

# Ranking Scalar Adjectives

Xiao Ling

lingxiao@seas.upenn.edu

## Abstract

Adjectives such as good, great, and excellent are similar in meaning but differ in intensity. Intensity ordering is useful in several NLP tasks, such as sentiment analysis and question answering. However this data is missing in most lexical resources such as dictionaries and WordNet. In this thesis we first review and replicate the state of the art in adjective ranking by de Melo and Bansal (2013) using monolingual corpora. Then we extend their approach in along two fronts: first by augmenting monolingual data with data derived from bilingual corpora for greater coverage, then we propose a variety of new approaches to the problem. Additionally we procure a new gold standard using mechanical turks and evaluate all methods on both gold standards. Notably, our best method achieves 75% pairwise accuracy and 0.61 Kendall's- $\tau$  on both gold standards, exceeding Bansal's method by a nontrivial amount, and approaches inter-annotator accuracy.

# 1 Introduction

This thesis investigates the efficacy of ranking scalar adjectives from corpus data using monolingual corpus and data derived from bilingual sources. We make three major contributions:

1. successfully reproduce the state of the art in adjective ranking using monolingual corpus
2. offer two alternative formulations of ranking adjectives using monolingual corpus, each one achieves parity with the state of the art.
3. successfully incorporate bilingual data with monolingual data, and outperforms the state of the art by a non-trivial amount.

Additionally, all code and data needed to reproduce the successful experiments in this thesis are distributed and can be found at: <https://github.com/lingxiao/adj-relation>.

This thesis is organized as follows: in the rest of this chapter we will motivate the problem of adjective ranking and give an overview of prior work. In chapter 2 we will discuss the two sources of data we consider in this thesis, as well as how the gold standards are procured. We will also discuss how the rankings are typically evaluated. Chapter 3 gives thorough review of the state of the art method, and implementation details that arose while we reproduced the method. In chapter 4 we offer a very different formulation of the problem at a high level. Chapter 5 will present a simple baseline in the spirit of this reformulation that performs surprisingly well. Chapter 6 refines the baseline using a variety of models, and finally in chapter 7 we combine the models and outperforms the state of the art by a significant amount. We conclude chapter 7 by recapping the underlying assumption of or

approach and offer new avenues for exploration. Chapters 8 and 9 follows logically from chapter 3, they are optional reading for those who wish to consider alternative points of view on the problem. Chapters 8 will gives a formulation that is on parity with the previous state of the art, although this model uses monolingual data only. Chapter 9 gives a detailed account of a failed attempt at ranking adjectives: this chapter is meant to amuse rather than to inform.

Now we will set up the problem of ranking adjectives. Linguistic scale is a set of words of the same grammatical category that can be ordered by their expressive strength or degree of informativeness (Sheinman and Tokunaga, 2009). Ranking adjectives over such a scale is a particularly important task in sentiment analysis, recognizing textual entailment, question answering, summarization, and automatic text understanding and generation. For instance, understanding the word “great” is a stronger indicator of quality than the word “good” could help refine the decision boundary between four star reviews versus five star one. However, current lexical resources such as WordNet do not provided such crucial information about the intensity order of adjectives.

Past work approached this problem in two ways: distributional and linguistic-pattern based. Kim and de Marneffe (2013) showed that word vectors learned by a recurrent neural network language model can determine scalar relationships among adjectives. Specifically, given a line connecting a pair of antonyms, they posited that intermediate adjective word vectors extracted along this line should correspond to some intensity scale determined by the antonyms. The quality of the extracted relationship is evaluated using indirect yes/no question answer pairs, and they achieved 72.8% pairwise accuracy over 125 pairs.

While distributional methods infer pairwise relationship between adjectives based on how they occur in the corpus separately, linguistic-pattern based approaches decides this relationship using their joint co-occurrence around pre-determined

patterns (Sheinman and Tokunaga, 2009; Schulam and Fellbaum, 2010; Sheinman et al., 2012) . For example, the phrase “good but not great” suggests good is less intense than great. These patterns are hand-curated for their precision and unsurprisingly enjoy high accuracy. However, they suffer from low recall because the amount of data needed to relate a pair of adjectives is exponential in length of the pattern, while such patterns are no less than four to five words long.

de Melo and Bansal (2013) addressed this data sparsity problem by exploiting the transitive property of partial orderings to determine unobserved pairwise relationships. They observed that in order to deduce an ordering over good, great, and excellent, it suffices to observe good is less than great, and great is less than excellent. Then by transitive property of the ordering we conclude good is also less than excellent. This fixed relationship among adjectives is enforced by a mixed integer linear program (MILP). Banal and de Melo tested their approach on 91 adjective clusters, where the average number of adjectives in each cluster is just over three, and each cluster is ranked by a set of annotators. They reported 69.6% pair-wise accuracy and 0.57 average Kendall’s tau. Now we will give a detailed explanation of how the monolingual and bilingual corpus is constructed, how the gold standards are curated, and how performance is evaluated.

## 2 Data and Evaluation

This section contains a detailed description of how the data is procured and pre-processed, as well as how the training and test sets are created.

### 2.1 Extracting Intensity Patterns from Monolingual Data

Previously, de Melo and Bansal (2013) showed that linguistic patterns connecting two adjectives reveal semantic intensities of these adjectives. Their work traces back to Sheinman Sheinman et al. (2012), who extracted the patterns by first compiling pairs of seed words where the relative intensity between each pair is known. Then they collected patterns of the form “a \* b” for each pair from an online search engine, where \* is a wildcard denoting one or more words, and word “a” is fixed to be weaker than word “b”. Sheinman then took the intersection of all wildcard phrases appearing between all pairs of words, thereby revealing a set of “weak-strong” patterns  $P_{ws}$  where words appearing in front of the pattern is always weaker than the word appearing behind. Table 1 shows the weak-strong patterns extracted by Sheinman. Bansal used a similar approach but used the Google N-gram corpus (?) as the source of patterns. Additionally, they also considered “strong-weak” patterns  $P_{sw}$  where words appearing in front of the pattern are stronger than those appearing behind. See table 1 for the set of strong-weak and weak-strong patterns mined by Bansal. Finally, during the course of the project, we found additional strong-weak patterns in the N-gram corpus that increased the accuracy of our results, they are found in table 2.

## 2.2 Extracting Adjectives associated with Intensity Patterns from Monolingual Data

We used the Google N-gram Web 1T 5-gram Version 1 publicly distributed by the Linguistic Data Consortium to replicate Bansal’s results. Because we aggressively downsized the N-gram corpus, a detail account of our process is given here. The entire N-gram corpus was first normalized by case folding and white-space stripping. Then for each linguistic pattern in tables 1 and 2, we grepped the corpus for key words appearing in each pattern. Both the grep commands and their corresponding grepped N-grams are located in the raw-data directory of project folder. The grepped N-gram corpus is several times smaller than the original corpus, thus dramatically increasing the number of experiments we can perform.

Next, we crawled the grepped corpus for the patterns found in tables 1 and 2. Specifically, for each pattern of form  $*P*$  and pairs of words  $a_1$  and  $a_2$ , we collect statistics for  $a_1Pa_2$  and  $a_2Pa_1$ . In a departure from Bansal’s method, we also collected statistics for  $*Pa_1$ ,  $*Pa_2$ ,  $a_1P*$ , and  $a_2P*$ , where  $*$  is allowed to be any string. Finally, we also count the occurrences of each pattern  $*P*$ .

## 2.3 Extracting Bilingual-induced Data from PPDB

While de Melo and Bansal (2013) and Sheinman et al. (2012) only considered patterns that relate pairs of adjectives to each other, we also considered adverbs and adverb phrases that occur in front of adjectives, thereby modifying their intensity. We hypothesized that the adverbs can be roughly separated into three classes: intensifying, deintensifying, and neutral. For example, we suspect the adverb “extremely” might intensify adjectives such as “good” in the phrase “extremely good”, while “slightly” would deintensify adjectives it modifies. In general however, neither the class in which adverbs belong to nor the degree in which they modify

adjectives are clear, thus both need to be learned from corpora.

The paraphrase database (PPDB) (?) makes this learning problem possible. This database maps English utterances to other English utterances of similar meaning, so that if  $(x_1, x_2)$  appears in the database, then we conclude  $x_1$  and  $x_2$  are paraphrases. Section 3.3 outlines a method that uses (adverb-adjective, adjective) and (adjective, adverb-adjective) pairs to simultaneously assign scalar values to both the adjectives and adverbs for the task of adjective ranking. We test our assignment on pairs of adjectives ranked by Amazon mechanical turks that also appear in the PPDB corpus. In order to reduce noise in the labels, we removed pairs of adjectives where there is no simple majority consensus among the turks.

## 2.4 Metrics

In this section we describe two ways in which ranking is evaluated: pairwise accuracy and Kendall's- $\tau$ . First we consider pairwise accuracy. If every adjective in each cluster is assigned a numerical ranking  $r(a_i)$ , then the label of each pair is defined as as:

$$L(a_i, a_j) = \begin{cases} > & \text{if } r(a_i) > r(a_j) \\ < & \text{if } r(a_i) < r(a_j) \\ = & \text{if } r(a_i) = r(a_j). \end{cases}$$

Given gold-standard labels  $L_G$  and predicted labels  $L_P$ , the pairwise accuracy of each cluster of adjectives is the fraction of pairs that are correctly classified:

$$PW = \frac{\sum_{i < j} \mathbb{1}(L_G(a_i, a_j) = L_P(a_i, a_j))}{\sum_{i < j} \mathbb{1}}$$

Now we describe Kendall's- $\tau$ , which captures rank correlation between the gold-standard and the predicted set. Kendall's- $\tau$  measures the total number of

pairwise inversions ?):

$$\tau = \frac{P - Q}{\sqrt{(P + Q + X)(P + Q + Y)}}. \quad (5)$$

P measures the number of concordant pairs and Q is the number of discordant pairs, X is the number of pairs tied in the gold ranking, and Y is number of ties in the predicted ranking. The pair  $(a_i, a_j)$  are:

- concordant if  $r_G(a_i) < r_G(a_j)$  and  $r_L(a_i) < r_L(a_j)$  or  $r_G(a_i) > r_G(a_j)$  and  $r_L(a_i) > r_L(a_j)$
- discordant if  $r_G(a_i) < r_G(a_j)$  and  $r_L(a_i) > r_L(a_j)$  or  $r_G(a_i) > r_G(a_j)$  and  $r_L(a_i) < r_L(a_j)$
- tied if  $r_G(a_i) = r_G(a_j)$  and  $r_L(a_i) = r_L(a_j)$  or  $r_G(a_i) = r_G(a_j)$  and  $r_L(a_i) = r_L(a_j)$

Since many of our ranking methods do not allow ties, we also consider a variant where the ties are not counted:

$$\tau' = \frac{P - Q}{n \cdot (n - 1)/2}, \quad (6)$$

here  $n$  is the number of adjectives in a cluster, and  $\frac{n \cdot (n-1)}{2}$  is the total number of unique pairs. In this case the predicted label is discordant w.r.t. gold if the label is flipped, or if the gold-standard pair is a tie. The overall efficacy of each ranking method is captured by finding the average kendall's tau score. Additionally, Bansal observed that sometimes the ordering of adjectives was clear but the annotators would disagreed about which end of the scale was the stronger one, thus absolute kendall's tau is also reported.



