

Understanding Scanned Receipts

Eric Melz

300 S Reeves Dr.
Beverly Hills, CA 90212
eric@emelz.com

Abstract

Tasking machines with understanding receipts can have important applications such as enabling detailed analytics on purchases, enforcing expense policies, and inferring patterns of purchase behavior on large collections of receipts. In this paper, we focus on the task of Named Entity Linking (NEL) of scanned receipt line items. Specifically, the task entails associating shorthand text from OCRd receipts with a knowledge base (KB) of grocery products. For example, the scanned item “STO BABY SPINACH” should be linked to the catalog item labeled “Simple Truth OrganicTMBaby Spinach”. Experiments that employ a variety of Information Retrieval techniques in combination with statistical phrase detection shows promise for effective understanding of scanned receipt data.

1 Introduction

Tasking machines with understanding receipts can have important applications such as enabling detailed analytics on purchases, enforcing expense policies, and inferring patterns of purchase behavior on large collections of receipts. In this paper, we focus on the task of Named Entity Linking (Hachey et al., 2012) of scanned receipt line items. Specifically, the task entails associating shorthand text from OCRd receipts with a knowledge base (KB) of grocery products. For example, the scanned item “STO BABY SPINACH” should be linked to the catalog item labeled “Simple Truth OrganicTMBaby Spinach”.

2 Related Work

A literature review reveals virtually no published work in this specific domain. While there is a body of work researching text extraction from scanned receipts (e.g. Huang et al., 2019), the work is primarily focused on Named Entity Recognition

(NER) instead of Named Entity Linking (NEL). That is, systems are considered successful if they can identify text items such as store locations, totals, etc, but they are not evaluated with respect to the interpretation of the extracted text.

Although no papers exist on linking scanned entities, there is literature in other areas that appear potentially relevant to the subject task. This includes work on general-purpose techniques for building abbreviation dictionaries, acquisition of medical abbreviations (e.g., “COPD” → “Chronic Obstructive Pulmonary Disorder”), and normalization of social media content (e.g., “ur coooool” → “you are cool”). The follow sections summarize a few papers in these areas.

2.1 Language Independent Acquisition of Abbreviations

(Glass et al., 2017) describe a language-independent technique for acquiring abbreviations and their expansions, by exploiting Wikipedia redirect and disambiguation pages. They begin by motivating the acquisition of abbreviations, noting that the explosion of social media has made the need for abbreviations increasingly important. They also note that a token such as “ACE” could have multiple expansions, including “accumulated cyclone energy” and “American Council on Education” in addition to the word “ace” (as in “Ace of spades”).

The authors present related work, noting that most of the previous work for abbreviation detection and expansion extraction has been in the domain of English biomedical text. A common strategy is to identify occurrences where an abbreviation is explicitly paired with its expansion for example through a pattern involving a parenthetical such as <short form> (<long form>) or <long form> (<short form>). Other approaches consider the contexts of short form and

long form occurrences, pairing short forms with long forms according to their distributional similarity by measuring the cosine of their context vectors. Another approach uses supervised learning, considering features such as string similarity and other characteristics of the short and long forms.

The authors work is based on previous work by (Jacquet et al., 2014) who describe a technique for mining abbreviations by making use of Wikipedia redirection pages. The authors observe that, due to the use of only redirect pages for the gold standard annotation, a shortcoming of the prior work is that each abbreviation only has a single expansion even though multiple different expansions are possible for some of the abbreviations. To remedy this shortcoming, the authors propose mining disambiguation pages in addition to redirect pages to gather multiple possible long-form expansions.

The authors mine redirect and disambiguation pages for abbreviations, while applying several rules such as (a) Short forms are restricted to ten characters or less, (b) At least half of the short-form characters must be upper case, and (c) The long-form must be at least twice as long as the short form, with at least two tokens. They generate candidate expansions and then score the expansions. Scoring occurs by computing features for synonym similarity, topic similarity, and surface similarity. Synonym similarity means that one term can be replaced with another while preserving the meaning of the sentence and is assessed using word embeddings using word2vec (Mikolov et al., 2013). Topical relatedness means that two terms occur in the same sorts of documents, and is assessed using Latent Semantic Analysis (Deerwester et al., 1990). Surface similarity is the overlap in the surface forms of the terms by computing the best possible alignment between a short form and a long form. The three similarity scores are combined using a logistic regression model.

