

TEAM MEMBERS: Yury A. Kuleshov, Pum Khai, Leying Zhou (Data Mine); Girija Bhamidipati, Thrishna Bhandari, Haleigh Gronwold, Pawandeep Kaur, Rthvik Raviprakash, Kayla Seeley (Capstone) **MENTORS**: Dr. Davide Gerbaudo, Dr. Sridhar Ramaswamy; Kevin R. Sellers (CAT); Dr. Jean F. Honorio Carrillo (Capstone) **TAs**: Jayanth Krishnamurthy (Data Mine), Tanmaya Udupa (Capstone)

INTRODUCTION

About Caterpillar & CAT Digital

- Cat Digital is a digital and technology business unit within Caterpillar Inc.
- \circ We're focused on using data, technology, advanced analytics and AI capabilities to help our customers build a better world.
- Cat Digital embodies the bitsand-bytes coming from pieces of machinery, customers and dealers, and is charged with ensuring data flows smoothly and safely through our platform and to applications.

Caterpillar Inc. is the world's leading manufacturer of construction and mining equipment, off-highway diesel and natural gas engines, industrial gas turbines and diesel-electric locomotives.

For nearly 100 years, we've been helping customers build a better, more sustainable world and are committed and contributing to a reduced-carbon future. Our innovative products and services, backed by our global dealer network, provide exceptional value that helps customers succeed.





PRODUCT LINE





Construction Equipment

Motivation & Problem

Currently, Caterpillar is using a clustering model based primarily on industry- and company-based rules to categorize invoices and parts. This strategy is effective for bootstrapping and is easy to understand and explain to their business partners. However, these rules are not necessarily backed by any scientific processes and trends/rules can change often. Hence, Caterpillar needs a more robust solution with proper validation to classify and categorize invoices from parts sales to customers.

Objective

The goal of this project is to explore a *novel document clustering algorithm* to be used with invoices. The new algorithm should leverage a similarity metric built upon the numerical and categorical features associated with the line items listed on the invoice documents.



Research Design Ideas

Since we want to better cluster the invoices, we are trying to determine *distance metrics* for this type of data specifically. \rightarrow The similarity metric should work with the "invoice" data type, namely a mix of numerical, categorical, and textual data. \rightarrow "Similarity metric" feature that given two invoices it can determine how similar or dissimilar they are.

Key Approaches

- Similarity metrics
- Tf-IDF clustering
- Natural Language Processing



.... Stage Looking for pattern via

correlation

Stage 4 Document for deliverable

Output









Equipment



Gas Engines





Locomotives









• Data Cleaning Principles

- Reduce the noise from negative price and quantity • Clarify electric machine component numbers Improve null values
- Delete non-Caterpillar parts in the dataset • Create dummy variables
- Data Construction
- Group by invoice_ids
- Group by major_class
- Carry out stratified sampling
- Data Format Create Sales_method as dummy variables

Models

Improvement of clustering algorithm • **Cosine similarity** is more useful than Euclidean

- distance.
- Principal Component Analysis did not seem to improve the models much in reducing the number of dimensions.
- To evaluate, we have been looking at the number of data points in each cluster to see how evenly they are spread and also looking at inter-cluster distance.

- Include more advanced clustering ideas such as Gower Distance, and k-medoids, etc.
- Further, to improve upon the algorithms/model with a computationally efficient metric in a way that millions of invoices can be processed by computers compared to traditional document clustering algorithms.

KEY FINDINGS & SOLUTIONS

• Data Description & Data Attributes

nns with a HASH() have been obfuscated for data confidentiality

- primary key to identify an invoice code to identify the dealership where the part invoice was issued
- date at which the invoice was issued how the parts were sold: OTC= over the counter, WO= work order with the repair performe B2B = business-to-business internet purchase (SERIALNUMBER) -- serial number (more precisely its prefix) of the machine(s) being repaired with
 - -- sales model of the machine(s) being repaired. code to identify a replacement part
- : name as given by the engineer (can be very specific/generic, have typos, abbr. etc.) SOURCE_OF_SUPPLY) -- whether the part is supplied by CAT (majority) or by another vendor (PART PRICING CODE) -- category of parts used to determine its price
- (MAJOR CLASS) -- coarse category of parts COMMERCIAL GROUP TYPE) -- intermediate category of parts
- AFTERMARKET OFFER) -- fine category of parts
- intermediate component where this part was installed ACHINE COMPONENT L3 -- fine component where this part was installed

• DBSCAN model (Figure 2) • K-Means with 3 clusters (Figure 3) • Hierarchical clustering (Figure 4)

CONCLUSIONS

FUTURE GOALS

1	-0.003	0.00077	0.0026	0.00017	-0.0017	-0.0015	-0.00055	0.00064	-2.6e-05	8e-05	-0.00061
0.003	1	0.021	0.032	-0.001	0.17	0.0067	0.009	-0.012	0.0052	-0.0065	-0.00037
- 0.00077	0.021	1	0.21	-0.00015	-0.049	0.014	0.016	-0.00018	0.029	0.0011	-0.00032
- 0.0026	0.032	0.21	1	0.00032	-0.026	0.0058	0.034	0.0037	0.033	0.0033	0.00022
- 0.00017	-0.001	-0.00015	0.00032	1	-0.00025	-0.015	0.011	0.027	0.011	0.00069	-0.00046
0.0017	0.17	-0.049	-0.026	-0.00025	1	-0.04	-0.042	-0.042	-0.061	-0.014	-0.00092
0.0015	0.0067	0.014	0.0058	-0.015	-0.04	1	-0.017	-0.016	0.1	0.002	-0.00012
0.00055	0.009	0.016	0.034	0.011	-0.042	-0.017	1	0.15	-0.16	0.0085	-0.0025
- 0.00064	-0.012	-0.00018	0.0037	0.027	-0.042	-0.016	0.15	1	0.16	-0.014	0.002
2.6e-05	0.0052	0.029	0.033	0.011	-0.061	0.1	-0.16	0.16	1	0.04	0.0022
- 8e-05	-0.0065	0.0011	0.0033	0.00069	-0.014	0.002	0.0085	-0.014	0.04	1	2.7e-06
-0.00061	-0.00037	-0.00032	0.00022	-0.00046	-0.00092	0.00012	-0.0025	0.002	0.0022	2.7e-06	1











DBSCAN

Figure 1. Heatmap (from 10,000 rows)

- No significant correlations between attributes most numerical attributes are hashed others are categorical • Sales model &
- sales_method are two specific attributes of relevance for data mining and business purposes since they provide insights into how big/small the repair was and the method in which it was accomplished



Figure 3. K-Means with 3 clusters

• Chosen based off an elbow curve

Figure 4. Hierarchical clustering Dendrogram produced using 6 clusters