During the course of this project we observed that it is possible to outperform certain gold standards under (5). This behavior is highly unexpected and it behooves the reader to consider this concrete example. Suppose our gold standard is:  $G = [[a, b], [c]]$ , read as:  $a$  is tied with  $b$ , and they both dominate  $c$ . Then relative to itself,  $G$  has three unique pairs:

$$pairs = [(a, b), (a, c), (b, c)],$$

two concordant pairs:  $P = [(a, c), (b, c)]$ , no discordant pairs, and one tied pair so that both  $X$  and  $Y$  in (5) are one. Thus Kendall's tau of  $G$  with respect to itself is:

$$\begin{aligned}\tau &= \frac{2 - 0}{\sqrt{(2 + 0 + 1)(2 + 0 + 1)}} \\ &= \frac{2}{\sqrt{9}} = \frac{2}{3}.\end{aligned}$$

Observe in the case of ties, the maximum Kendall's tau is less than 1. Next consider the ranking  $A = [[a], [b], [c]]$ , read as  $a$  dominates  $b$  and  $c$ , while  $b$  dominates  $c$ . Once again we have two concordant pairs  $[(a, c), (b, c)]$  but no discordant pairs by definition,  $X = 1$  and now  $Y = 0$  because  $A$  does not have ties. Thus  $A$  with respect to  $G$  is:

$$\begin{aligned}\tau &= \frac{2 - 0}{\sqrt{(2 + 0 + 1)(2 + 0)}} \\ &= \frac{2}{\sqrt{6}} > \frac{2}{3}.\end{aligned}$$

Therefore an algorithms that ranks the adjectives in correct order without ties can actually outperform the gold standard against itself if the gold ranking does have ties. In the interest of fair comparison, we also report how well gold performs against itself in table 5.

<b>Strong-Weak Patterns</b>	<b>Weak-Strong Patterns</b>
not * (,) just *	* (,) but not *
not * (,) but just *	* (,) if not *
not * (,) still *	* (,) although not *
not * (,) but still *	* (,) though not *
not * (,) although still *	* (,) (and,or) even *
not * (,) though still *	* (,) (and,or) almost *
* (,) or very *	not only * but *
not * (,) just *	not just * but *

Table 1: Bansal and de Melo’s linguistic patterns. Note the syntax (and,or) means either one of “and” or “or” are allowed to appear, or not appear at all. Similarly, (,) denotes a comma is allowed to appear. Additionally, articles such as “a”, “an”, and “the” are may also appear before the wildcards. Wildcards matches any string.

<b>Strong-Weak Patterns</b>	<b>Weak-Strong Patterns</b>
* (,) unbelievably *	very * (and,or) totally *
* not even *	* (,) yet still *
	* (,) (and,or) fully *
	* (,) (and,or) outright *

Table 2: The weak-strong patterns were found by Sheinman. We mined for the strong-weak patterns from google N-gram corpus.

### 3 MILP Method using Monolingual Corpus

This this chapter we reproduce the state of the art in adjective ranking by Bansal and de Melo de Melo and Bansal (2013).

#### 3.1 Problem Formulation

Bansal and de Melo use pairwise co-occurrence of adjective pairs around the paraphrases found in table 1 to infer the relative strength between the adjectives. Because this data is missing for most pairs, pairwise rankings are computed when possible, otherwise Bansal and de Melo used the transitive property of partial rankings to infer the missing relationships. Pairwise ranking is computed as:

$$score(a_1, a_2) = \frac{(W_1 - S_1) - (W_2 - S_2)}{cnt(a_1) \cdot cnt(a_2)},$$

where:

$$W_1 = \frac{1}{P_1} \sum_{p \in P_{ws}} cnt(p(a_1, a_2))$$

$$W_2 = \frac{1}{P_1} \sum_{p \in P_{ws}} cnt(p(a_2, a_1))$$

$$S_1 = \frac{1}{P_2} \sum_{p \in P_{sw}} cnt(p(a_1, a_2))$$

$$S_2 = \frac{1}{P_2} \sum_{p \in P_{sw}} cnt(p(a_2, a_1)),$$

with:

$$P_1 = \sum_{p \in P_{ws}} cnt(p(*, *))$$

$$P_2 = \sum_{p \in P_{sw}} cnt(p(*, *)).$$

Observe  $W_1$  measures the likelihood of encountering the phrase  $a_1 p a_2$  conditioned on the fact that the corpus is composed entirely of phrases of form  $*p*$ ; a similar interpretation holds for  $W_2$ ,  $S_1$ , and  $S_2$ . Furthermore,  $(W_1 - S_1) - (W_2 - S_2)$  is positive when  $a_1$  occurs more often on the weaker side of the intensity scale relative to  $a_2$ , hence  $score(a_1, a_2)$  is an *cardinal* measure how much weaker  $a_1$  is relative to  $a_2$ . The denominator  $cnt(a_1) \cdot cnt(a_2)$  penalizes high absolute value of the numerator due to higher frequency of certain words, thus normalizing the score over all pair of adjectives, and therefore global comparison is well defined over some cardinal scale. Finally, observe that  $score(a_1, a_2) = -score(a_2, a_1)$ .

Given pairwise scores over a cluster of adjectives where a global ranking is known to exist, Bansal then aim to recover the ranking using mixed integer linear programming. Assuming we are given  $N$  input words  $A = \{a_1, \dots, a_N\}$ , the MILP formulation places them on a scale  $[0, 1]$  by assigning each  $a_i$  a value  $x_i \in [0, 1]$ . The objective function is formulated so that if  $score(a_i, a_j)$  is greater than zero, then we know  $a_i$  is weaker than  $a_j$  and the optimal solution should have  $x_i < x_j$ . The entire formulation is reproduced below:

**Maximize**

$$\sum_{i,j} (w_{ij} - s_{ij}) \cdot score(a_i, a_j)$$

**s.t**

$$\begin{aligned}
d_{ij} &= x_j - x_i & \forall i, j \in \{1, \dots, N\} \\
d_{ij} - w_{ij}C &\leq 0 & \forall i, j \in \{1, \dots, N\} \\
d_{ij} + (1 - w_{ij})C &> 0 & \forall i, j \in \{1, \dots, N\} \\
d_{ij} + s_{ij}C &\geq 0 & \forall i, j \in \{1, \dots, N\} \\
d_{ij} - (1 - s_{ij})C &< 0 & \forall i, j \in \{1, \dots, N\} \\
x_i &\in [0, 1] & \forall i, j \in \{1, \dots, N\} \\
w_{ij} &\in \{0, 1\} & \forall i, j \in \{1, \dots, N\} \\
s_{ij} &\in \{0, 1\} & \forall i, j \in \{1, \dots, N\}.
\end{aligned}$$

Note  $d_{ij}$  captures the difference between  $x_i$  and  $x_j$ ,  $C$  is a very large constant greater than  $\sum_{i,j} |\text{score}(a_i, a_j)|$ . If the variable  $w_{ij} = 1$ , then we conclude  $a_i < a_j$ , and vice versa for  $s_{ij}$ . The objective function encourages  $w_{ij} = 1$  for  $\text{score}(a_i, a_j) > 0$  and  $w_{ij} = 0$  otherwise. Furthermore, note either  $s_{ij}$  or  $w_{ij}$  can be one, thus the optimal solution does not have ties. Bansal then extended the objective to incorporate synonymy information over the  $N$  adjectives, defined by  $E \subseteq \{1, \dots, N\} \times \{1, \dots, N\}$ . The objective is now to maximize:

$$\sum_{(i,j) \notin E} (w_{ij} - s_{ij}) \cdot \text{score}(a_i, a_j) - \sum_{(i,j) \in E} (w_{ij} + s_{ij}) \cdot C,$$

while the constraints remain unchanged. The additional set of terms encourages both  $s_{ij}$  and  $w_{ij}$  to be zero if both  $a_i$  and  $a_j$  are in  $E$ , thus the optimal solution may contain synonyms.

Method	Pairwise Accuracy	Avg. $\tau$	Avg. $ \tau $	Avg. $\tau'$	Avg. $ \tau' $
Inter-Annotator Agreement	78.0%	0.67	0.76	N/A	N/A
Gold Standard	100.0%	0.90	0.90	0.90	0.90
MILP reported	69.6%	0.57	0.65	N/A	N/A
MILP with synonymy reported	78.2%	0.57	0.66	N/A	N/A
MILP reproduced	68.0%	0.55	0.64	0.41	0.54
MILP with new patterns reproduced	69.6%	0.57	0.66	0.44	0.56

Table 3: Main results using Bansal’s patterns. Note  $\tau$  refers to kendall’s  $\tau$  with ties, while  $\tau'$  refers to the variant where ties are not considered.

## 3.2 Results

In this section we report both the results of Bansal’s method as stated in their paper, and the results we reproduced. Table 3 reports all results from the different approaches. Note that our reproduction of Bansal’s MILP is accurate with acceptable errors. However, we were not able to reproduce MILP with synonymy accurately because Bansal relied on synonyms marked by annotators, these annotations were not released in the public code base accompanying the paper. Instead we attempted to replicate the experiment using synonyms used by word-net, and observed a marked decrease in accuracy across all measures. However since the synonym set are presumably different, these results cannot be meaningful and thus are not reported.

## 4 Reformulation

In this brief chapter, we give an overview of how we formulated the problem of recovering the “most likely” total ordering of adjectives. However most of the interesting details are missing, and this chapter serves as a motivation for the rest of the thesis. Suppose we have a cluster of  $n$  items,  $\mathcal{C} = \{s_1, \dots, s_n\}$ , then we can form the set of all permutations of  $\mathcal{C}$  by:

$$\Omega = \{\omega \in \Pi(\mathcal{C}) : \Pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}\}.$$

Now we need to place a distribution over  $\Omega$  so that:

$$\begin{aligned} \Pr[\omega] &= \Pr[s_1 < \dots < s_n] \\ &= \Pr[s_1 < s_2 \text{ and } \dots \text{ and } s_1 < s_n \text{ and } s_2 < s_3 \dots s_2 < s_n \text{ and } \dots s_{n-1} < s_n] \\ &= \prod_{i \in \{1, \dots, n\}, i < j} \Pr[s_i < s_j], \end{aligned}$$

where the last statement follows from the independence assumption among pairwise comparisons. If we assume this  $\Pr$  exists and can be estimated, then we can pick the most likely ordering  $\omega^*$  with:

$$\arg \max_{\omega \in \Omega} \Pr[\omega].$$

Some remarks are in order.

*Remark 4.1.* The independence assumption is a fair characterization of how the data is generated: we surmise when people decide if one word is stronger than another in everyday speech, they are not imagining how they assigned strength to other words in the cluster the last time they spoke of them. Note however the test set may be generated in a very different manner. If the comparisons are done pairwise by turkers, then this approximates the setting of every day speech where

the decisions are independent. However if the annotator is given all  $n$  words at once, and asked to rank them, then there is strong sequential dependence between the pairwise comparisons. That is to say if the turker decides that “good” is less intense than ”great”, then this will inform how he/she places the word “better”. Furthermore, the sequence in which the turker makes decisions may lead to different orderings, since the conditional probability that  $\Pr[s_i < s_j | s_i < s_k]$  may not be the same as  $\Pr[s_i < s_j | s_i < s_l]$  for  $k \neq l$ . However, we do not observe this sequence of decisions and therefore cannot decide if our estimate is close to the true distribution. Thus the independence assumption is, in some sense, the most conservative approach.