The authors compare their system with a previous system developed by (Schwartz and Hearst, 2003) that extracts abbreviations using parentheses based patterns. The metric used to compare systems is Area Under the Precision/Recall curve. Without the scoring extensions, the 2 systems are comparable: the Schwartz and Hearst system has an AUC of 0.359 and the Candidate System has an AUC of 0.324. However, by adding the alignment and embedding scoring extensions, the Can-

didate Systems performance improves to an AUC of 0.480.

2.2 Clinical Abbreviation Expansion

(Liu et al., 2015) describe a system for identifying clinical abbreviation expansions. They note that abbreviations are heavily used in medical literature and documentation. In notes written by physicians, high workloads and time pressure intensify the need for using abbreviations. This is especially true within intensive care medicine, where its crucial that information is expressed in the most time efficient manner to provide time-sensitive care to critically ill patients. Within the arena of medical research, abbreviation expansion using NLP can enable knowledge discovery and has the potential to improve quality of care.

The author's system works as follows. Word embeddings are trained using word2vec (Mikolov et al., 2013). The material used to train embeddings consists of medical texts such as articles, journals, and books, in addition to hand-written Intensive Care notes. To generate expansions for abbreviations in the hand-written notes, abbreviations are extracted from the notes, and then matched against a domain-specific abbreviation knowledge base. From this list of expansions, embedding vectors are retrieved for the abbreviation and candidate expansion. A similarity score is computed for each (abbreviation, expansion) pair, producing a ranked list of candidates expansions.

To test the performance of the system, a ground-truth dataset is produced by having physicians manually expand and normalize the handwritten notes. The authors compare their model against several baselines. For example, one baseline chooses the highest rated candidate expansion in the domain specific knowledge base. Comparing accuracy of the authors system against the baselines results in a 50%+ increase. For example the rating baseline has an accuracy of 21% and the authors system has an accuracy of 83%.

2.3 Social Media Text Normalization

(Lourentzou et al., 2019) present a Sequence to Sequence (Seq2Seq) model for normalizing social media text. They observe that social media texts have an enormous amount of variation, and that text normalization systems that rely on surface or phonetic representations may be ill-equipped to handle such variability. To rectify this situation, they propose a hybrid word-character Seq2Seq

model with attention. This type of model has been successfully applied to tasks such as machine translation, and has promise for text normalization.

The authors frame the task of text normalization as mapping an out-of-vocabulary (OOV) non-standard word to an in-vocabulary (IV) standard word that preserves the meaning of the sentence. The non standard forms in user generated content include misspellings (defenitely \rightarrow definitely), phonetic substitutions (2morrow \rightarrow tomorrow), shortening (convo \rightarrow conversation), acronyms (idk \rightarrow i dont know), slang (“low key”, “woke”), emphasis (coooooool \rightarrow cool), and punctuation (doesnt \rightarrow doesn’t).

The authors note that lexicon-based approaches are not able to handle social media text properly. String similarity, such as edit distance, does not work on non-standard words where the number of edits is large, for example abbreviation. Additionally, systems that rely on candidate generation and scoring are limited in that they are not able to handle multiple normalization errors at once, e.g., spelling errors on an acronym. The authors suggest that using end-to-end neural models, particularly Seq2Seq models can deal with these shortcomings.

The authors train a bidirectional word-based Seq2Seq model to translate unnormalized text to normalized texts. OOV words are trained using a character-based Seq2Seq model. The dataset is enhanced by synthetically generating negative examples based on common normalization transformations. The network is trained on source sequences and target sequences. An example source is “got exo to share, u interested? Concert in hk !”, with a corresponding target of “got extra to share, are you interested? Concert in hong kong !”.

The authors present results for several variations of the model, including a word-level Seq2Seq model and the hybrid word-char Seq2Seq model. The best score is an F_1 score of 83.94 on the hybrid word-char Seq2Seq model.

