# ANALYSIS ON PUNCTUATIONS IN ONLINE REVIEWS

Tongzhou Wang    CS C281A Final Project

## OVERVIEW

- Explore information in punctuations in Steam online reviews.

- Two different directed graphical models are analyzed using a publicly available dataset.

- Gather insights about punctuations and structures in online reviews.

## MOTIVATION

Many text analyses focus mainly on the words. However, punctuations in texts obviously also carry huge amount of information. In fact, using punctuations as features has advantages including:

1. insensitive to typos,

2. less sensitive to language used,

3. expressing strong emotions, and

4. limited in types, resulting in much smaller sized models.

Want to investigate the relations among punctuations in Steam online game review dataset.

Aforementioned advantages of punctuation features are particularly useful in this dataset because online game reviews

1. tend to have many typos,

2. can be in many different languages,

3. sometimes express strong emotions, e.g. anger, excitement, etc.

## METHODS

- Preprocess data s.t. each sentence ends in proper punctuation.

- Separate the entire review dataset into positive ones and negative ones.

- Fit two models on positive data and negative data respectively.

- Analyze different punctuation relations in positive and negative data.

## MODEL

Sentences in reviews usually serve for different purposes, e.g. introduction, listing pros & cons, scoring, etc.

Directed model with sentence type:

1. Each sentence are considered as having one of $m$ sentence types.

2. Sentences types in a review are from a Markov chain $(\pi_s, A_s)$.

3. For each sentences, given its sentence type $s_i$, punctuations are from Markov chain $(\pi_p^{s_i}, A_p^{s_i})$.
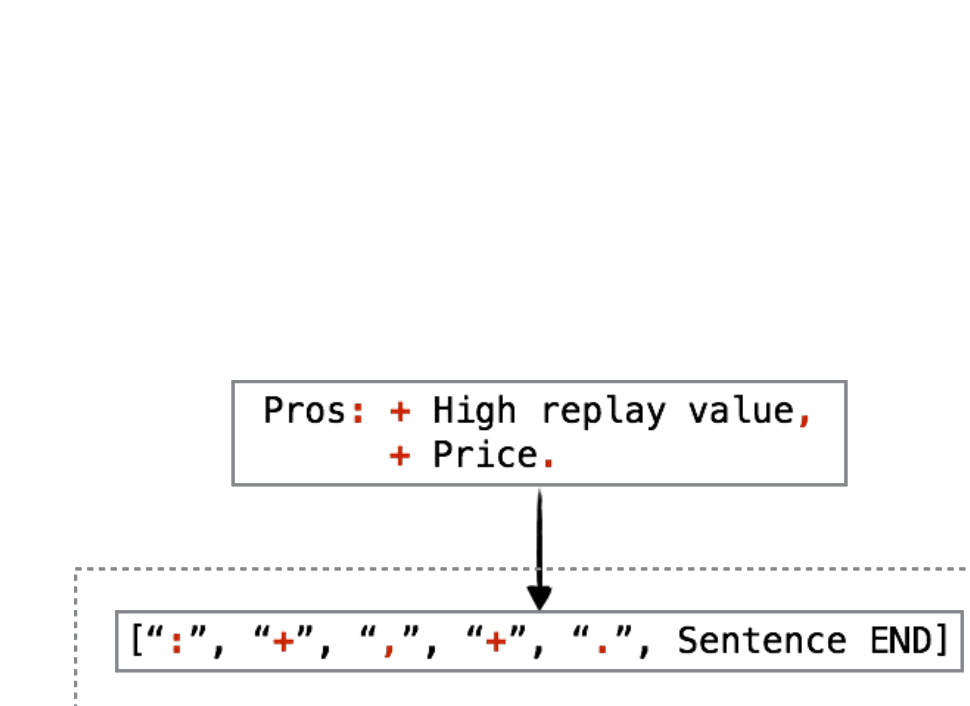
This model

- Essentially an HMM but with $P(\text{punctuations}|\text{sentence type})$ parametrized by a Markov chain

- Can use EM algorithm to approximately find MLEs.