*Remark 4.2.* The complexity of a naive computation of the  $\arg \max$  operation is  $n^2n!$  in the size of the cluster  $\mathcal{C}$ . In theory this operation is prohibitively expensive and there is literature tackling just this problem. In practice our clusters are no more than 5 items large, so it is very manageable even on my two year old laptop.

*Remark 4.3.* Another possible formulation for a distribution over  $\Omega$  is to assume that there exists a distribution  $\mathcal{D}$  over  $\Omega$ , and nature (or man) selects  $\omega$  according to  $\mathcal{D}$ , and then reveal some pair from  $\omega$  according to another distribution  $\mathcal{D}'$  over the set of all true comparisons implied by  $\omega$ :  $\{s_i < s_j : i, j \in \{1, \dots, n\}, i < j\}$ . Again, we cannot measure either of these distributions from the data we have, so this formulation is curious but not helpful.

*Remark 4.4.* Some readers may find the index notation confusing, so let us be very clear on what we are enumerating in  $\Pi$  and  $\prod$ . Given vertices  $s, t$ , and  $r$ , there are  $3!$  elements in  $\Omega$  enumerated by  $\Pi$ :



$$s < t < r$$

$$s < r < t$$

$$t < s < r$$

$$t < r < s$$

$$r < s < t$$

$$r < t < s.$$

These are the set of all possible orderings of the three vertices. Next for each  $\omega$  in  $\Omega$ , the indices of  $\prod$  enumerates all possible consistent orderings implied by this  $\omega$ . For example, if  $\omega := s < t < r$ , then we have  $O(n^2)$  comparisons:

$$s < t$$

$$s < r$$

$$t < r.$$

Note a contradictory possibility such as  $s > t, s > r, t < r$  is not enumerated, so it could never be picked. Unlike Mohit, we resolve contradiction by not considering it. Finally, since we assume the decisions above are independent, the probability of this  $\omega$  is:

$$\Pr[s < t < r] = \mathbf{Pr}[s < t] \cdot \mathbf{Pr}[s < r] \cdot \mathbf{Pr}[t < r].$$

In the next chapter, we offer some ways of determining what the probability should be over, and how to estimate  $\mathbf{Pr}[s_i < s_j]$  on different data sets. Unless it is explicitly noted, we recover the total ordering from pairwise probabilities using the *argmax* expression above.

## 5 Naive Baselines

### 5.1 Introduction

In this chapter we present a two simple baselines to estimate  $\Pr[s < t]$ .

### 5.2 “Random” Baseline

If we assume the sign  $<$  is a Bernoulli variable so that either  $s < t$  or  $t < s$ , then we can assign a uniform probability of  $\frac{1}{2}$  to each outcome. Thus all  $\omega \in \Omega$  will have equal probability. This presents an interesting problem to pick the maximum value, since in practice  $\Omega$  is represented by a list, Python will pick the maximum over a list of equal valued items by selecting the first element in the list, therefore the ordering of this list is of profound importance. We could default to randomization, but this introduces some uncertainty in our baseline; picking deterministically is a must to prevent pollution from a bad seed. If we sort  $\Omega$  *lexicographically*, then there may not be any bias given some arbitrary cluster of words (speaking loosely). But if there is any relationship between intensity of words and their surface form, then this sorting will introduce a huge bias. This is exactly the case in any cluster with base-comparative-superlative words, sorting alone guarantees 85% pairwise accuracy on the set with only base-comparative-superlative clusters. The solution is to sort and then *reverse* the list, which will create an equally strong bias against the base-comparative-superlative pairs. But this will ensure any gains we have in accuracy come from how well we incorporate information, not sorting. In conclusion, the baseline is still contaminated by a bias, but it is the “worst possible contamination.” We hope the reader finds this argument convincing. Finally, we use the same sort and max function in all future methods so the bias will be consistent.

We test this baseline on three data sets: the original data set annotated by Mohit,

	N-gram		
Test set	Pairwise	Avg. $\tau$	Avg. $ \tau $
Mohit	43.0%	-0.04	0.42
Turk	42.0%	-0.07	0.62
BCS	16.0%	-0.69	0.99

Table 4: Random baseline. Note how poorly “randomness” performed on the base-comparative-superlative baseline simply because the lists are sorted lexicographically, while the other sets perform close to random as expected. A high absolute  $\tau$  is concerning because this suggests that enough of the cluster in the test sets are of length 2 so as to make absolute  $\tau$  look deceptively high.

the set we constructed using mechanical turks, and finally a set of base comparative superlative adjectives found in the PPDB data base. All results are presented in table 1. Readers who are satisfied with how we constructed the baseline may skip the remarks, those who are unsatisfied may skim and critique.

*Remark 5.1.* Readers who are committed to a truly random baseline have to contend with the fact that certain methods only have a slight edge relative to others, it is important the measures reflects this edge. Since many clusters are only two or three items long, different randomized outcomes could result in a  $\tau$  of 1.0, 0.333, or  $-1.0$ . Therefore randomization will overwhelm any gains in method quality over different runs of the method. In this sense, a deterministic lexicographical ordering presents the best solution to make different methods comparable during development.

### 5.3 Pointwise Estimation

The next simplest baseline to estimate  $\Pr[s < t]$  in the spirit of pointwise estimation. Similar to the baseline above, we have two possible events:  $\Omega = \{s < t, s > t\}$ , and we observe a sequence of comparisons between  $s$  and  $t$ :  $\mathcal{S} = \{s < t, s < t, \dots, s > t \dots\}$ , we can ask what is the probability that the next element we will observe is  $s < t$ . This is a Bernoulli distribution with parameter  $p$  and it is well known that the most likely  $p$  is simply:

$$\Pr_{\Omega}[s < t] = \frac{|\{s < t \in \mathcal{S}\}|}{|\mathcal{S}|}.$$

Because this is a baseline, if  $\mathcal{S}$  is empty then we default to  $\Pr[s < t] = \frac{1}{2}$ , that is to say “we don’t know.”

### 5.4 Pointwise Estimation Results

In this section we analyze the results for where the pointwise estimation baseline fails. The goal of the pointwise baseline is to two fold: (1) assess how much information can be gained by considering pairwise comparisons alone, and (2) understand how the three data sets differ. We estimate the probabilities over three data sets, the N-gram data set composed of all occurrences of  $s\mathcal{P}t$  for specified patterns  $\mathcal{P}$ , the PPDB set of paraphrases, and finally a naive combination of the two sets where the data is simply “thrown together.” All results are presented below in table 1. Readers who do not like to “lay in bed with the data” may skip to the conclusion, those who are not too smart to count may read each paragraph in detail.

We examine each annotated set over all three data sets. First let us examine the Mohit’s test sets over the N-gram data. Over all, all the clusters have some corpus evidence for at least one pair of comparisons. Note similar to Mohit’s algorithm, we correctly place “first, . . . , eight” in the correct order even though we

only observe data for some of the comparisons. Furthermore we note that sometimes multiple orderings in  $\Omega$  will have equal probabilities, this is not the case here, we picked the unique  $\omega^*$ . In the case of “close, near, intimate” we only observe data for one out of three possible comparisons, and correctly placed the near to be less than intimate. Now we have two possibilities: “near < close < intimate” and “close < near < intimate”. Since the list of options are sorted lexicographically and “near” comes after “close” lexicographically, we picked “near < close < intimate”. In “real, solemn, serious, grave”, there is strong n-gram corpus evidence that serious is less intense than solemn, thus the order is flipped. Finally, 9 cluster have a negative  $\tau$ , all of which have strong corpus evidence supporting the ordering under the measure we defined.

Now we examine Mohit’s test over the PPDB data, note it performed only slightly better than the random baseline, suggesting the PPDB corpus has marginal value add. Notably, the cluster “first, . . . , eight” performed poorly because there is no data in the PPDB graph for these words. The cluster “close, near, intimate” now has a negative  $\tau$  because there are four paraphrases in the PPDB corpus suggesting “close < near”:

1. quite close  $\rightarrow$  near
2. real close  $\rightarrow$  near
3. really close  $\rightarrow$  near
4. so close  $\rightarrow$  near
5. too close  $\rightarrow$  near
6. very close  $\rightarrow$  near,

and no evidence suggesting “near < close.” In the cluster “real, solemn, serious, grave” we output an ordering that is of the same pairwise accuracy and  $\tau$  value as the N-gram case, but making different mistakes: flipping the order of “real” and “solemn” in this case but not because of evidence because there is none, and “solemn” is after “real” lexicographically. Overall, 47 out of the 79 clusters in the Turk set did not have any data for any pairs of comparisons, this is close to 60% of the data set. In all the clusters with a negative  $\tau$ , all but one case was due to complete lack of data. A similar story is true in BCS data set, where our method performs better than the “random” baseline because every time there is evidence for the words, the < sign is flipped to the correct ordering. Finally, we simply note that Mohit’s data set did comparably well on the PPDB + N-gram set, this is not surprising since the many of the Mohit’s words do not make an appearance in the PPDB data set.

Now we consider the Turk set over all data sets. On the N-gram data, the Turk set performed no better than random. Over 70% of the clusters in the Turk set do not have any observations in the N-gram set. Thirty five clusters in this set have negative  $\tau$ ’s, and 29 of which are due to lack of data. In other cases, a negative  $\tau$  is either due to corpus evidence contradicting the gold set (two cases) or due to ties in the gold set, which our model does not account for. Next we consider the Turk set over the PPDB data, note that 68% pairwise accuracy is a respectable showing if we recall that Mohit achieved 69.2% accuracy with the MILP formulation on his dataset. Over all 38 of the clusters have data for every possible comparison between words ( $O(n^2)$  comparisons), however 23 of these clusters have only two words. Five clusters have no observation for any of the links, curiously three out of these five clusters also only contains two words. For the curious the clusters are:

1. uncomfortable, embarrassed

2. shitty, awful
3. sturdy, intact
4. vast, plenty, abundant
5. tough, formidable, daunting.

Suffice it to say the last cluster is a fair description of writing this thesis, we hope the second cluster does not describe the experience of *reading* this thesis. Finally 14 clusters have negative  $\tau$ 's, two of these clusters have no data (the second and third cluster from the list above, for those who are wondering). The rest have corpus evidence that contradicts annotators. For example, the annotators ranked “hardworking < tough < tenacious”, while the algorithm ranked “tough < tenacious < hardworking” because in the PPDB corpus we observe:

1. very tough  $\rightarrow$  tenacious
2. very tough  $\rightarrow$  hardworking
3. pretty tough  $\rightarrow$  hardworking
4. really tough  $\rightarrow$  hardworking
5. so tough  $\rightarrow$  hardworking.

