Introduction to Sentiment Analysis

How to get meaning out of text?

"I bought a Canon G12 camera six months ago. I simply love it. The picture quality is amazing. The battery life is also long. However, my wife thinks its too heavy for her." (1) I bought a Canon G12 camera six months ago.

- (2) I simply love it. (+ve)
- (3) The picture quality is amazing. (+ve)
- (4) The battery life is also long. (+ve)
- (5) However, my wife thinks it is too heavy for her. (-ve)

What does it mean to say an opinion?

- An opinion has two key components: a **target** g and a **sentiment** s on the target, (g, s),
 - where <u>g</u> can be any entity or aspect of the entity on which an opinion has been expressed, and s is a positive, negative, or neutral sentiment or a numeric sentiment rating.
- Positive, negative, and neutral are called sentiment or opinion orientations.

A very simple approach to get sentiment scores...

What if we assign the words a sentiment score. Positive (1), neutral (0) and Negative (-1). Would that make sense?

"The(0) picture(1) quality(0) is(0) amazing(1)." → 2
(We can say this text is positive)

What about punctuations, context, boosters, dampers...

"The(0) picture(1) quality(0) is(0) not(0) amazing(1)." → 2 (We can say this text is positive???)

So this simple algorithm does not work, we can not just look at individual words, (emoticons) and decide the overall opinion (valence) of the sentence.

This is the problem of sentiment analysis.

The Problem of Sentiment Analysis

- Unlike factual information, sentiment and opinion have an important characteristic, namely, they are **subjective**.

- 1. Firstly of all, different people may have different experiences and thus different opinions.
- 2. Different people may see the same thing in different ways because everything has two sides.
- 3. Different people may have different interests and/or different ideologies.

Sentiment

<u>Sentiment</u> is the underlying feeling, attitude, evaluation, or emotion associated with an opinion. It is represented as a triple,

(y, o, i)

where **y** is the <u>type of the sentiment</u>, o is the orientation of the sentiment, and i is the intensity of the sentiment.

Because sentiment classification is a text classification problem, any existing supervised learning method can be directly applied, such as naïve Bayes classification or support vector machines

VADER (Valence Aware Dictionary and sEntiment Reasonor)

A lexicon and **rule-based sentiment analysis** tool specifically designed for social media and short text.

Built to recognize not just standard words, but also slang, emoticons, acronyms, and punctuation commonly used online.

It uses a <u>predefined</u> dictionary of words and phrases with sentiment intensity scores assigned to each one.

How VADER works?

1. Preprocessing:

- Tokenizing the input text.
- Identifies known words, slang, emojis, punctuation, etc.

2. Lexicon Scoring:

- Each word has a pre-assigned valence score from the VADER lexicon.
- Adjusts scores for **degree modifiers** (e.g., "very", "extremely").

How VADER works?

- 3. Rule-based Adjustments:
 - Emphasis modifiers:
 - Capitalization, punctuation (e.g., "!!!").
 - Negation handling (e.g., "not good").
 - Emoji and slang interpretation.
- 4. Sentiment Calculation:
 - Outputs four scores:
 - Positive, Negative, Neutral, and Compound (normalized overall score between -1 and 1).

Main Functions in VADER Sentiment Analysis

SentimentIntensityAnalyzer(): Initializes the analyzer using the built-in lexicon.

.polarity_scores(text): Returns sentiment scores (positive, negative, neutral, compound) for the given text.

.make_lex_dict() (internal use): Converts a VADER lexicon file into a dictionary of word scores.

Output of polarity_scores()

The polarity_scores() function returns a dictionary with four keys representing sentiment metrics.

- pos: Proportion of text that conveys positive sentiment (value between 0 and 1).
- neu: Proportion of text that is neutral in sentiment (value between 0 and 1).
- neg: Proportion of text that conveys negative sentiment (value between 0 and 1).
- compound: A normalized score from -1 (most negative) to +1 (most positive) that summarizes the overall sentiment.

The compound score is often used for final sentiment classification (e.g., > 0.05 is positive, < -0.05 is negative).

Thank You!!