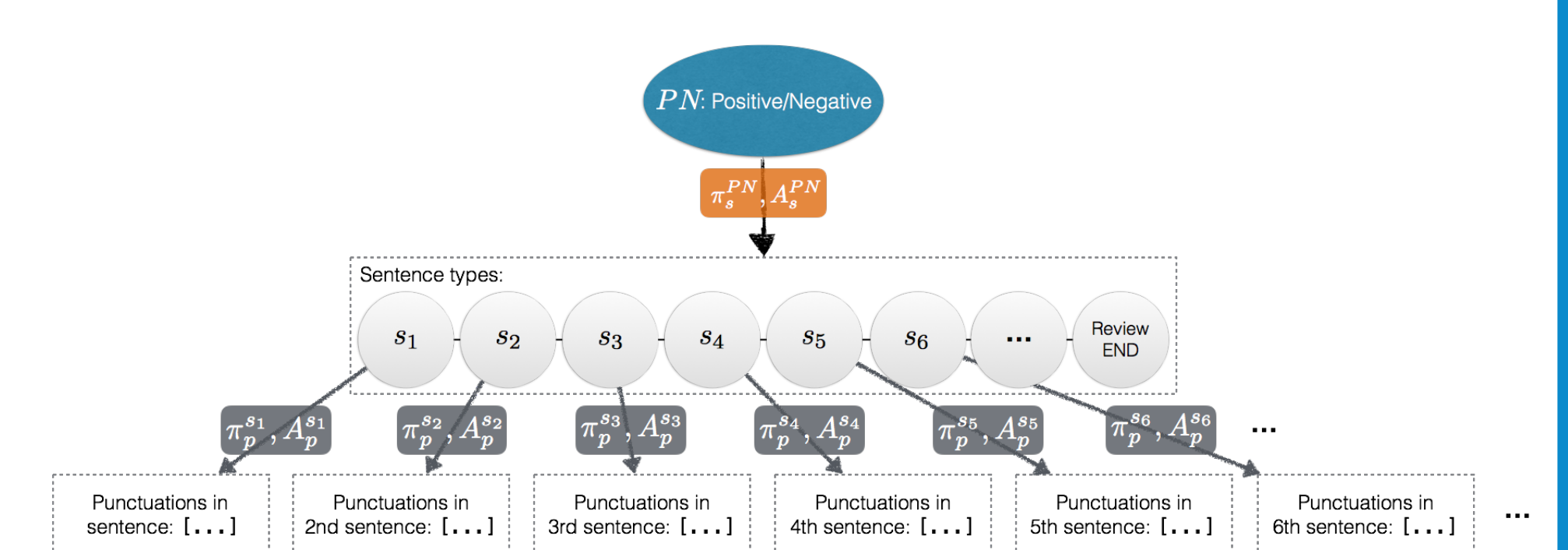- Can use Forward-backward algorithm for hidden variable marginals in EM.
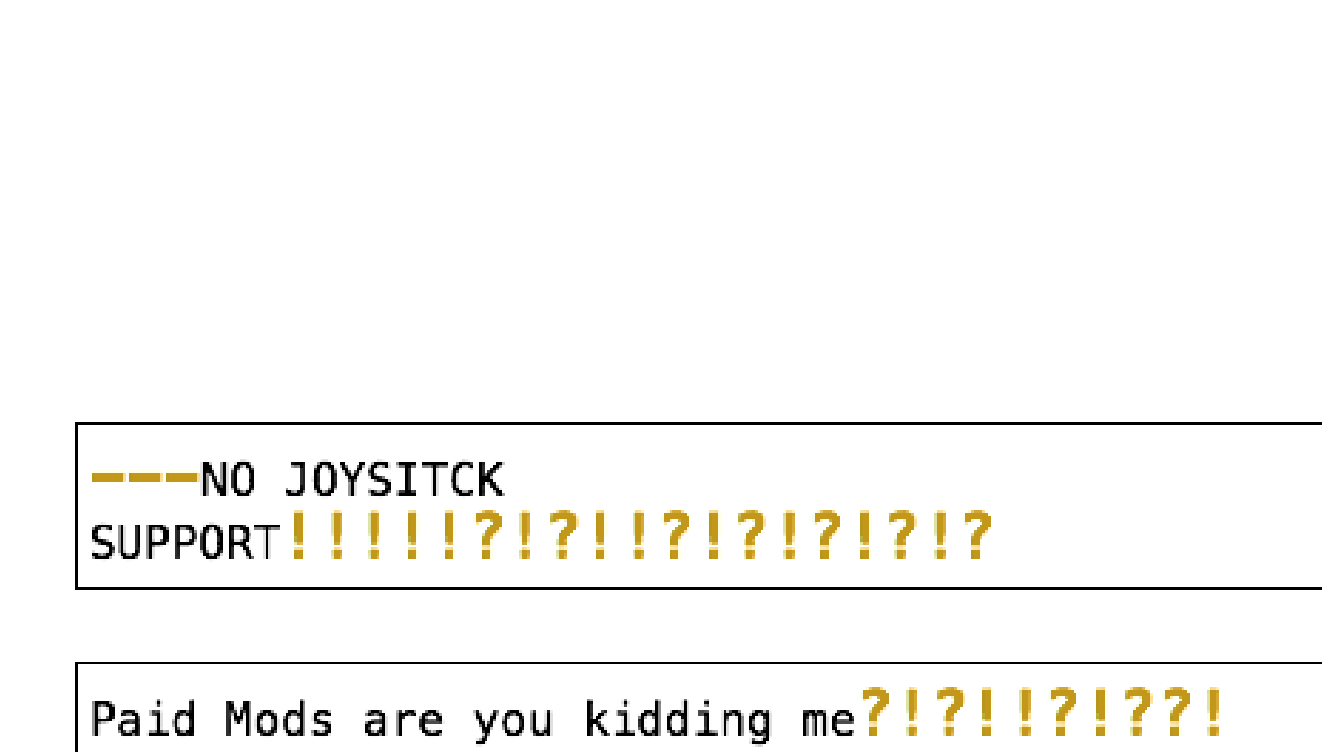
## VISUALIZATIONS AND RESULTS