We leave it to the reader’s imagination for the interpretation of < in this setting, and which quality should rank higher under <. Finally we examine the Turk set over the PPDB + N-gram data set, where our estimation performed as well as we did on Mohit’s set. This is encouraging because this suggests that although the two data sets performed drastically differently over the PPDB and N-gram data sets alone, they performed comparably well on the combined data sets. Which

confirms our suspicion that the lack of data is the cause of the discrepancy, not how the sets are curated.

Finally we examine the base comparative superlative set. When test on the N-gram data set, 189 out of 285 clusters have no data according to the N-gram corpus, 64 clusters have data for every pair of comparison. One hundred and eighty two clusters have negative  $\tau$ 's, 24 of which have data for one comparison out of the  $O(n^2)$  possible comparison between words, in all cases corpus evidence placed the words in the correct order, but lexicographical ordering ensured that the overall order of words still has a negative  $\tau$ . For example, in the cluster “short, shorter, shortest” we observe  $\Pr[short < shorter] = 0.99$  but have no observations for the other pairs, so the overall ranking is  $shortest < short < shorter$ , again due to lexicographical ordering.

Now we examine the BCS set on the PPDB data set. Pairwise accuracy is appreciably higher on this data set because only 101 clusters have no data, while 159 have data between all  $O(n^2)$  possible comparisons. Ninety five clusters have negative  $\tau$ 's, however six of these clusters have data supporting the ranking. The ranking output by our methods are:

1. more < many
2. more < some
3. more < much
4. better < well
5. latest < later
6. poorest < poorer.



The fact that “more” makes such a strong showing is curious, in the first cases, there exists just one edge from “more” to “many”: “any more  $\rightarrow$  many.” In the second case there is just one edge again: “a few more  $\rightarrow$  some”. In both cases there are no edges going the other way. In the third cases, there are six edges suggesting more is less intense than much:

1. a lot more  $\rightarrow$  much
2. considerably more  $\rightarrow$  much
3. little more  $\rightarrow$  much
4. lot more  $\rightarrow$  much
5. significantly more  $\rightarrow$  much
6. substantially more  $\rightarrow$  much,

and just one edge suggesting more is more intense than much: very much  $\rightarrow$  more. A sharp eyed reader might immediately ask how connected are the two vertices relative to each other, to satiate your curiosity we report the values here: more has 147 total neighbors, 42 in-neighbors and 115 out-neighbors. Much on the other hand only has 78 neighbors, 23 in-neighbors and 66 out-neighbors. In other words, one vertex may appear to dominate another if we consider the number of edges alone, but may not dominate if we take their overall connectivity in the graph into account. We will use this particular observation to improve our results later.

Finally we close the chapter by examining the BCS labeled set on the PPDB + N-gram data. Incorporating N-gram data has the immediate consequence that one more cluster now has data for all possible comparisons among adjectives, progress

comes in the most incremental steps. Ninety nine clusters still have no data whatsoever, and 93 clusters still have negative  $\tau$ 's, 81 of which are because there is no data for any pairs. Only 5 of these negative clusters have data between all pairwise comparisons, they are have already been listed above. Notably, “more” is now correctly classified as less intense than “much” due to overwhelming N-gram evidence: 67386 edges point from “much” to “more”, while only 2834 edges point the other way.

In conclusion, we observe three kinds of mistakes:

1. mistakes from no data between pairs
2. mistakes from data between pairs that contradicts the gold set
3. mistakes from predicting some ordering, when in fact the words are tied.

We write off mistake number two as bad luck from a poor sample. We also write off mistake three since we will not be developing models that predict synonyms, which is a different task all together. In the next few chapters, we will focus on addressing mistake one: missing data. A very important observation is that in order to increase the overall accuracy, it suffices to have a model that gives us a slight edge over random baseline most of the time in the case for when there is no data. That is we require a model so that:

- if  $s > t$  then we need  $\Pr[s > t] \geq \frac{1}{2} + \epsilon$
- if  $s < t$ , then then we have  $\Pr[s > t] \leq \frac{1}{2} - \epsilon$ ,

with high probability.

*Remark 5.2.* What exactly does high probability entail in our setting? There appears no right answer here,  $\frac{1}{n^2}$  in the number of pairwise comparisons may be ideal,

	N-gram			PPDB			PPDB + N-gram		
Test set	Pairwise	Avg. $\tau$	Avg. $ \tau $	Pairwise	Avg. $\tau$	Avg. $ \tau $	Pairwise	Avg. $\tau$	Avg. $ \tau $
Mohit	<b>72.0%</b>	<b>0.56</b>	<b>0.65</b>	46.2%	0.02	0.46	71.3%	0.53	0.66
Turk	47.0%	0.04	0.62	68.5%	0.49	0.70	<b>71.0%</b>	<b>0.55</b>	<b>0.72</b>
BCS	38.0%	-0.24	0.92	65.5%	0.30	0.94	<b>66.0%</b>	<b>0.33</b>	<b>0.95</b>

Table 5: Results across all datasets. Observe how N-gram graph only performed slightly better than the base line on base-comparative-superlative dataset. A similar story holds for the Turk set. However on Mohit’s set we already manage to achieve a higher or comparable accuracy across all measures on the N-gram set than what Mohit did in his TACL paper.

but perhaps is too much to ask. Next we could ask for  $\frac{1}{n}$ , that is in a data set with 300 unknown pairwise comparisons, we need to achieve 96% pairwise accuracy. Again speaking from intuition this is a very tall order. Thus we will aim for  $\frac{1}{\log(n)}$ , so for 300 unknown pairs we need to achieve 84% pairwise accuracy. This will be the hurdle we aim to achieve for the rest of the thesis. So in a sense, we are actually aiming for “reasonable probability” (my definition, not found in literature), not “high probability.”

## 6 Regression

### 6.1 Introduction

In this section we will find models that approximate the ordering of the adjectives when no comparisons between them are observed. Towards this end, we will use logistic regression and elastic net regression to learn what it means for  $s$  to be less than  $t$ . We construct a variety of feature representations using adverbs and/or phrases that co-occur with the adjectives for this task, and divide our annotated data into train, validation, and test sets to assess their efficacy.

### 6.2 Data Set

Recall we have three sources of comparisons, N-gram data alone, PPDB data alone, and the combination of N-gram and PPDB data. For each data set, we extracted the pairs of adjectives where no data exists in corpus. The number of pairs for each annotated set for each data set is displayed in table 1. We will learn a different model for each data set (N-gram, PPDB, PPDB + N-gram) and assess their efficacy using their respective validation and test sets.

### 6.3 Literature Review

Note in our graph we have 610 adverbs and phrases, and since we are interested in representing an adjective using its co-occurrence with adverbs/phrases, the feature representation could be large than the number of examples. Furthermore, this representation will be very sparse because most adjectives do not co-occur with most adverbs/phrases at all. Regularization will be necessary to prevent over-fitting. We consider two regularized models:  $l_1/l_2$ -penalized regression, and elastic net regression. In this section we will give a brief review of the two models.

Elastic net regression is a natural fit for our setting because it gives us the ability to select features and control sparsity. In this section we will give a brief overview of elastic net regression. Suppose our data set is  $(\mathbf{X}, \mathbf{y})$  so that  $\mathbf{X}$  is the  $n \times p$  design matrix, where each input  $x$  is represented by the appropriate feature vector  $\mathbf{x} \in \mathbb{R}^p$ , and let  $\mathbf{y}$  be the  $n$ -dimensional response vector  $\mathbf{y} = (y_1, \dots, y_n)$ . We assume  $\mathbf{y}$  is generated by this process:

$$\mathbf{y} = \mathbf{X}\beta + z,$$

where  $z$  is a zero mean Gaussian noise factor. Then elastic net regression will recover the estimated  $\hat{\beta}$  where:

$$\hat{\beta} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\| + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1.$$

Roughly, the  $l_1$  penalty encourages a sparse solution where only a few variables in  $\mathbf{x}$  participate in predicting  $\mathbf{y}$ , while the  $l_2$  penalty encourages “grouping” so that more than a few variables in  $\mathbf{x}$  participates.

Now we review logistic and regularized logistic regression for binary outcomes. Given a  $n \times p$  design matrix  $\mathbf{X}$ , logistic regression models the vector of probability  $\mathbf{p}$  by:

$$\log \frac{\mathbf{p}}{1 - \mathbf{p}} = \mathbf{X}^T \beta,$$

and we see that  $\mathbf{p}$  is:

$$\mathbf{p} = \frac{\exp(\mathbf{X}^T \beta)}{1 + \exp(\mathbf{X}^T \beta)}.$$

Again given the binary outcome vector  $\mathbf{y} \in \{0, 1\}^n$ , the loss function  $\mathcal{L}$  is:

$$\mathcal{L}(\beta) = \log \mathbf{p} + (1 - \mathbf{y})^T \log(1 - \mathbf{p}),$$

we can find the best  $\beta$  by:

$$\hat{\beta} = \arg \min_{\beta} \mathcal{L}$$

In our setting, we experiment with two penalties on  $\beta$ :  $l_2$ -penalty or ridge logistic regression, and  $l_1$ -penalty or LASSO logistic regression. In ridge logistic regression, we simply add the  $l_2$ -penalty to the objective function:

$$\mathcal{L}_{ridge}(\beta) = \mathcal{L} - \frac{1}{2}\lambda \|\beta\|_2^2.$$

Similarly, for LASSO logistic regression, the objective function is:

$$\mathcal{L}_{LASSO}(\beta) = \mathcal{L} - \frac{1}{2}\lambda \|\beta\|_1^2.$$

Again the ridge penalty encourages grouping of all variables, while LASSO encourages sparsity.

## 6.4 Problem Formulation

Now we will formulate our problem in terms of the two models we just introduced.

In our setting, we will define:

$$y = \begin{cases} 1 & s < t \\ 0 & otherwise. \end{cases}$$

And for each pair of adjectives  $s$  and  $t$ . Additionally, we will need to find a corresponding feature representation so that:

$$\mathbf{x} = g(\phi(s), \phi(t)),$$

for some function  $\phi$  and  $g : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^p$ , where  $m \leq p$ . In a departure from notation of previous chapter,  $s$  now refers to the string representation of the word, while  $\phi(s)$  is the corresponding vector representation. If we have this model  $\hat{\beta}$  and the appropriate  $g$  and  $\phi$ , then we can use this definition.

**Definition 6.1.** Given words  $s$  and  $t$ , and their representation  $\mathbf{x} = g(\phi(s), \phi(t))$ , and let:

$$\hat{y} = \widehat{\beta}^T \mathbf{x},$$

where  $\widehat{\beta}$  is the elastic net model, then we can define:

$$Pr[s < x] = \begin{cases} \frac{1}{2} + \epsilon & \hat{y} < \delta \\ \frac{1}{2} - \epsilon & otherwise, \end{cases}$$

for an appropriate threshold  $\delta$ . Again we discard the value of  $\hat{y}$  and define the probability by fiat. In the case of penalized logistic regression, we use the actual probability output by the model.

## 6.5 Feature Representations

In this section we describe two broad sets of features we use, and their associated function  $g$ . In the first set of features, we will represent the adjective by the frequency of adverbs incident and/or outgoing from this adjective. In the second set of features we will represent the adjective by adjectives that are its neighbors. In an attempt to avoid confusion, we will denote the first set of features  $\phi(s)$ , while the second set will be  $\nu(s)$ .

First we list the  $\phi$ 's.

