003377563293124246), ('track', 0.003377563293124246), ('album', 0.003377563293124246), ('cent', 0.003377563293124246), ('tryin', 0.003377563293124246)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']
[['('acceptable', 0.018181818181818182), ('could', 0.09095295536791315), ('found', 0.09095295536791315), ('appropriate', 0.09090909090909091), ('representation', 0.09090909090909091), ('requested', 0.09090909090909091), ('resonance', 0.09090909090909091), ('cent', 0.09090909090909091), ('tryin', 0.09090909090909091)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']
[['('cent', 0.043573128097359984), ('like', 0.031227385374834665), ('got', 0.02033405949745823), ('album', 0.015250719865303939), ('record', 0.014524344177363407), ('know', 0.013071895424836602), ('get', 0.009440940782741677), ('centcopyright', 0.009440940782741677), ('tryin', 0.009440940782741677)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']
[['('cent', 0.05135549595539336), ('top', 0.03162202382137466), ('chart', 0.03125), ('artist', 0.01971969343885415), ('best', 0.01971780225027996), ('lp', 0.01636904761904762), ('hot', 0.015997293985245548), ('rap', 0.015625548136881182), ('centcopyright', 0.015625548136881182), ('tryin', 0.015625548136881182)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']
[['('account', 0.03140853146853147), ('news', 0.024475525420194832), ('rs', 0.02447552447556591), ('digital', 0.020984224109224108), ('manage', 0.02097902897902898), ('rolling', 0.017499427655801966), ('stone', 0.017499427655801966), ('centcopyright', 0.017499427655801966), ('tryin', 0.017499427655801966)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']
[['('cent', 0.056206509519136236), ('girl', 0.05386416862047927), ('hold', 0.051522254336218085), ('let', 0.05152225128988892), ('bag', 0.05152224824355972), ('document', 0.02576112412177986), ('http', 0.02187728337236534), ('centcopyright', 0.02187728337236534), ('tryin', 0.02187728337236534)]]
Selected top 10 terms: ['50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin']

Knowledge Base: 50 cent:Curtis James Jackson III (born July 6, 1975), known professionally as 50 Cent, is an American rapper, actor, television producer, and businessman.
Knowledge Base: massacre:The Massacre is the second studio album by American rapper 50 Cent.
Knowledge Base: g unit:G-Unit (short for Guerilla Unit) was an American hip hop group formed by longtime friends and East Coast rappers 50 Cent, Tony Yayo, and Lloyd Banks.
Knowledge Base: hook:A hook is a musical idea, often a short riff, passage, or phrase, that is used in popular music to make a song appealing and to catch the ear of the listener.
Knowledge Base: club:an association or organization dedicated to a particular interest or activity.
```

Top 10 terms: '50 cent', 'massacre', 'g unit', 'hook', 'club', 'released', 'record', 'stone', 'albums', 'tryin'

Example of Chatbot Dialog:

**User:** Can you finish the rest of this lyric? “Hate it or love it underdogs on top:

**Chatbot:** “And I'm gon' shine, homie, until my heart stop”

Alan Perez  
AXP200075  
CS4395.001

**User: What is 50 cent's best selling album?**

**Chatbot: "Get Rich or Die Tryin"**

**User: What group was 50 cent a part of in the 2000s?**

**Chatbot: g unit**