3 Data

For this task, we need a dataset which includes scanned receipt product mentions (e.g., “BRHD CHEESE”) and the corresponding product entities (e.g., “Boar’s Head Monterey Jack with Jalapeno Pre-Sliced Cheese”). A brief web search revealed that no such publicly available dataset exists. To

obtain a dataset, we built our own by scraping a grocery store website that contains purchase data. Specifically, we use our personal loyalty account with Ralph’s (a subsidiary of Kroger) to obtain representations of scanned receipts along with corresponding web pages that contain fully-resolved entities. As an example, Figure 1 shows an instance of a receipt.

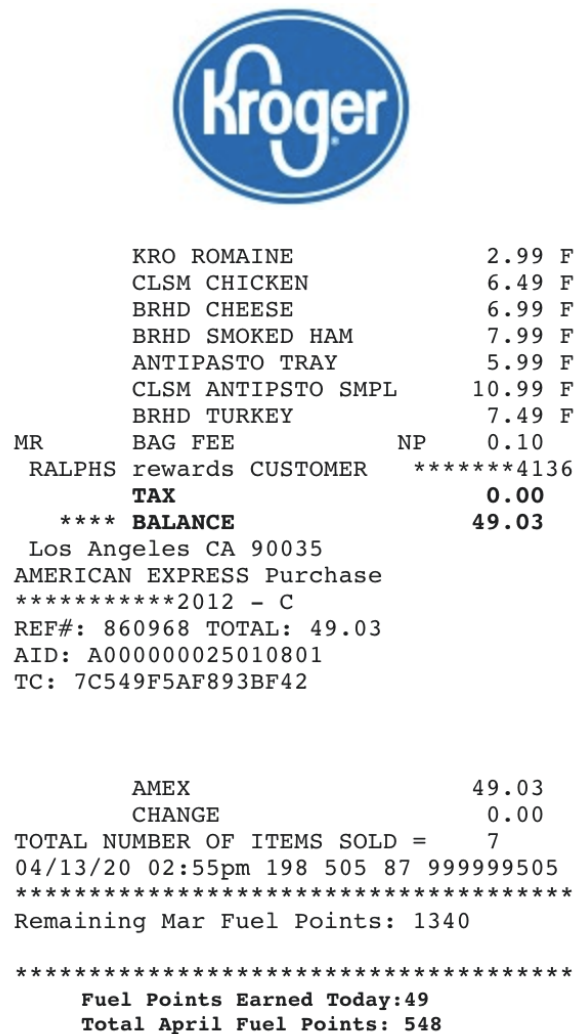


Figure 1: Scanned Receipt

Figure 2 shows part of the corresponding web page which contains linked representations of the purchased items.

We scrape both the text content of the raw receipts and the user-friendly web rendering, then join the raw receipt data with the corresponding web data. This produces a JSON structure per receipt. A sample of the JSON is shown in Figure 3. The “raw” field represents the product mention, and the “web” field represents the label associated with the entity. The “id” field is scraped from the