1. In-neighbor only. So that for each adverb  $v$ :

$$\phi(s)_v^{in} = \begin{cases} n & \text{there are } n \text{ edges pointing to } s \text{ from all neighbors with the adverb } v \\ 0 & otherwise. \end{cases}$$

2. Out-neighbor only. So that for each adverb  $v$ :

$$\phi(s)_v^{out} = \begin{cases} n & \text{there are } n \text{ edges pointing from } s \text{ to all neighbors with the adverb } v \\ 0 & \text{otherwise.} \end{cases}$$

3. Concatenation of in and out neighbor. So that  $\phi(s) = (\phi^{in}, \phi^{out})$ . In this case we also considered  $\phi(s) = (\phi^{in}, -\phi^{out})$ , where  $-\phi^{out}$  is scalar multiplication of  $-1$  with all entries of  $\phi^{out}$ .
4. Element wise addition of in and out neighbor. So that  $\phi(s) = \phi^{in} + \phi^{out}$ .
5. Element wise subtraction of in and out neighbor. So that  $\phi(s) = \phi^{in} - \phi^{out}$ .

Furthermore, for each  $\phi^{in}$  and  $\phi^{out}$ , we can vary the number of adverbs in the vector. We sort the adverbs by frequency of appearance and pick the top 10, 20, 30, 40, 50, 100, 200, 300, 400, 500, and all 605 adverbs/phrases.

Now we consider the  $\nu$ 's. In this case we find all neighbors of  $s$  and represent each neighbor by the frequency of adverbs between the neighbor and our vertex. We list the  $\nu$ 's.

1. In-neighbor only. So that for each neighbor  $t$ :

$$\nu(s)_t^{in} = \begin{cases} n & \text{there are } n \text{ edges pointing to } s \text{ from } t \\ 0 & \text{otherwise.} \end{cases}$$

2. Out-neighbor only. So that for each neighbor  $t$ :

$$\nu(s)_t^{out} = \begin{cases} n & \text{there are } n \text{ edges pointing from } s \text{ to } t \\ 0 & \text{otherwise.} \end{cases}$$

3. Concatenation of in and out neighbor. So that  $\nu(s) = (\nu^{in}, \nu^{out})$ . Again we also experiment with  $\nu(s) = (\nu^{in}, -\nu^{out})$ .



4. Element wise addition of in and out neighbor. So that  $\phi(s) = \phi^{in} + \phi^{out}$ .
5. Element wise subtraction of in and out neighbor. So that  $\phi(s) = \phi^{in} - \phi^{out}$ .
6. Bernoulli representation of in and out neighbor. So that for each neighbor  $t$  we have:

$$\nu(s)_t^{\mathcal{B}} = \begin{cases} \frac{|\{(t,s) \in \mathbf{E}\}|}{|\{(t,s) \in \mathbf{E}\}| + |\{(s,t) \in \mathbf{E}\}|} & \text{there is at least one edge from } t \text{ to } s \\ \frac{1}{10} & \text{otherwise.} \end{cases}$$

In all cases, we sort the neighbors so that the most connected neighbor to  $s$  is at the top, and so on. Finally we take the top  $n$  neighbors of  $s$  as its feature representation, where  $n$  is also an experimental parameter. If an adjective has less than  $n$  neighbors, then we pad the vector with 0. Note how we smooth out the definition of  $\nu$  above so that the padding can be distinguished from a case where there are vertices from  $s$  to  $t$ , but not vice versa.

Finally, we explore a variety of  $g$ 's:

1. element wise addition.  $\mathbf{x} = \phi(s) + \phi(t)$ ,
2. element wise subtraction.  $\mathbf{x} = \phi(s) - \phi(t)$ ,
3. concatenation.  $\mathbf{x} = (\phi(s), \phi(t))$ ,
4. dot product.  $\mathbf{x} = \phi(s) \cdot \phi(t)$ ,

and similarly for  $\nu$ . After  $g$  has been applied, we also experimented with using the raw counts, versus normalizing the entries of  $\phi$  so they are between 0 and 1. Normalization appears to make sense when the N-gram data is combined with the PPDB data, since the N-gram data is often times three orders of magnitude larger than the PPDB data. We normalized the raw  $n \times p$  design matrix with  $n$  samples and  $p$  features in one of two ways.

	Mohit	Turk	BCS
N-gram	550	586	556
PPDB	750	182	294
PPDB + N-gram	408	170	290

Table 6: The base-comparative-superlative pairs form the training set for each data set, the Turk pairs form the validations set, while Mohit’s pairs will be the test set. Note in almost all cases but one, the test set is actually larger than the training set.

1. Normalize by mean and variance. Suppose each feature is a unique multinomial random variable  $\mathbf{x} \in \{1, \dots, N\}$  for some suitable  $N$ , then we can normalize this variable by:

$$\tilde{x} = \frac{x - \mathbf{E}[\mathbf{x}]}{\text{var}(\mathbf{x})}.$$

2. Normalize by maximum value. So we have  $\tilde{x} = \frac{x}{\max(x_1, \dots, x_n)}$

## 6.6 Results

In addition to varying the feature representations, we also varied the emphasis over each of the two penalty terms in the objective function. For the sake of brevity, we report results from the two best performing models in each method only: (1)  $l_1$ -penalized logistic regression model with Bernoulli representation of the top 10 most connected neighbors and (2) elastic-net regression where the feature vector is the Bernoulli representation of the top 50 most connected neighbors and. Note this model is learned using the data set with PPDB and N-gram comparisons and base-comparative-superlative gold sets, validated using the Turk set, and tested using Mohit’s set.

	Elastic Net Regression		$l_1$ -Logistic Regression	
<b>Gold Set</b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>
BCS	77.0%	0.54	86.0%	0.72
Turk	61.0%	0.22	61.0%	0.22
Mohit	71.0%	0.44	72.0%	0.44

Table 7: . Results for the two best models. Note the performance on the validation and test sets are comparable between the two models. It is also curious to note that the models performed better on the test set than the validation set. In fact, no model performed better than 70% on the validation set.

All models were implemented using Python’s Scikit-learn library. The best performing  $l_1$ -penalized logistic regression model has a constant  $C$  of 0.4. While the best performing elastic net model has an  $\alpha$  of 0.9, and  $l_1$  of 0.1.

One natural question we can ask is how well the model performs on the pairs of adjective with no direct comparisons, versus how well the baseline performs on the set of pairs with direct comparisons. Table 3 displays the baseline one pairs with direct comparisons. Overall, we see that the model does 9% worse on the Turk set relative to the baseline using PPDB and N-gram data, and does only 2% worse relative to the baseline on Mohit’s clusters.

	N-gram		PPDB		PPDB + N-gram	
Test set	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$
BCS	100.0%	1.00	97.0%	0.93	97.0%	0.94
Turk	78.0%	0.57	69.5%	0.38	70.0%	0.40
Mohit	84.0%	0.67	48.2%	-0.04	74.3%	0.49

Table 8: Results across all datasets for pairs of gold standards for which direct comparison exists. Note how performance is higher when using N-grams alone across all gold standards. Note how PPDB data alone fails to do better than random on Mohit’s clusters, even though direct direct comparisons exists for the pairs in these clusters. Suggesting the data that exists for Mohit’s pairs are too noisy to give useful information.

## 7 Logistic Regression Regression

### 7.1 Introduction

In this section we present a remedial solution to combine our regression models with the pointwise-estimation baseline model, and a more principled solution. Since we are combining the models, our gold standards will now be the original annotated sets composed of with size between two to ten.

### 7.2 Remedial Solution

First we give brief reminder for our baseline model. We defined two possible events:  $\Omega = \{s < t, s > t\}$ , and after observing a sequence of comparisons between  $s$  and  $t$ :  $\mathcal{S} = \{s < t, s < t, \dots, s > t \dots\}$ , we can ask what is the probability that the next element we will observe is  $s < t$ . This is a Bernoulli distribution with parameter  $p$  and it is well known that the most likely  $p$  is simply:

$$\mathbf{Pr}[s < t] = \frac{|\{s < t \in \mathbf{S}\}|}{|\mathbf{S}|}.$$

In the baseline, if  $\mathbf{S}$  is empty then we defaulted to  $\mathbf{Pr}[s < t] = \frac{1}{2}$ .

Now we present the remedial solution. Recall in the previous chapter we defined this probability value for the  $\hat{y}$  output by elastic net regression:

$$\mathbf{Pr}[s < x] = \begin{cases} \frac{1}{2} + \epsilon & \hat{y} < \delta \\ \frac{1}{2} - \epsilon & \text{otherwise,} \end{cases}$$

while we used the actual probability value  $p$  output by the logistic regression model. In the remedial solution, we use the elastic definition defined above, and in the case of logistic regression, we actually discard the value of  $p$  and define:

$$\mathbf{Pr}[s < x] = \begin{cases} \frac{1}{2} + \epsilon & p > \frac{1}{2} \\ \frac{1}{2} - \epsilon & \text{otherwise.} \end{cases}$$

This captures our intuition that the prediction output by the model is less accurate than that of the actual data. Additionally, we also constructed a version of the Turk and Mohit’s clusters where ties are removed. We reasoned that since our models are designed to predict ordering, while ties can be interpreted as synonyms, clusters generated without ties may give a more “fair” representation of how well the models perform. Results are displayed below.

### 7.3 Solution with Beta Prior

In this section we provide a more formal variant of the remedial solution. The heart of the problem is that we have some prior belief about the likelihood that one adjective is weaker than another, and an updated belief after observing some data, be it direct comparison or estimation from a model. Since we are modeling each edge a Bernoulli variable with parameter  $\theta$  ranging over  $[0, 1]$ , the prior is then

	Elastic Net Regression		$l_1$ -Logistic Regression	
<b>Gold Set</b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>
BCS	90.0%	0.81	93.0%	0.85
Turk	75.0%	0.62	74.0%	0.61
Turk no-tie	81.0%	0.63	81.0%	0.62
Mohit	74.0%	0.61	74.0%	0.61
Mohit no-tie	76.0%	0.52	76.0%	0.53

Table 9: . Results for the two best models combined with pointwise estimation baseline in the remedial fashion. Note how two models performs comparable across all gold sets. In addition, not the gold clusters with no ties enjoyed a higher pairwise accuracy but suffer a lower  $\tau$  value.

a distribution over the Bernoulli  $\theta$ , this is the Beta distribution. In this next few paragraphs, we give a brief overview of Beta-Binomial model, in particular how it applies to our problem.

It is well known that the prior the binomial and Bernoulli likelihood function is the Beta distribution with paramters  $\beta_1, \beta_2 \in \{1, \dots\}$ , where we have:

$$\begin{aligned}
\Pr[\theta|\beta_1, \beta_2] &= \frac{\theta^{\beta_1-1}(1-\theta)^{\beta_2-1}}{\int_0^1 \mu^{\beta_1-1}(1-\mu)^{\beta_2-1} d\mu} \\
&= \frac{\Gamma(\beta_1 + \beta_2)}{\Gamma(\beta_1)\Gamma(\beta_2)} \theta^{\beta_1-1}(1-\theta)^{\beta_2-1}.
\end{aligned}$$

The exact form of the  $\Gamma$  function is beyond the scope of this introduction but the reader may select any introductory book on statistics for a refresher.

