



Motivation

- Our objective: create a model that learns to **extract revenue information** from richly formatted financial 10-Q filings
- Complicated task because attributes and relations are expressed in a combination of textual, structural, tabular and visual signals

Dataset

Get Data from SEC

- PostgreSQL
- HTML format
- 2,007 reports
- Text spans: +60 million

Structural model



Defining training candidates

Figure 10.1
CONDENSED CONSOLIDATED STATEMENTS OF OPERATIONS (Unaudited)
in thousands, except number of shares which are included in thousands and per share amounts

	2018	2017
Net sales	79,164	79,164
Products	12,570	8,226
Services	66,594	70,938
Total net sales	79,164	79,164

Cost of sales

Products	48,200	50,475
Services	4,861	3,966
Total cost of sales	53,061	54,441
Gross margin	26,103	24,723

• Candidate - (Date, Revenue) pair

• Potential # candidates - 60 mm * 60 mm

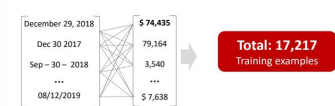
• Total: 3.5e+15

↑

• Too many so we filter

Hard-filtering – limit # of candidates

- Logical, tabular, format, content, linguistic, RegEx rules to limit # of candidates



Weak supervision

- Manual labeling is unfeasible, so we use data programming (Snorkel)
- **Labeling functions (LFs)** evaluate the relation between the mentions of each candidate

Example labeling functions (we have 14 in total)

```
def lf_vertical_align(c):
    (period, revenue) = c
    per = period.context.sentence
    rev = revenue.context.sentence
    if per.col_start == rev.col_start and per.col_end == rev.col_end: return TRUE
    else: return ABSTAIN

def lf_early_in_table(c):
    (period, revenue) = c
    per = period.context.sentence
    rev = revenue.context.sentence
    if per.row_start == 0 and rev.row_start == 12: return TRUE
    else: return ABSTAIN
```

GAN for label generation

- Apply LFs to unlabeled data, resulting in a label matrix Λ . Then encode generative model $p_w(\Lambda, Y)$ using three factor types: labeling propensity, accuracy, and pairwise correlations of LFs

$$\phi_{i,j}^{\text{Lab}}(\Lambda, Y) = \mathbb{1}\{\Lambda_{i,j} \neq \emptyset\} \quad p_w(\Lambda, Y) = Z_w^{-1} \exp\left(\sum_{i=1}^m w_i^T \phi_i(\Lambda, Y)\right)$$

$$\phi_{i,j}^{\text{Acc}}(\Lambda, Y) = \mathbb{1}\{\Lambda_{i,j} = y_i\}$$

$$\phi_{i,j,k}^{\text{Corr}}(\Lambda, Y) = \mathbb{1}\{\Lambda_{i,j} = \Lambda_{i,k}\} \quad \hat{w} = \arg \min_w -\log \sum_Y p_w(\Lambda, Y)$$

Creating weak labels

GAN marginal probabilities

Custom marginal probabilities

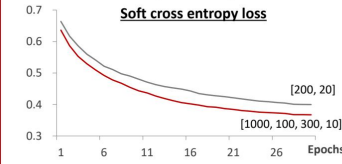
Will try both training marginals to choose the best performing one

Feature matrix

- **107,000 features:** One-hot vectors of surrounding words, style features (font, size, bold, caps, etc.), NLP structure, length, lemma sequences, row and column headers, table characteristics, page location, object hierarchy, html tags, and others

Results

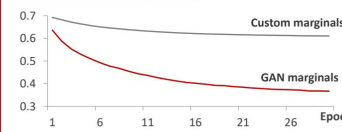
- Chosen model: Shallow [1000, 100, 300, 10] neural network with a soft cross entropy loss function using Adam optimizer



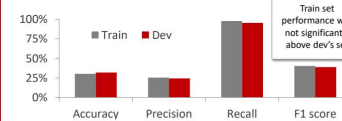
	larger NN		Smaller NN	
	Development	Test	Development	Test
Accuracy	33.8%	35.7%	27.2%	31.0%
Precision	24.8%	28.0%	23.2%	27.2%
Recall	95.3%	87.8%	96.3%	92.8%
F1 score	39.4%	42.4%	F1 score	37.4%

[Performance vs true labels - obtained from Capital IQ]

Choice of marginals



Overfitting?



Conclusions

- Through a combination of hard filters and weak supervision, our model was able to pinpoint a handful of revenue + date candidates out of billions of potential pairs
- F1 performance was slightly below 40%, with a baseline performance of 0% for the broad dataset and 20% for the filtered one
- Although performance was decent, more than half of the model's performance was in the hard-filtering portion (20%).

Expansions

- Use fewer hard filters and train model with more candidates (over 1mm ideally)
- To avoid computation unfeasibility, trim features from 107,000 to less than 10,000
- Overall, creating a structural model is computationally inefficient for a task like this. Table and page extraction based on heuristics paired with a machine learning model would be more efficient

References

[1] Sen Wu, Luke Hsiao, Xiao Cheng, Braden Hancock, Theodoros Rekatsinas, Philip Levis, and Christopher Ré. 2018. Fonduer: Knowledge Base Construction from Richly Formatted Data. In Proceedings of 2018 International Conference on Management of Data, Houston, TX, USA, June 10–15, 2018 (SIGMOD'18), 16 pages.

[2] Alexander Ratner, S. H. (November 2017). Snorkel: Rapid Training Data Creation with Weak Supervision. Stanford University.

Team members

- Andrea Aguirre, andreaar@stanford.edu
- Roberto Seminario, rseminar@stanford.edu