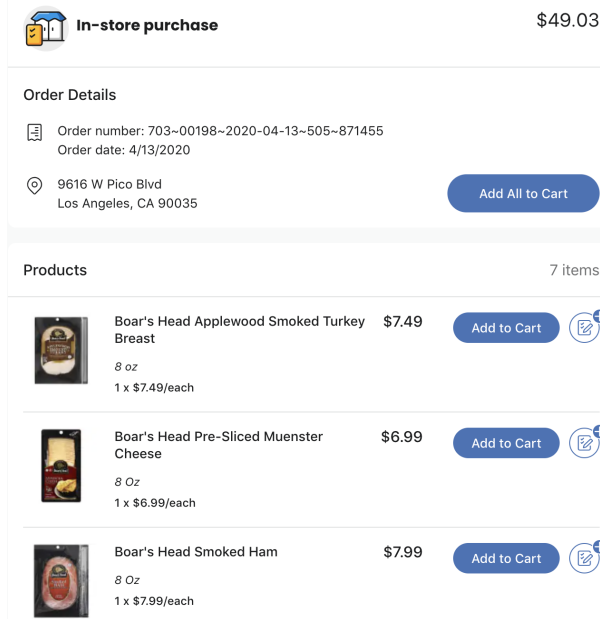


Figure 2: Web Receipt

web page and can be used as a succinct identifier for the entity.

```
{
  "web": "Avocado - Extra Large",
  "raw": "AVOCADO HASS",
  "id": "0000000004770"
},
{
  "web": "Beyond Meat Hot Italian Plant-Based Sausage",
  "raw": "BYND SSG HT ITLN",
  "id": "0085262900475"
},
{
  "web": "Eggland's Best Organic Grade A Large Brown Eggs",
  "raw": "EGGLANDS BEST EGGS",
  "id": "0071514171682"
},
{
  "web": "Fancy Feast Grilled Beef Feast in Gravy Wet Cat Food",
  "raw": "FFST CAT FOOD",
  "id": "0005000004070"
},
{
  "web": "Fancy Feast Minced Turkey Feast in Sauce Wet Cat Food",
  "raw": "FFST CAT FOOD",
  "id": "00050000043494"
},
}
```

Figure 3: JSON representation of joined “raw” and “web” data

The dataset consists of 65 scraped receipts, producing 711 non-unique line items, and 296 unique line items. All data and code for these experiments are available on Github (Eric Melz, 2020).

4 Methodology

To evaluate model performance, we gather unique mentions and measure the accuracy of predicting entities. This can be conceived as a multi-class classification task where the entities to be pre-

dicted are the classes. An alternative metric would be to use a macro-average F_1 score, but this is overkill for this specific experiment setup since there is a uniform distribution across classes: each class is represented by exactly one test instance. To be concrete, heres an example. Suppose we have the following 2 unique mentions:

- BRHD CHEESE
- AVOCADO

Further suppose that the entity for BRHD CHEESE was correctly predicted as Boar’s Head Monterey Jack with Jalapeno Pre-Sliced Cheese, and the prediction for AVOCADO yielded nothing. The first prediction is a “hit” and the second is a “miss”. Dividing the total hits by the total number of predictions, we obtain an accuracy of $(1 + 0) / 2 = 0.5$.

Note that the mention representations contain much less information than the entity labels. In the above example, BRHD CHEESE is matched with “Boar’s Head Monterey Jack with Jalapeno Pre-Sliced Cheese”, but also could have been matched to “Boar’s Head Spicy Cheddar Cheese”, or a number of other types of Boars Head cheese. To account for this ambiguity, a prediction is counted as a hit if it is any of the possible resolutions of the product mention. In the previous example, both of the long descriptions would be considered hits for BRHD CHEESE.

The dominant modeling paradigm for the entity linking system we use is Information Retrieval. A baseline model indexes entity labels using the Lucene (Apache Software Foundation) IR engine. Lucene provides a toolkit of tokenizers and token analyzers, enabling many strategies for matching text. The most basic setup uses strict matching on tokens, providing a good baseline model. Subsequent experiments tune of the by selecting more sophisticated IR and NLP techniques such as using wildcard queries, phrase detection, etc.

5 Experiments

In fhte following sections, we report a series of experiments that leverage Lucene. We index the web versions of entities and query the index with the raw line items derived from scanned receipts.

Lucene’s default scoring formula is BM25 (Robertson, 2009). BM25 is based on a bag-of-words approach. The score of a document D given a query Q which contains the words

q_1, \dots, q_n is given by:

$$\text{score}(D, Q) = \frac{\sum_{i=1}^n \text{IDF}(q_i) \cdot f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{\text{avgdl}})} \quad (1)$$

where $f(q_i, D)$ is q_i 's term frequency in the document D , $|D|$ is the length of the document D in words, IDF is inverse document frequency, calculated as $\log(1 + (N - n + 0.5)/(n + 0.5))$, where n is the number of documents that q_i appears in, and N is the total number of documents, and avgdl is the average document length in the text collection from which documents are drawn. k_1 and b are free parameters. k_1 represents a term saturation parameter, controlling the equation's sensitivity to incrementally increasing term frequencies, and b is a length normalization parameter, controlling the equation's sensitivity to the length of a phrase. For the experiments reported in this paper, $k_1 = 1.2$ and $b = 0.75$.