Now after observing  $n$  coin tosses with  $h$  heads and  $t$  tails for  $h + t = n$ , the posterior probability over  $\theta$  given some prior setting of  $\beta_1$  and  $\beta_2$  is:

$$\begin{aligned}
\Pr[\theta|h, t, \beta_1, \beta_2] &= \frac{\Pr[h|n, \theta] \Pr[\theta|n, \beta_1, \beta_2]}{\Pr[h|n, \beta_1 + \beta_2]} \\
&\propto \theta^{h+\beta_1-1} (1-\theta)^{t+\beta_2-1},
\end{aligned}$$

note the posterior distribution is also a beta distribution. Now we have the distribution  $\mathbf{Pr}[\theta|h, t\beta_1, \beta_2]$ , we can return the pointwise estimation setting and ask what is the likelihood the next toss lands heads, this is exactly the posterior mean:

$$\begin{aligned} E[\theta|h, t\beta_1, \beta_2] &= \int_0^1 \theta \mathbf{Pr}[\theta|h, t\beta_1, \beta_2] d\theta \\ &= \frac{\beta_1 + h}{\beta_1 + \beta_2 + n}. \end{aligned}$$

Not how the last line appeals strongly to intuition and therefore can be easily used: the expected outcome of the next toss given the prior is simply the prior tosses plus the tosses observed from data.

In our setting, we fix the ratio of  $\beta_1$  and  $\beta_2$  so that the prior probability is exactly 1/2, thus reflecting our ignorance. Note this is consistent with our ad-hoc setting in the remedial solution. The exact values of  $\beta_1$  and  $\beta_2$  is a hyperparameter to be tuned, in practice we set  $\beta_1 = \beta_2 = 1$ . The coin tosses observed from data and the model are also hyperparameters. Note the more confident we are with the data, the larger the values of  $h$  and  $t$  should be with respect to  $\beta_1$  and  $\beta_2$ . We experimented with a variety of values, and settled on the following settings for  $h$  and  $t$ :

1. If there is an observation, then we use the raw comparison counts between the adjectives as  $h$  and  $t$
2. If there is no observation so we are using the model, we set  $h$  to be the probability that the model predicts less than, and  $t = 1 - h$ .

In informal terms, we are confident in the quality of direct comparisons, if they can be observed, and not very confident in the prediction of the model. Results for Beta-Binomial model is presented below. All in all, the best model uses the

beta-binomial model to combine direct observations with  $l - 1$ -penalized logistic regression model, the regression model uses top the coin toss probability of 10 most connect neighbors as features. This model achieved 75% pairwise accuracy on Mohit’s data set and the Turk set, and a Kendall’s  $\tau$  value of 0.61 and 0.62. After adjusting for ties, the pairwise accuracy on Mohit’s data set was 76%, which approaches the inter-annotator accuracy of 78%, while the pairwise accuracy on the Turks set was 82% after adjusting for ties.

On the other hand, we observe that although Bansal’s method performs well on the gold set they procured, it performs substantially worse across all other data sets: the  $\tau$ ’s are close to zero and the pairwise accuracy is comparably low. Such surprisingly low results warrants further inspection. The low score actually hides two sources of error: error from absolute no data for any pair of adjectives, and error from wrong rankings when there is some data. In order to disentangle the two sources of error, we constructed a subset of the Turk set and the base-comparative-superlative sets where at least one pair of words from each cluster must have some data. Recall the original Turk set has 79 clusters, the new set only has 23 clusters, that is to say less than 30% of the Turk clusters appear in the N-gram data set. The story is similarly bleak for the base-comparative-superlative gold set: only 64 out of 238 clusters has any data at all, this is just 27%. After adjusting for data sparsity, we see that that Bansal’s method achieves 65% accuracy on the Turk set, which is more in line with how it performed on the Mohit gold set. On the base-comparative-superlative set, Mohit’s method only achieved 56% accuracy, this is not as high as the other gold set, but certainly better than random. We also report results for our regression method using N-gram and PPDB data for the two reduced gold sets, and observe the results are comparable two their performance on the full gold set as expected. All in all, it is clear that the increase in performance is primarily due to better coverage due to PPDB data.



	Elastic Net Regression		$l_1$ -Logistic Regression		MILP	
<b>Gold Set</b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>
BCS	90.0%	0.81	<b>93.0%</b>	<b>0.85</b>	18.0%	0.02
Turk	75.0%	0.62	<b>75.0%</b>	<b>0.62</b>	25.0%	0.13
Mohit	74.0%	0.61	<b>75.0%</b>	<b>0.61</b>	69.6%	0.57
Turk no-tie	81.0%	0.63	<b>82.0%</b>	<b>0.63</b>	19.0%	0.12
Mohit no-tie	76.0%	0.52	<b>76.0%</b>	<b>0.53</b>	68.0%	0.46
BCS has-data	100.0%	1.00	100.0%	1.00	56.0%	0.13
Turk has-data	66.0%	0.47	66.0 %	0.47	65.0%	0.46

Table 10: . The first two columns show results for the two best models combined with pointwise estimation baseline using Beta-Binomial model. The third column displays Mohit’s MILP method using N-gram data only. The results shows that  $l_1$ -logistic regression outperformed elastic net regression on most data sets by a small (possibly insignificant) margin, otherwise they are equivalent. In particular, observe how logistic regression performs just as well on Mohit’s set as it does on the Turk set. Furthermore, both models outperform MILP by a non-trivial amount on all gold sets. Finally, note how well the MILP method performs on Mohit’s gold cluster, versus how poorly it performs on other gold standards. The bottom two rows report results across all methods for the subsets of the gold clusters where there is some N-gram data.

## 7.4 Conclusion and Future Work

We will close this chapter by reiterating our assumption, and offer some possible avenues of exploration. First, we review the fundamental hypothesis underlying all methods we have seen so far: there is a positive relationship between how words are used on average in the context of intensifiers such as adverbs or N-gram patterns, and how an annotator ranks pairs of such words in isolation. The fact that all successful methods we have seen so far achieved close to or above 70% pairwise accuracy confirms this hypothesis. However, we must caution those who wish to improve upon this results, inter-annotator accuracy on Mohit's set is 78%. In fact, many of the rankings are subjective and we do not believe it is meaningful to achieve, for example 90%, on the gold standards procured by Mohit and us.

In fact, consider this cluster:

- interesting
- entertaining
- fascinating
- intriguing
- amusing
- exciting,

the gold set suggests

interesting < intriguing < amusing = entertaining < fascinating = exciting,

but an equally plausible ranking is:

interesting < amusing < entertaining < fascinating < intriguing < exciting.

Penalized elastic net regression output:

interesting < entertaining < fascinating < intriguing < amusing < exciting.

We leave it to the reader to judge which cluster makes the most sense, if any. This example suggests an immediate improvement upon our work: procure better clusters so the results are more meaningful. There are two ways to improve the quality of each cluster: prune the clusters for synonyms, and resist the temptation for large clusters. The first suggestion is self explanatory, since none of the models account for synonyms, clusters that contain words annotators would consider to be synonyms will unfairly depress the score. Eschewing large clusters is important because the larger the cluster, the more likely polysemy will

## 7.5 ILP over N-gram Patterns

In this chapter, we offer an alternate formulation of recovering ranking using monolingual data using integer linear programming (ILP). This chapter stands apart from the previous two chapters, and follows chapter 3 directly.

## 7.6 Problem Formulation

Recall in Mohit's work,  $score(a_i, a_j)$  is zero for most pairs of adjectives due to lack of data, we are motivated to find alternate ways of approximating this value using single sided patterns of form:  $\{a_i p^*, a_j p^*, *p a_i, *p a_j\}$ . Loosely speaking, even if we do not observe any  $a_i p a_j$  for some weak-strong pattern  $p$ , we can still

approximate the likelihood of observing this string using the frequency in which  $a_i$  appears in front of the the pattern  $p$ , and the frequency  $a_j$  appears behind the pattern  $p$ . Once we approximate the likelihood for  $a_j p a_i$  over weak-strong patterns  $p$ , and similarly for all strong-weak patterns, we can infer whether adjective  $a_i$  is weaker than  $a_j$  by determining whether  $a_i$  is more likely to appear on the weaker side of each phrase. This intuition is naturally expressed in the heuristic:

$$score(u) = \frac{cnt(u p_{sw} *) + cnt(* p_{ws} u)}{cnt(u p_{ws} *) + cnt(* p_{sw} u)},$$

where:

$$cnt(u p_{sw} *) = \sum_{v \in \mathbf{V}} \sum_{p \in P_{sw}} cnt(u p v).$$

This score captures the proportion of times  $u$  dominates all other words through the patterns given, relative to the proportion of times  $u$  is dominated by all other words through the pattern. The results of this ranking is displayed in table 6 in the line “Markov heuristic.”

Now we make the intuition precise. Suppose we have a simple language  $\mathbf{L}$  made up only of phrases of the form “word pattern word” for every word in the unigram set  $\mathbf{V}$  and every pattern in table 1, that is we have:

$$\mathbf{L} = \{u p v \mid u, v \in \mathbf{V}, p \in \mathbf{P}_{sw} \cup \mathbf{P}_{ws}\}.$$

If we can confidently approximate likelihood of each phrase from  $\mathbf{L}$  based on N-gram corpus evidence alone then we are done. But because data is sparse, we must fill in the missing counts by assuming the phrases in  $\mathbf{L}$  is generated by this Markov process involving two random variables,  $V$  whose value range over vocabulary  $\mathbf{V}$ , and  $P$  ranging over the patterns in table 1. The speaker selects a word  $u$  according to some distribution  $\mathcal{D}_V$  over  $\mathbf{V}$ , then conditioned on the fact that  $u$  is drawn, a phrase  $p$  is drawn according to the conditional distribution  $\mathcal{D}_{P|V}$ .

