

Feature Engineering for SpaceEval

Seth Dworman

17 December 2014

1.0 Task Description

ISO-Space is a rich annotation markup language for capturing spatial and motion based information in natural language text. One of the overarching goals in creating the ISO-Space gold standard is to provide supervised data for machine learning of annotation language. That is, given unstructured natural language text, an algorithm should generate the ISO-Space markup automatically. This task by itself is rather difficult due to the diversity of spatial language, the richness of the ISO-Space specification, and the quality, nature, and distributions of the gold standard data. Simple bag-of-word approaches are not sufficient as spatial language is itself relatively sparse even in the gold standard. In addition, identification of spatial mentions (=text extents which should be classified into an ISO-Space tag) is not sufficient for generating the full ISO-Space markup, as each tag itself has its own attribute values, some of which are not trivial and must also be inferred by an algorithm. Lastly, once all mentions are determined, it is necessary to link them correctly via the ISO-Space links.

2.0 Feature Engineering

We propose that bag-of-word approaches will not suffice for learning the ISO-Space language or developing a suitable generalized algorithm. To get a reasonable algorithm, we suggest several features to supplement a bag-of-words approach. We follow the **statistical language model representation hypothesis** (known as LMRH), which claims that features that are informative for one NLP task are also informative on others (Huang, Ahuja, Downey, Yang, Guo, and Yates 2014). The features we choose to select for are part-of-speech tags (POS tags), named-entity recognition tags (NER tags), and semantic labels (Sparser labels) from Sparser, a handwritten shallow semantic based parser (McDonald).

3.0 Training Instances

The ISO-Space specification works at the sentential level, with only possibly METALINKs being intersential relations, as they are essentially coreference links limited to spatial mentions. We will ignore this for now and assume that an algorithm generating ISO-Space markup would operate with a single sentence as its input (with larger documents simply iterated over). Below we have reproduced an example sentence, *I have taken a boat up Amazon, crossed dirt roads over 15,000 ft Andean passes, luckily escaped bandits in southern Mexico, felt the wind and rain of of Patagonia, and dodged buses in the capital cities of nearly every country.*, already tokenized for us from the gold standard set. The column named **Label** represents the token’s ISO-Space tag value, with NONE indicating it is not a spatial mention. We interpret **Label** as the true class label y we wish to predict given a token. This example is taken from `Tokenized/RFC/BuenosAires.xml`.

Token	Label	Token	Label
I	SPATIAL_ENTITY	southern	NONE
have	NONE	Mexico	PLACE
taken	MOTION	,	NONE
a	MOTION_SIGNAL	felt	NONE
boat	MOTION_SIGNAL	the	NONE
up	MOTION_SIGNAL	wind	SPATIAL_ENTITY
the	NONE	and	NONE
Amazon	PATH	rain	SPATIAL_ENTITY
,	NONE	of	SPATIAL_SIGNAL
crossed	MOTION	Patagonia	PLACE
dirt	NONE	,	NONE
roads	PATH	and	NONE
over	MOTION_SIGNAL	dodged	MOTION
15,000	MEASURE	buses	SPATIAL_ENTITY
ft	MEASURE	in	SPATIAL_SIGNAL
Andean	NONE	the	NONE
passes	PATH	capital	NONE
,	NONE	cities	PLACE
luckily	NONE	of	SPATIAL_SIGNAL
escaped	NONMOTION_EVENT	nearly	NONE
bandits	NONE	every	NONE
in	SPATIAL_SIGNAL	country	PLACE
		.	NONE

We consider each token from a sentence to be a possible spatial mention, and augment a simple bag-of-words approach with POS tags, NER tags, and Sparser labels. POS tags are relatively straightforward and come from Stanford NLP’s POS tagger (Toutanova, Klein, Manning, and Singer 2003), which achieves a state-of-the-art performance using the Penn TreeBank POS labels. NER tags come from a small set of labels : **O(utside)**, **person**, **organization**, and **location**. We also make use of Stanford NLP’s NER tagger (Finkel, Grenager, and Manning 2005). Our intuition is that tokens which are identified as a NER label are more likely to be involved in some kind of spatial relation (especially NER **location**), given that the domain of the gold standard is travel blogs. Finally, we use Sparser edge labels, which come from Sparser, a shallow semantic parser (McDonald). Sparser’s labels are semantics based and form an open set. We believe that Sparser’s labels can be useful for finding spatial mentions which are places, paths, measurements, and spatial signals, whose correct identification can bootstrap finding other spatial mentions (e.g. spatial named entities). Below we reproduce part of the original sentence (the first two clauses), augmented with POS, NER, and primary Sparser labels. Recall that the **Label** column should be interpreted as the class label y we wish to predict given a token.

Token	POS	NER	Sparser	Label
I	PRP	O	SINGLE-CAPITALIZED-LETTER	SPATIAL__ENTITY
have	VBP	O	HAVE	NONE
taken	VBN	O	NONE	MOTION
a	DT	O	INDIVIDUAL	MOTION__SIGNAL
boat	NN	O	INDIVIDUAL	MOTION__SIGNAL
up	RP	O	UP	MOTION__SIGNAL
the	DT	O	NAME	NONE
Amazon	NNP	LOCATION	NAME	PATH
,	NONE	O	NONE	NONE
crossed	VBD	O	CROSS	MOTION
dirt	NN	O	DIRT	NONE
roads	NNS	O	PATH-TYPE	PATH
over	IN	O	OVER	MOTION__SIGNAL
15,000	CD	O	NONE	MEASURE
ft	JJ	O	NONE	MEASURE
Andean	JJ	O	NAME	NONE
passes	NNS	O	PASS-KIND	PATH
,	NONE	O	NONE	NONE
luckily	RB	O	LUCKILY	NONE
escaped	VBD	O	ESCAPE-EVENT	NONMOTION__EVENT
bandits	NNS	O	BANDIT	NONE
in	IN	O	IN	SPATIAL__SIGNAL
southern	JJ	O	DIRECTION	NONE
Mexico	NNP	LOCATION	COUNTRY	PLACE

We note that Sparser also provides other labels besides the primary label shown ; in fact all the field values of the Sparser **Edge** data structure are available as features, though some of these may not turn out to be useful.