To develop and optimize the retrieval process, we worked with an index derived from a single receipt. In each section, we report scores against an index based on the development receipt. The final results report results against the entire set of 65 receipts.

5.1 Baseline

The first experiment used out-of-the-box Lucene. This setup tokenized the index and queries using the *StandardAnalyzer*, which does superficial processing on terms such as lowercasing them.

Using this approach, we achieve an accuracy of 0.623. This basic approach is capable of matching entities with identical scanned representations, but falls short when the scanned representation contains abbreviations. For example, the query "KRO WATER" will consider the entity "Fiji Water" an equally good match as "Kroger Water" even though "KRO" is an abbreviation for "Kroger".

5.2 Wildcards

To help solve the problem of missed abbreviations, we introduced a wildcard technique. First, we construct a dictionary of all the words present in the web entity entries. When constructing a query from the scanned representation, we eliminate all words that are present in the dictionary. For the remaining terms, we rewrite those terms

to match indexed terms that contain all the letters in the raw terms. For example, the term "KRO" will be rewritten as in the query as "K*R*O*" and hence match "Kroger". Note that in Lucene, wildcard matches do not use term frequency or inverse document frequency but count as 1.0 if there is a match between the wildcard and some indexed term, and 0 otherwise. Hence, a query with 3 wildcard terms and two matches would receive a score of exactly 2.0.

Using the wildcard technique, we boost the accuracy to 0.841, a substantial improvement of the baseline of 0.623. The wildcard approach does seem to make some progress towards solving the abbreviation problem, however, there are still cases where this simple approach fails. For example, if an abbreviation spans multiple terms, the wildcard approach will not be able to find a match. For example, the query "P*R*S*L* TOMATOES" will not match the entity "Private Select Tomatoes", since wildcards are only matched against single terms: "P*R*S*L*" matches neither "Private" or "Select".

5.3 Mashed Wildcards

To attempt to solve the multi-word abbreviation problem, we introduced a new field into the indexed documents, called "mashed_terms". This field concatenates all terms into a set of mashed terms, allowing each word in the set to serve as the prefix of the mashed term. For example, the terms "the quick brown" would be rewritten into the `mashed_terms` field as "thequickbrown quickbrown brown". This approach matches multi-word abbreviations. For example "P*R*S*L*" will now match "privateselect". Using the Mashed Wildcards approach, accuracy improves to .884.

5.4 Phrases

Examination of errors using the Mashed Wildcard technique reveals that mashing terms results in some strange false positives. For example, "OCEANS HALO BROTH" matches "Pero Organic Green Beans". This is due to the fact that the non-dictionary term "OCEANS" matches the mashed term "OrganiCgreenBEANS" (caps indicate matching characters). This type error suggests that we should be more selective when generating mashed terms. Ideally, we would like to curate semantically meaningful phrases instead of long, arbitrary sequences of words. To do this, we can analyze the entity labels and identify strongly

collocated bigrams and trigrams, and use concatenated versions of those instead of mashed terms.

An information-theoretically motivated measure for discovering interesting collocations is *Pointwise Mutual Information* (Manning and Schütze, 1999)

$$\begin{aligned} I(x', y') &= \log_2 \frac{P(x', y')}{P(x')P(y')} \\ &= \log_2 \frac{P(x'|y')}{P(x')} \\ &= \log_2 \frac{P(y'|x')}{P(y')} \end{aligned} \quad (2)$$

Using this technique, we generate a list of semantically significant bigrams and trigrams. For example, the top 10 detected bigrams are:

- advanced whitening
- alfaros artesano
- alfresco pasture-raised
- ang chck
- antipasto italiano
- arm hammer
- arrabbiata fra
- aretsano bakery
- athenos crumbled
- atkins indulge

And the top 10 trigrams are

- alfaros artesano bakery
- ang chck pty
- antipasto italiano wildbrine
- arm hammer peroxi
- arrabbiata fra diavolo
- artesano bakery bun
- atkins indulge chocolate
- bagel sesame bagels
- bagels zia italiana
- bf ang chck