Finally, conditioned on the fact that  $p$  is drawn, a word  $v$  is sampled from  $\mathcal{D}_{V|P}$ . The attentive reader will observe that this crude model does not respect word order. In the phrase “not great (,) just good”, our model would generate the phrase “good not just great”. Surprisingly this model works well enough to outperform Bansal’s method. Now the probability of a phrase is:

$$\mathcal{D}_L = \frac{\mathcal{D}_V \mathcal{D}_{P|V} \mathcal{D}_{V|P}}{Z},$$

where  $Z$  is an appropriate normalization constant. But since we are only interested comparing the relative likelihood of phrases,  $Z$  does not need to be computed. So we have:

$$\mathcal{D}_L = Pr[u p v] \propto Pr[u]Pr[p|u]Pr[v|p], \quad (1)$$

where:

$$\begin{aligned} Pr[V = u] &= \frac{cnt(u)}{cnt(*)} \\ Pr[P = p|V = u] &= \frac{cnt(u p *)}{cnt(u *)} \\ Pr[V = v|P = p] &= \frac{cnt(* p v)}{cnt(* p *)} \end{aligned}$$

where  $cnt(*) = \sum_{x \in \mathbf{V}} cnt(x)$ . The first distribution is approximated by the one-gram corpus, the second and third distribution by four and five grams. In the interest of not computing normalization constant whenever possible, we put the following crude bound on  $cnt(u *)$ :

$$cnt(u *) = \sum_x count(u x) \leq cnt(u),$$

where  $x$  is ranges over all suffixes of length three or four. So (1) becomes:

$$Pr[u \ p \ v] \propto \frac{cnt(u \ p \ *) \cdot cnt(* \ p \ v)}{cnt(*) \cdot cnt(* \ p \ *)}. \quad (2)$$

Now define the probability that  $u$  is stronger than  $v$  under  $\mathcal{D}_L$  as:

$$\begin{aligned} Pr[u > v] &= Pr[u \ P_{sw} \ v \ \text{or} \ v \ P_{ws} \ u] \\ &= Pr[u \ P_{sw} \ v] + Pr[v \ P_{ws} \ u] \\ &= \sum_{p \in P_{sw}} Pr[u \ p \ v] + \sum_{p \in P_{ws}} Pr[v \ p \ u], \end{aligned}$$

and similarly for  $v > u$ . We decide  $u$  is stronger than  $v$  if:

$$\begin{aligned} Pr[u > v] &\geq Pr[v > u] \\ \implies & \\ \frac{cnt(u P_{sw} *) \cdot cnt(* P_{sw} v)}{cnt(*) \cdot cnt(* P_{sw} *)} + \frac{cnt(v P_{ws} *) \cdot cnt(* P_{ws} u)}{cnt(*) \cdot cnt(* P_{ws} *)} & \\ \geq & \\ \frac{cnt(v P_{sw} *) \cdot cnt(* P_{sw} u)}{cnt(*) \cdot cnt(* P_{sw} *)} + \frac{cnt(u P_{ws} *) \cdot cnt(* P_{ws} v)}{cnt(*) \cdot cnt(* P_{ws} *)} & \\ \implies & \\ \frac{cnt(u P_{sw} *) \cdot cnt(* P_{sw} v)}{cnt(* P_{sw} *)} + \frac{cnt(v P_{ws} *) \cdot cnt(* P_{ws} u)}{cnt(* P_{ws} *)} & \\ \geq & \\ \frac{cnt(v P_{sw} *) \cdot cnt(* P_{sw} u)}{cnt(* P_{sw} *)} + \frac{cnt(u P_{ws} *) \cdot cnt(* P_{ws} v)}{cnt(* P_{ws} *)}. \quad (3) \end{aligned}$$

Note the normalization constant  $cnt(*)$  drops out, and there is a qualitative

symmetry in (3) that echos intuition. Since (3) does not output cycles, ranking is done by topological sort; results are reported in table 6 under “Markov pairwise approximate.” Next, we combine the approximate value from (3) with those directly observed in the corpus. If for some adjective pair  $(u, v)$  we observe any one of the following values:  $u p_{sw} v$ ,  $u p_{ws} v$ ,  $v p_{sw} u$ , or  $v p_{ws} u$ , then we can define the probability that  $u > v$  under  $\mathcal{D}_L$  as:

$$Pr[u > v] = \frac{cnt(u P_{sw} v) + cnt(v P_{ws} u)}{Z} \quad (4)$$

where

$$Z = cnt(u P_{sw} v) + cnt(v P_{ws} u) \\ + cnt(v P_{sw} u) + cnt(u P_{ws} v),$$

and

$$cnt(u P_{sw} v) = \sum_{p \in P_{sw}} cnt(u p v).$$

Now to rank  $u$  against  $v$ , we compute (4) if possible, otherwise we approximate the probability that  $u > v$  using (3). Since the ordering over each pair of adjectives is decided separately using (3) or (4), cycles do exist and transitivity must be enforced. First, we consider a simple integer linear programming formulation, given  $N$  adjectives in a cluster where a ranking is known to exist, and define:

$$P_{uv} = \frac{Pr[u > v]}{Pr[v > u]},$$

so that if  $u \geq v$  under  $\mathcal{D}_L$  then  $P_{uv} \geq 1$ :

**Maximize**

$$\sum_{u,v \in \{1, \dots, N\}} P_{uv} \cdot s_{uv} + P_{vu} \cdot (1 - s_{vu})$$

**s.t**

$$(1 - s_{uv}) + (1 - s_{vw}) \geq (1 - s_{uw}),$$

$$\forall u, v, w \in \{1, \dots, N\}.$$

Thus the objective encourages  $s_{uv} = 1$  if  $u > v$ , and  $s_{uv} = 0$  otherwise. Overall the objective gives precedent to those pairs where  $u$  dominates  $v$  the most, while the constraints enforce transitive properties of the ordering. This mixed approximation of  $Pr[u > v]$  and *ILP* gives the best results on Bansal’s data, see tables 4 and 5 under “Markov pairwise mixed ILP.” Lastly, in the interest of exploring the trade off between precision versus sparsity of data, we use our approximation of frequency of  $u$   $p$   $v$  in Bansal’s formulation. Define:

$$\begin{aligned} score(u, v) = & Pr[v P_{sw} u \text{ or } u P_{ws} v] \\ & - Pr[u P_{sw} v \text{ or } v P_{ws} u], \end{aligned}$$

so that  $score(u, v) > 1$  if  $u < v$  under  $\mathcal{D}_L$ , thereby conforming to Bansal’s formulation of the score function in terms of sign. See “Markov pairwise approximate MILP.”

## 7.7 Results

Table 11 reports all results from Mohit’s MILP approach, and the ILP reformulation. now we discuss the advantages and drawbacks of each method, specifically



<b>Gold</b>	<b>MILP</b>	<b>Markov ILP</b>
(cool, chilly) < unfriendly < hostile	unfriendly < cool < hostile < chilly	cool < chilly < unfriendly < hostile
strong < intense < terrible < overwhelming < violent	strong < intense < terrible < overwhelming < violent	intense < strong < terrible < violent < overwhelming
high < higher < soaring	soaring < high < higher	soaring < high < higher

Table 11: The top row displays an example where Bansal’s MILP fails to output the correct order ( $\tau = -0.18$ ) but Markov Mixed ILP output the correct order modulo ties ( $\tau = 0.91$ ). The middle row is an example where Bansal’ MILP correctly predicted the ranking despite sparse data, only six out of twenty pairs had any N-gram hits. Using Markov assumption the missing data was filled in, but at the cost accuracy. The bottom row shows an instance where both methods fail because there is overwhelming copora evidence that higher is more intense than soaring, revealing the limitations of the pattern-based approach.

the settings in which each method succeeded over another. Looking at table 12 line “Markov pairwise mixed ILP”, we see that the integer linear programming formulation with missing data approximated by (3) enjoy the highest accuracy across all measures. Thereby validating the hypothesis that the Markov condition is sufficiently reasonable assumption to approximate frequency of unseen phrases. Furthermore, it is interesting that “Markov pairwise approximate” performs almost as well as Bansal’s MILP method, even though it relies on the approximation given by (3) alone. This indicates that the linguistic patterns specified is an accurate enough of a predictor of relative adjective strength, so that even crude approximations are powerful predictors. In terms of data complexity, we stress that both (3) and (4) can be calculated from the same data collected by Bansal to compute the score in section 3.1, suggesting that the aforementioned data holds more information about the ordering of adjectives than is used by Bansal’s MILP method. Most notably, we outperform Bansal’s method without relying on annotator labels of synonyms, which a very restrictive condition indeed.

Now we wish to highlight some specific instances where “Markov pairwise mixed ILP” did well relative to Bansal’s MILP method, and where it faltered. In the interest of fairness, we compare the output of our methods from table 6 only. Bansal’s MILP method output 11 rankings where Kendall’s tau score was less than zero, while Markov pairwise mixed ILP only output six. Notably, five out of six bad ILP clusters also appear in Bansal’s method, one example of this is given in the bottom row of table 4. In all these cases there is strong corpora evidence leading to the wrong conclusion, but in 94% of the cases linguistic-pattern points in the right direction. In conclusion, not only is our integer linear program formulation simpler than that of Bansal’s, we also outperform the state of the art across all measures in a nontrivial manner. We do so without relying on additional annotations during ranking, which is a very restrictive assumption. Finally, the sample complexity of

Method	Pairwise Accuracy	Avg. $\tau$	Avg. $ \tau $	Avg. $\tau'$	Avg. $ \tau' $
Inter-Annotator Agreement	78.0%	0.67	0.76	N/A	N/A
Gold Standard	100.0%	0.90	0.90	0.90	0.90
MILP reproduced	68.0%	0.55	0.64	0.41	0.54
Markov heuristic	65.0%	0.43	0.61	0.31	0.52
Markov pairwise approximate	70.0%	0.53	0.63	0.41	0.54
Markov pairwise mixed ILP	72.0%	0.57	0.64	0.44	0.54
Markov MILP	70.0%	0.53	0.65	0.41	0.56

Table 12: Main results using Bansal’s patterns. Note  $\tau$  refers to kendall’s  $\tau$  with ties, while  $\tau'$  refers to the variant where ties are not considered.

our approach is comparable to that of Bansal’s.

## 8 Ranking Using Network Centrality

In this chapter, we report a failed attempt at ranking adjectives using rank centrality.

### 8.1 Introduction

This measure differs notably from the previous chapter in that (1) the previous chapters considers local information, here we consider the whole graph; and (2) there is large precedent in literature for these measures.

### 8.2 Motivation

We hypothesize that both the relative strength of vertices incident on  $s$  versus those of  $t$ , and the relative strength of vertices pointed to by  $s$  and  $t$  are important. For example, suppose both "good" and "great" have equal number of edges incident upon it, but great is a paraphrase of "much better", which we know to be stronger than "good", then we should rank "great" as more intense than "good". What we are interested in measuring is the "centrality" of an adjective. In fact there is a family of measures constructed for just this task, and all involves defining some random walk over the graph and computing the stationary distribution over all vertices under this walk. In this chapter we consider two in particular: PageRank and Personalized PageRank. PageRank outputs a total ordering over all vertices, ranked by their centrality measure. Personalized PageRank, on the other hand, does not assume a unique ordering of all words in the graph, instead it computes how important some vertex  $t$  is given the perspective of a particular vertex  $s$ . And note the ranking of  $t$  may be different for each  $s$ . Now we will introduce some notation and review these measures.

### 8.3 Notation

Consider the directed graph  $\mathbf{G}$  with vertices  $\mathbf{V}$  and edges  $\mathbf{E}$ , the vertices in this graph be denoted by  $s$ ,  $t$ , or  $r$ . An edge from  $s$  to  $t$  will be denoted  $(s, t)$ . The adjacency matrix  $\mathbf{A}$  of the graph is defined as:

$$\mathbf{A}_{st} = \begin{cases} 1 & \text{if } (s, t) \\ 0 & \text{otherwise.} \end{cases}$$

Note  $\mathbf{A}_{st}$  is the  $st$ 'th entry of matrix  $\mathbf{A}$ . Furthermore, the degree matrix  $\mathbf{D}$  is the diagonal matrix with:

$$\mathbf{D}_{st} = \sum_t \mathbf{A}_{st}.$$

Finally we can construct the probability transition matrix  $\mathbf{W}$  with:

$$\mathbf{W} = \mathbf{D}^{-1}\mathbf{A}.$$

Note for every  $st$ 'th entry in  $\mathbf{W}$  we have:

1.  $\mathbf{W}_{st} \in [0, 1]$
2.  $\sum_t \mathbf{W}_{st} = 1$ .

