

1 Regression

1.1 Introduction

In the previous two chapters we attempted to construct generic measures that may work for every graph, but in the end certainly did not perform on this graph. In this chapter we will construct a measure specifically for this dataset. Towards this end, we will use 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.

1.2 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.

1.3 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 regres-

sion. 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 β by:

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

In our setting, we experiment with two penalties on β : 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.

1.4 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 ϕ 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 ϕ , then we can use this definition.

Definition 1.1. Given words s and t , and their representation $\mathbf{x} = g(\phi(s), \phi(t))$, and let:

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

then we can define:

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

for an appropriate threshold δ . Again we discard the value of \hat{y} and define the probability by fiat.

1.5 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 ϕ 's.

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

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

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 ϕ^{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 ϕ^{in} and ϕ^{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 ν '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 ν '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 ν 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 ν . After g has been applied, we also experimented with using the raw counts, versus normalizing the entries of ϕ 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 1: 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 $\tilde{x} = \frac{x}{\max(x_1, \dots, x_n)}$

1.6 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 50 most frequent adverbs incident on the vertex and outgoing from the vertex. 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	
Gold Set	Pairwise	Avg. τ	Pairwise	Avg. τ
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 2: . 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 α of 0.9, and l_1 of 0.1.