Using the ngram technique, we improve accuracy to .899

5.5 Fuzzy Phrases

It turns out that abbreviations are not the only match problem that we need to contend with. Sometimes raw representations are represented as plural, whereas the web representations are represented as singular. For example, the web representation could be “artichoke” with a corresponding raw representation of “ARTICHOKES”.

| Technique | Single Receipt | All Receipts |
|--------------|----------------|--------------|
| Baseline | 0.62 | 0.47 |
| Wildcard | 0.84 | 0.72 |
| Mashed | 0.88 | 0.76 |
| Ngrams | 0.89 | 0.77 |
| Fuzzy Ngrams | 0.93 | 0.79 |

Table 1: Receipt Item Linking Results.

A query of “A*R*T*I*C*H*O*K*E*S*” will not match “artichoke”. To deal with this, we introduce fuzzy matching, which will match terms within a small edit distance of the query term. In the above example, the query “ARTICHOKE” will be rewritten to “A*R*T*I*C*H*O*K*E*S* ARTICHOKES”. The second term will correctly match the “artichoke” entity.

Using fuzzy terms in conjunction with ngrams, accuracy improves to 0.928.

6 Results

We ran each IR method on an index of all 65 receipts. Table 1 shows the accuracies for the single receipt and all receipt cases. The results show improve with the introduction of each new IR technique. Unsurprisingly, the results for all receipts is lower than that of a single receipt. This is primarily to do the expansion of choices introduced by more data. Recall that we are only judging a hit by looking at the top search result. It is likely that in many examples, the true match is still contained in one of the top-k matches, where k is a relatively small number.

7 Future Work

blah

8 Conclusion

blah

References

- Apache Software Foundation. [Lucene](#).
- S. Deerwester, S.T. Dumais, G.W. Furnas, T.K. Landauer, and R.A. Harshman. 1990. Indexing by latent semantic analysis. *Journal of the American Society for Information Science* 41, pages 391–407.
- Eric Melz. 2020. [Understanding receipts github](#).
- Michael R. Glass, Md. Faisal Mahbub Chowdhury, and Alfio Massimiliano Gliozzo. 2017. [Language](#)

independent acquisition of abbreviations. *CoRR*, abs/1709.08074.

Ben Hachey, Will Radford, Joel Notham, Matthew Honnibal, and James R. Curran. 2012. [Artificial intelligence, wikipedia and semi-structured resources evaluating entity linking with wikipedia](#). *Artificial Intelligence*, 194(11):130–150.

Z. Huang, K. Chen, J. He, X. Bai, D. Karatzas, S. Lu, and C. V. Jawahar. 2019. Icdar2019 competition on scanned receipt ocr and information extraction. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 1516–1520.

Guillaume Jacquet, Maud Ehrmann, and Ralf Steinberger. 2014. Clustering of multi-word named entity variants: Multilingual evaluation. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, Reykjavik, Iceland. European Language Resources Association (ELRA).

Yue Liu, Tao Ge, Kusum Mathews, Heng Ji, and Deborah McGuinness. 2015. [Exploiting task-oriented resources to learn word embeddings for clinical abbreviation expansion](#). In *Proceedings of BioNLP 15*, pages 92–97, Beijing, China. Association for Computational Linguistics.

Ismiini Lourentzou, Kabir Manghnani, and Chengxiang Zhai. 2019. [Adapting sequence to sequence models for text normalization in social media](#). *CoRR*, abs/1904.06100.

Christopher D. Manning and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA, USA.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. [Distributed representations of words and phrases and their compositionality](#). In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 3111–3119. Curran Associates, Inc.

S. Robertson. 2009. [The Probabilistic Relevance Framework: BM25 and Beyond](#). *Foundations and Trends® in Information Retrieval*, 3(4):333–389.

Ariel S. Schwartz and Marti A. Hearst. 2003. [A simple algorithm for identifying abbreviation definitions in biomedical text](#). In *Pacific Symposium on Biocomputing*, pages 451–462.

A Appendices

blah

B Supplemental Material

blah blah blah