Following the tradition of [citation], we will use  $\pi_s$  to denote a distribution over all vertices in  $\mathbf{G}$ . We also overload the letter  $s$  to denote both the vertex in  $\mathbf{V}$  and a distribution over vertices where the entry  $s$  has probability one, and all other vertices have probability zero. Now we will define PageRank and Personalized PageRank. In the interest of highlighting their similarities, we will use similar language in both definitions.

## 8.4 PageRank

PageRank of vertex  $s$  is the probability that a random walk of some length starting from an arbitrary vertex will terminate at  $s$ . Formally,  $X_i$  be a random variable ranging over distributions over vertices in  $\mathbf{G}$ , then the ordered sequence  $(X_0, X_1, \dots, X_L)$  is a random walk of length  $L$  starting from  $X_0 = s$ .  $L$  follows a geometric distribution where  $\Pr[L = l] = (1 - \alpha)^l \alpha$ . That is to say the random walker starts at some arbitrary vertex and with probability  $\alpha$ , transition to an out-neighbor  $t$  according to  $W_{st}$ , and with probability  $1 - \alpha$ , “teleport” to an arbitrary vertex on the graph. The PageRank vector  $\pi$  is the solution to the expression:

$$\pi = \widetilde{\mathbf{W}}\pi,$$

where  $\widetilde{\mathbf{W}}$  is the “Google” matrix:

$$\widetilde{\mathbf{W}} = \alpha \mathbf{W} + (1 - \alpha) \frac{1}{n} \mathbf{J},$$

and  $\mathbf{J} = \mathbf{1} \cdot \mathbf{1}^T$  is the all one matrix. Note the expression  $\alpha \mathbf{W}$  above captures the random surfing behavior, while the second expression is the teleportation factor. The PageRank of a vertex  $s$  is the probability this random walk terminates at  $s$ :  $\pi[s] = \Pr[X_L = s]$

## 8.5 Personalized PageRank

A personalized PageRank of vertex  $t$  relative to  $s$  is defined as the probability that a random walk of the appropriate length starting from  $s$  will terminate at  $t$ . Similar to PageRank, let  $X_i$  be a random variable ranging over distributions over vertices in  $\mathbf{G}$ , then the ordered sequence  $(X_0, X_1, \dots, X_L)$  is a random walk of length  $L$  starting from  $X_0 = s$ .  $L$  follows a geometric distribution where  $\Pr[L = l] = (1 - \alpha)^l \alpha$ . In other words, the random walk starts at  $s$  and with probability

$1 - \alpha$  continue to a random neighbor of the current vertex; and with probability  $\alpha$  terminates at the current vertex  $s$ . Again at each vertex  $s$ , the random neighbor  $t$  is chosen with probability  $\mathbf{W}_{st}$ . The personalized PageRank vector  $\pi_s$  with respect to  $s$  is the solution to the expression:

$$\pi_s = \alpha s + (1 - \alpha)\pi_s \mathbf{W}.$$

The PPR of vertex  $t$  with respect to  $s$  is the probability we terminate at  $t$ :

$$\pi_s[t] = \mathbf{Pr}[X_L = t].$$

Solving for  $\pi$  and  $\pi_s$  is well studied but beyond the scope of this thesis, the interested reader should refer to [citation] for a thorough treatment.

## 8.6 Pairwise Comparison Using Centrality Measures

In this brief section we define how to assign  $\mathbf{Pr}[s < t]$  using the two centrality measures defined above.

**Definition 8.1.** Given two adjectives  $s$  and  $t$ ,  $s$  is less intense than  $t$  under PageRank over the graph if  $\pi[s] < \pi[t]$ . Furthermore by fiat, we have for some suitable  $\epsilon$ :

$$\mathbf{Pr}[s < t] = \begin{cases} \frac{1}{2} - \epsilon & \pi[s] < \pi[t] \\ \frac{1}{2} + \epsilon & \text{otherwise} \end{cases}$$

**Definition 8.2.** Given two adjectives  $s$  and  $t$ ,  $s$  is less intense than  $t$  under Personalized PageRank over the graph if  $\pi_s[t] > \pi_t[s]$ . And similar to above, we can define:

$$\Pr[s < t] = \begin{cases} \frac{1}{2} - \epsilon & \pi_s[t] > \pi_t[s] \\ \frac{1}{2} + \epsilon & \text{otherwise} \end{cases}$$

Note in both cases we “forget” the actual probability induced by the random walk, and define the probability by fiat. This is due to our previous claim in the baseline chapter that a slight edge suffices. The second reason is mostly practical: due to the quality of our data, in practice it is prudent to not assign too much credence to the values output by the centrality measures beyond the gross direction they indicate.

## 8.7 Constructing the Transition Matrix

This section outlines two ways of constructing the transition matrix  $\mathbf{W}$  from the multi directed graph. We only consider one variation here: whether to include the number of edges between vertices. Let us denote  $\tilde{\mathbf{E}}$  as the set of edges in the multi-directed graph so that  $(s, t, v) \in \tilde{\mathbf{E}}$  signifies an edge from  $s$  to  $t$  labeled with adverb  $v$ . In the first construction, we ignore the number of edges between two vertices  $s$  and  $t$ , and define the probability of transitioning from  $s$  to  $t$  as one over the number of out-neighbors of  $s$ :

$$\mathbf{W}_{st} = \frac{1}{|\{(s, t, v) : (s, t, v) \in \tilde{\mathbf{E}}\}|}. \quad (1)$$

We will call this construction the out-neighbor construction. In the second construction it is the number of edges from  $s$  to  $t$ , over the total number of out edges of  $s$ :

$$\mathbf{W}_{st} = \frac{|\{(s, t, v) \in \tilde{\mathbf{E}}\}|}{|\{(s, t, v) : (s, t, v) \in \tilde{\mathbf{E}}\}|}. \quad (2)$$

This construction will be referred to as the out-edge construction



## 8.8 Results

Similar to the previous chapter, we test on the set of pairs for which no data exists. And once again we remind the reader that we wish to be in the low to mid 80 percent for pairwise accuracy across all test sets. The results are reported in the tables below. We varied the teleportation constant  $\alpha$  from 0.1 to 0.9 in 0.1 increments. In practice we noted  $\alpha = 0.2$  performed best for PageRank, and  $\alpha = 0.8$  performed best for Personalized PageRank. In the interest of space only results from the best  $\alpha$  are reported. Scanning across the tables, we see that the centrality measures do outperform the random baseline, but their performance as a whole do not meet our minimum 80% criteria. However, the two measures performed significantly differently across the data sets, and this difference merit some brief comments. The two most salient conclusions are:

1. PageRank out performed Personalized PageRank as a whole, falsifying our hypothesis that assuming there is no total ordering over adjectives would improve results.
2. out-edge construction out performed out-degree construction in some cases but not others, suggesting that the exact construction of the adjacency matrix does not matter.

Detailed comments are found underneath each table. All in all, PageRank appears to show promise on some labeled sets, while Personalized PageRank does not perform. Overall, both measures are unsatisfactory.

	N-gram no data		PPDB no data		PPDB + N-gram no data	
Test set	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$
Mohit	47.0%	-0.05	48.0%	-0.03	48.0%	-0.04
Turk	40.0%	-0.19	37.0%	-0.25	41.0%	-0.18
BCS	14.0%	-0.73	15.0%	-0.70	15.0%	-0.70

Table 13: Results across all datasets for uniform baseline coupled with reverse lexicographic sorting, reprinted here for convenience.

	N-gram		PPDB		PPDB + N-gram	
Test set	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$
Mohit	48.0%	-0.03	<b>64.0%</b>	<b>0.28</b>	55.0%	0.10
Turk	<b>57.0%</b>	<b>0.13</b>	53.0%	0.05	52.0%	0.04
BCS	73.0%	0.45	<b>77.0%</b>	<b>0.54</b>	70.0%	0.39

Table 14: PageRank using graph constructed using out-edge expression. Note BCS performed best on PPDB, this is not surprising since this is how the PPDB dataset is generated. The Turk set performed best on the N-gram set, which is curious since the Turk set was constructed with respect to the PPDB adjective set. Finally we see Mohit’s set performed best on the PPDB data set, but it appears as if there is a slight negative relationship between PageRank values and adjective strength on the N-gram set.

	N-gram		PPDB		PPDB + N-gram	
Test set	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$	Pairwise	Avg. $\tau$
Mohit	49.0%	-0.02	63.0%	0.25	<b>59.0%</b>	<b>0.19</b>
Turk	<b>58.0%</b>	<b>0.16</b>	51.0%	0.01	54.0%	0.08
BCS	<b>72.0%</b>	<b>0.43</b>	63.0%	0.25	72.0%	0.45

Table 15: PageRank with out-neighbor construction. Note in theory this construction should be more susceptible to “confusing” signals due to polysemy, however in practice it actually performs better on some annotated sets.

	N-gram		PPDB		PPDB + N-gram	
<b>Test set</b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>
Mohit	60.0%	0.20	39.0%	-0.22	63.0%	0.25
Turk	35.0%	-0.30	62.0%	0.23	58.0%	0.15
BCS	29.0%	-0.42	61.0%	0.22	57.0%	0.14

Table 16: Personalized PageRank with out-edge construction. Note it performs significantly worse than PageRank. In particular there appears to be a negative relationship between PPR and adjective strength on some data sets, and positive relationship on others. Furthermore, observe the negative correlations appear on the N-gram dataset where we do not observe paths over several vertices, this makes PPR highly unreliable. It also appears for Mohit’s set on the PPDB data, which was not constructed with Mohit’s gold set in mind.

	N-gram		PPDB		PPDB + N-gram	
<b>Test set</b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>	<b>Pairwise</b>	<b>Avg. <math>\tau</math></b>
Mohit	61.0%	0.22	39.0%	-0.21	62.0%	0.24
Turk	34.0%	-0.31	59.0%	0.19	49.0%	-0.01
BCS	32.0%	-0.37	63.0%	0.25	68.0%	0.35

Table 17: Personalized PageRank with out-neighbor construction. Note in general it is not worse than the out-edge construction. Finally, the negative correlation appear here in the same test set and gold set.

## References

- Gerard de Melo and Mohit Bansal. 2013. Good, great, excellent: Global inference of semantic intensities. *Transactions of the Association for Computational Linguistics*, 1:279–290.
- Joo-Kyung Kim and Marie-Catherine de Marneffe. 2013. Deriving adjectival scales from continuous space word representations. In *EMNLP*, pages 1625–1630.
- Peter F Schulam and Christiane Fellbaum. 2010. Automatically determining the semantic gradation of german adjectives. In *KONVENS*, pages 163–167. Cite-seer.
- Vera Sheinman and Takenobu Tokunaga. 2009. Adjscales: Visualizing differences between adjectives for language learners. *IEICE TRANSACTIONS on Information and Systems*, 92(8):1542–1550.
- Vera Sheinman, Takenobu Tokunaga, Isaac Julien, Peter Schulam, and Christiane Fellbaum. 2012. Refining wordnet adjective dumbbells using intensity relations. In *Sixth International Global Wordnet Conference*, pages 330–337.