

Determining Environmental Context from Fictional Narratives

Abstract

Experiencing a narrative can be made more enjoyable and powerful by increasing the level of the immersion in the story. This can be achieved through contextual cues, such as sounds and images. In order to provide these cues automatically, we attempt to detect parts of the environment described in a fictional narrative. We present three machine learning approaches to this problem, compared to a rule-engineered baseline, and evaluate them over a dataset of fictional texts. In addition, we analyse the performance of this system over the output of a text-to-speech system taking the same narratives as input. This audio dataset is made available. The system is able to recognise environmental context and respond appropriately.

1 Introduction

Non-verbal cues are important to improving the experience of a narrative. For example, adding sound changes mood, and is useful for immersion and positive narrative experiences (Ermi and Mäyrä, 2005; Madden and Logan, 2009; Huiberts, 2010). An immersive environment is most engaging when participants believe their actions affect the environment; reliable and responsive cues such as sounds are important for this (Bobick et al., 1999). However, providing such cues can be an expensive manual process, and many narratives do not have accompanying contextual media such as background sounds or images. An automatic approach could be taken to building such context. The research question we attempt to answer is: how can we automatically identify in narratives environments and events that can be associated with external cues?

Taking narratives as input, we attempt to determine what the appropriate environmental component of the context is, which can then be used to provide cues (such as ambient sounds). There is a diverse range of potential environments from which cues can be generated, even in a constrained corpus of narratives. In addition, effects (e.g. sound or audio cues) are required for surprise events described in a story, like an approaching horse or a thunderstorm. This is similar to streamed topic extraction (Allan, 2002; Preotiuc-Pietro et al., 2012). We determine a corpus, frame the task by giving a specific set of environments for which cues are available, and then evaluate system performance over this dataset. Finally, we create a new language resource of audio recordings and automatic transcriptions, used to evaluate the system on input closer to that it might receive when operating in real-time.

2 Method

2.1 Dataset

Our ideal requirements are that a dataset be structured as instances, each being a focused passage of text, with most containing descriptions of environments or events. For this dataset, we used paragraphs from a set of fantasy role-playing books. Each paragraph is set in a distinct location, each possibly combined with the description of an event, making manual environment annotation simple and distinct. In comparison, the bounds of passages of text relevant to a particular cue are harder to determine in running narrative text from e.g. a fictional novel. The set of potential environmental items in this genre is:

Environments: Mountain; Hill; Forest; Swamp; Meadow; Road; Town; Crowd; Tavern; Underground.

Ambience: Windy; Blizzard; Rain; Lightning; Stream; River; Campfire;

Night.

Events: Trotting horse; Galloping horse; Thunder

Paragraphs were labelled with one or more of these labels by a human annotator. In total, 203 paragraphs were labelled, with 13884 words.

2.2 Spoken dataset

An eventual use of this system is to provide automatic sound effects in real-time for a story that is read out loud by a human narrator. As a result, it should operate well on the output of a speech recognition system. Paragraphs were read by an English native speaker, and recorded. A speech recognition system (Lamere et al., 2003) then interpreted these readings and generated a textual version of each paragraph. The labels used in the text input corpus were then associated with these outputs. This constitutes the transcribed dataset.

2.3 Baseline

As a baseline, we use a gazetteer trigger words that match the name of the cue. If a word occurs in a candidate passage, then that cue is triggered. For example, if there is a cue named “swamp”, the passage “He marches through the swamp, sweating from his brow” will trigger that cue using the baseline system.

2.4 Features and Classification

We take a variety of approaches. Firstly, we try a bag-of-words representation using a multinomial Bayes classifier. In addition, we represent paragraphs as average vectors in n-dimensional space by taking the cosine product of embeddings of words in the paragraph, and then learn a binary SVM for each cue based on these representations. We used Collobert and Weston[] embeddings.

The genre used to generate these had a significant effect on result quality. For example, using embeddings generated from two billion Google News articles, the top-ranking non-punctuation match for the sentence “*The going underfoot becomes muddier, until eventually you reach an area where bulrushes tower over your head.*” was the phrase *Maria Sharapova* – not an intuitively relevant result. News is intrinsically constrained in topic and style, and out-of-domain for fictional narratives. Therefore, we used embeddings learned from the multi-topic Brown corpus, which

includes large amounts of narrative and fictional text.

LD feature extraction (Lui and Baldwin, 2011) + nbayes; LD motivated by diversity in subsets of environments encountered per story (e.g. some stories set mostly in woodland, others in a town)

We note that false positives have a higher penalty than false negatives. It is disruptive to receive the wrong cue. For example, hearing battle noises when one should hear a babbling brook (false positive) is detrimental, and worse than hearing nothing (false negative). In order to bias classifications in this direction, we experiment with the SVM Cost parameter (?) and also evaluate with both F1 and F2.

3 Results

Our evaluation is using precision and recall. This is not a plain multi-way classification task; spurious cues and missed cues are both possible, as are multiple cues per text passage.

4 Related Work

ML doc classification (Sebastiani, 2002).

NN good at binary doc classification (Derczynski, 2006).

SVM doc classification (Isa et al., 2008).

5 Conclusion

References

- James Allan. 2002. Introduction to topic detection and tracking. In *Topic detection and tracking*, pages 1–16. Springer.
- Aaron F Bobick, Stephen S Intille, James W Davis, Freedom Baird, Claudio S Pinhanez, Lee W Campbell, Yuri A Ivanov, Arjan Schütte, and Andrew Wilson. 1999. The kidsroom: A perceptually-based interactive and immersive story environment. *Presence: Teleoperators and Virtual Environments*, 8(4):369–393.
- Peter Darvill-Evans, Steve Jackson, and Ian Livingstone. 1989. *Portal of Evil*. Puffin.
- Leon Derczynski. 2006. Machine learning techniques for document selection. Master’s thesis, University of Sheffield.
- Laura Ermi and Frans Mäyrä. 2005. Fundamental components of the gameplay experience: Analysing immersion. *Worlds in play: International perspectives on digital games research*, 37.
- Sander Huiberts. 2010. *Captivating sound the role of audio for immersion in computer games*. Ph.D. thesis, University of Portsmouth.

- Dino Isa, Lam Hong Lee, V Kallimani, and Rajprasad Rajkumar. 2008. Text document preprocessing with the Bayes formula for classification using the support vector machine. *Knowledge and Data Engineering, IEEE Transactions on*, 20(9):1264–1272.
- Steve Jackson. 1984. *Scorpion Swamp*. Puffin.
- Paul Lamere, Philip Kwok, William Walker, Evandro B Gouvêa, Rita Singh, Bhiksha Raj, and Peter Wolf. 2003. Design of the CMU Sphinx-4 decoder. In *Proc. INTERSPEECH*.
- Ian Livingstone. 1987. *Crypt of the Sorcerer*. Puffin.
- Marco Lui and Timothy Baldwin. 2011. Cross-domain feature selection for language identification. In *Proc. IJCNLP. ACL*.
- Neil Madden and Brian Logan. 2009. Collaborative narrative generation in persistent virtual environments. In *Proc. AAAI*.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Daniel Preotiuc-Pietro, Sina Samangooei, Trevor Cohn, Nicholas Gibbins, and Mahesan Niranjan. 2012. Trendminer: An architecture for real time analysis of social media text. In *Proceedings of the workshop on real-time analysis and mining of social streams*.
- Fabrizio Sebastiani. 2002. Machine learning in automated text categorization. *ACM computing surveys (CSUR)*, 34(1):1–47.
- Luke Sharp, Ian Livingstone, and Steve Jackson. 1989. *Fangs of Fury*. Puffin.