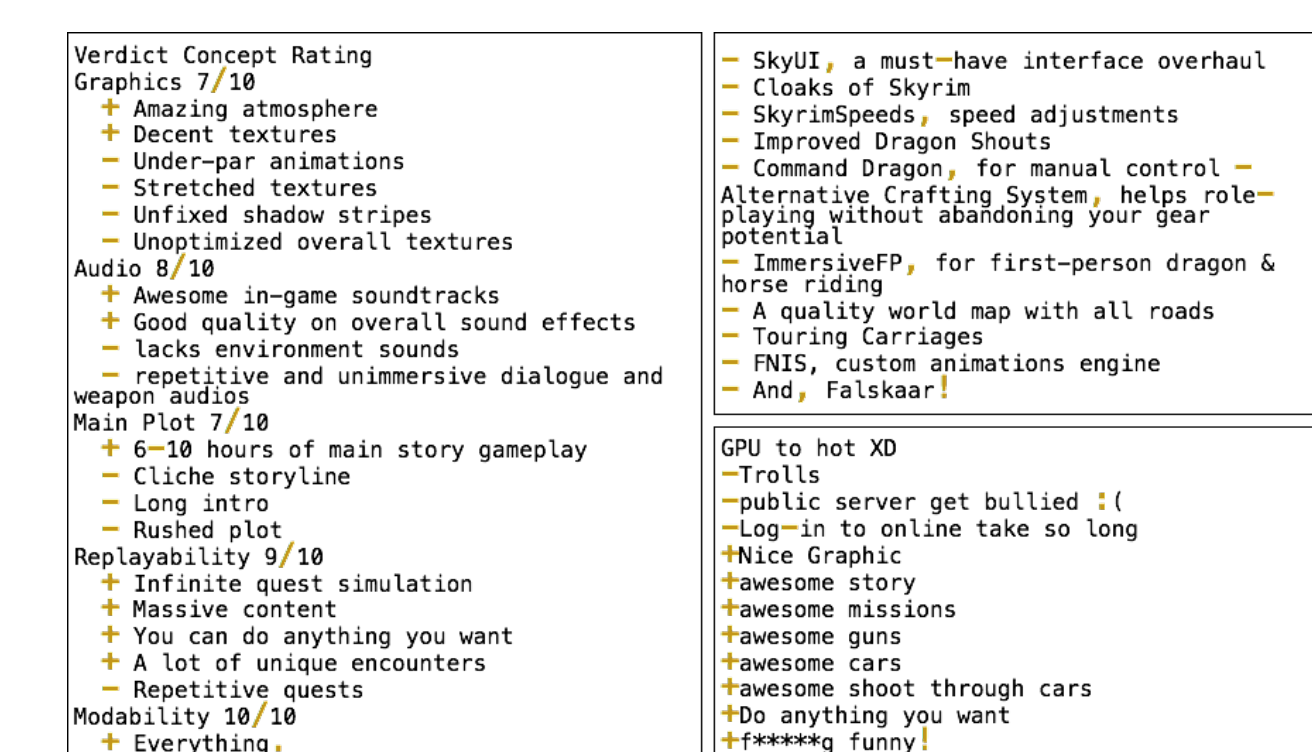**(a)** Hypothetical sentence types in a review example



