

## **Data Appendix:**

### **episode1\_clean**

- This is a cleaned data file of Episode 1, The Phantom Menace, movie script
- The cleaning code turns the original SW01\_txt data file into one column containing every word in the script, cleaned of any punctuation (excluding spaces), multiple spaces, capitalization, numbers, and screenplay markers
  - Each new line in the original script was replaced with a space in order to put the entire new cleaned script into one column that we could run our analysis on
  - Unit of Observation: text from script of Episode 1
- On this cleaned file we performed both the VADER sentiment analysis and NRC lexicon analysis

### **episode2\_clean**

- This is a cleaned data file of Episode 2, Attack of The Clones, movie script
- The cleaning code turns the original SW02\_txt data file into one column containing every word in the script, cleaned of any punctuation (excluding spaces), multiple spaces, capitalization, numbers, and screenplay markers
  - Each new line in the original script was replaced with a space in order to put the entire new cleaned script into one column that we could run our analysis on
  - Unit of Observation: text from script of Episode 2
- On this cleaned file we performed both the VADER sentiment analysis and NRC lexicon analysis

### **episode3\_clean**

- This is a cleaned data file of Episode 3, Revenge of the Sith, movie script
- The cleaning code turns the original SW\_txt data file into one column containing every word in the script, cleaned of any punctuation (excluding spaces), multiple spaces, capitalization, numbers, and screenplay markers
  - Each new line in the original script was replaced with a space in order to put the entire new cleaned script into one column that we could run our analysis on
  - Unit of Observation: text from script of Episode 3
- On this cleaned file we performed both the VADER sentiment analysis and NRC lexicon analysis

### **episode4\_clean**

- This is a cleaned data file of Episode 4, A New Hope, movie script
- The cleaning code turns the original SW04\_txt data file into one column containing every word in the script, cleaned of any punctuation (excluding spaces), multiple spaces, capitalization, numbers, and screenplay markers
  - Each new line in the original script was replaced with a space in order to put the entire new cleaned script into one column that we could run our analysis on
  - Unit of Observation: text from script of Episode 4
- On this cleaned file we performed both the VADER sentiment analysis and NRC lexicon analysis

### **episode5\_clean**

- This is a cleaned data file of Episode 5, The Empire Strikes Back, movie script

- The cleaning code turns the original SW05\_txt data file into one column containing every word in the script, cleaned of any punctuation (excluding spaces), multiple spaces, capitalization, numbers, and screenplay markers
  - Each new line in the original script was replaced with a space in order to put the entire new cleaned script into one column that we could run our analysis on
  - Unit of Observation: text from script of Episode 5
- On this cleaned file we performed both the VADER sentiment analysis and NRC lexicon analysis

#### **episode6\_clean**

- This is a cleaned data file of Episode 6, The Return of the Jedi, movie script
- The cleaning code turns the original SW06\_txt data file into one column containing every word in the script, cleaned of any punctuation (excluding spaces), multiple spaces, capitalization, numbers, and screenplay markers
  - Each new line in the original script was replaced with a space in order to put the entire new cleaned script into one column that we could run our analysis on
  - Unit of Observation: text from script of Episode 6
- On this cleaned file we performed both the VADER sentiment analysis and NRC lexicon analysis

#### **episode7\_clean**

- This is a cleaned data file of Episode 7, The Force Awakens, movie script
- The cleaning code turns the original episode7.txt data file into one column containing every word in the script, cleaned of any punctuation (excluding spaces), multiple spaces, capitalization, numbers, and screenplay markers
  - Each new line in the original script was replaced with a space in order to put the entire new cleaned script into one column that we could run our analysis on
  - Unit of Observation: text from script of Episode 7
- On this cleaned file we performed both the VADER sentiment analysis and NRC lexicon analysis