**(b)** Punctuation data example of a sentence



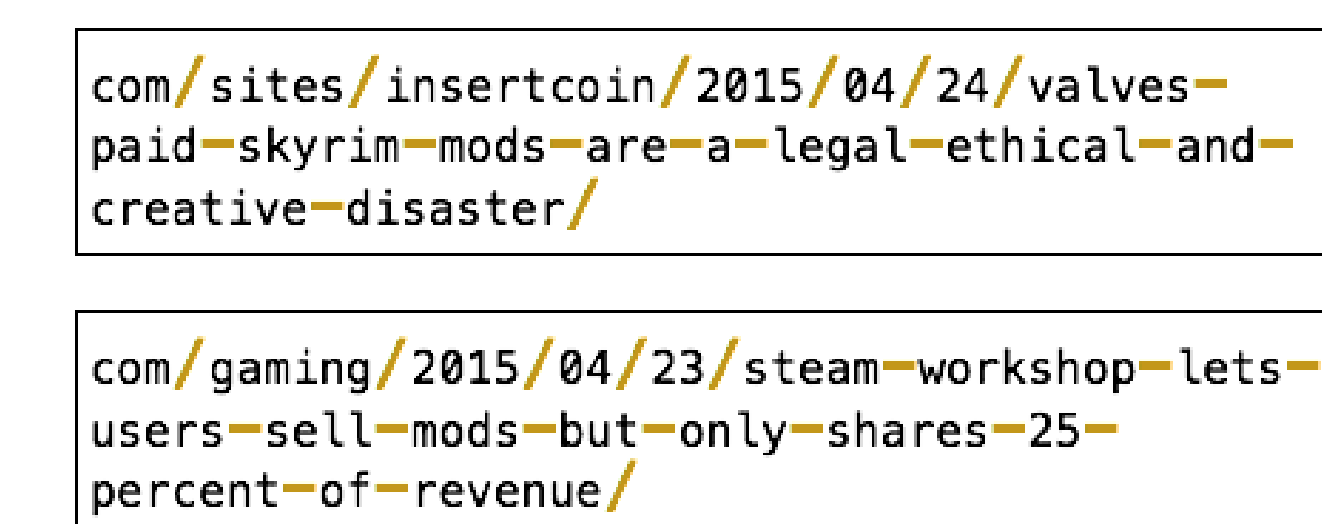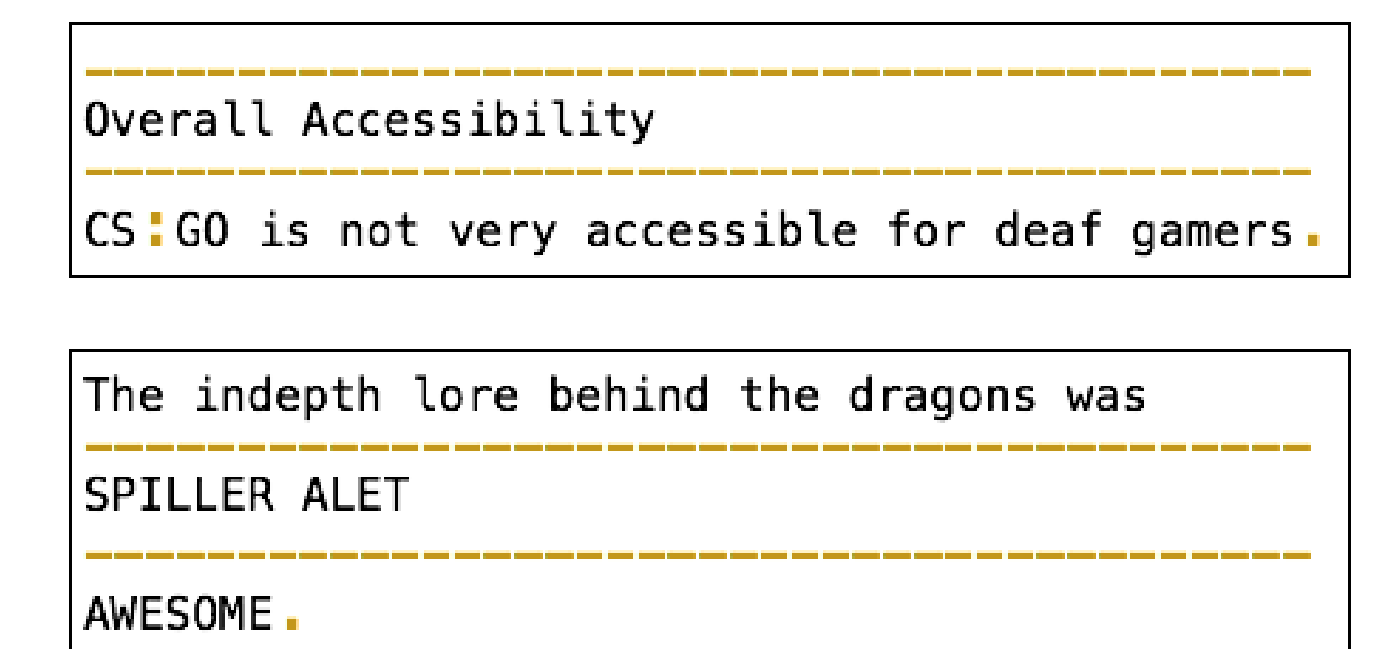**(c)** Model visualization



**(d)** Sentence type **1** of **negative** reviews examples



**(e)** Sentence type **2** of **positive** reviews examples



**(f)** Sentence type **4** of **positive** reviews examples



**(g)** Sentence type **0** of **negative** reviews examples



**(h)** Sentence type **1** of **positive** reviews examples



**(i)** Sentence type **5** of **positive** reviews examples

## OBSERVATIONS AND ANALYSIS

- Punctuation data are very informative.
- Sentence types successfully capture punctuation patterns:
  - **Negative** type **1**: pattern of anger.
  - **Positive** type **2**: pros and cons.
  - **Positive** type **4**: formatting pattern.

- Transition probability also gives interesting insights:
  - **Negative** reviews have high prob. from type **0** (url pattern) to REVIEW END.
  - **Positive** reviews have high prob. from type **1** (!!!) to type **5** (more !!!!!!).

## FUTURE WORK

1. Positive/negative prediction using the two fitted models.

2. Experiment and analyze best choice for $m$ number of sentence types.

## RESOURCES AND REFERENCES

[1] Mulholland, *Steam Review Datasets*, https://github.com/mulhod/steam_reviews.

[2] Gu, Leon, *EM and HMM*.