
SALT: Sales Autocompletion Linked Business Tables Dataset

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Foundation models, particularly those that incorporate Transformer architectures,
2 have demonstrated exceptional performance in domains such as natural language
3 processing and image processing. Adapting these models to structured data, like
4 tables, however, introduces significant challenges. These difficulties are even more
5 pronounced when addressing multi-table data linked via foreign key, which is
6 prevalent in the enterprise realm and crucial for empowering business use-cases.
7 Despite its substantial impact, the research focusing on such linked business tables
8 within enterprise settings remains a significantly important yet underexplored
9 domain. To address this, we introduce a curated dataset sourced from an Enterprise
10 Resource Planning (ERP) system, featuring extensive linked tables. This dataset is
11 specifically designed to support research endeavors in table representation learning.
12 By providing access to authentic enterprise data, our goal is to potentially enhance
13 the effectiveness and applicability of models for real-world business contexts.

14 1 Introduction

15 Deep learning has made substantial strides in areas like text understanding, language translation,
16 image classification, and object detection. These advancements are largely driven by foundational
17 models trained on diverse datasets and self-supervised training techniques, especially those that
18 incorporate Transformer architectures. However, using these models on structured, tabular data,
19 essential for enterprise business operations, poses unique challenges. These challenges become more
20 pronounced with multi-table configurations consisting of large tables interconnected by foreign keys
21 and comprising extensive business datasets, a setup to which we refer to as *linked business tables*.
22 Such setups are common in real-world business scenarios.

23 The challenges in applying foundational models to linked business data are primarily twofold: al-
24 gorithmic and data-related. Algorithmically, a significant challenge is adapting models that were
25 originally designed for unstructured internet data to handle structured data effectively - see (Grinsztajn
26 et al., 2022) for a comprehensive discussion. This process requires a sophisticated integration of
27 structural knowledge and the unique characteristics of linked business data, which is inherently more
28 complex and interconnected than straightforward internet-scraped table data.

29 One major limitation in the current landscape is the absence of realistic, enterprise-linked multi-table
30 datasets at scale. Existing table datasets often originate from HTML pages and do not accurately
31 represent the complexity and dynamics of expansive database tables used in active enterprise sys-
32 tems (Bodensohn et al., 2024). Moreover, obtaining large, clean, and high-quality datasets for
33 structured tabular applications presents difficulties (Hulsebos et al., 2023; Van Breugel and Van
34 Der Schaar, 2024), particularly in enterprise settings where data privacy, confidentiality, and commer-
35 cial interests restrict data access. This lack of suitable public datasets leads to significant domain
36 adaptation challenges and shifts in data distribution, which pose difficulties for many existing mod-
37 els (Fey et al., 2024).

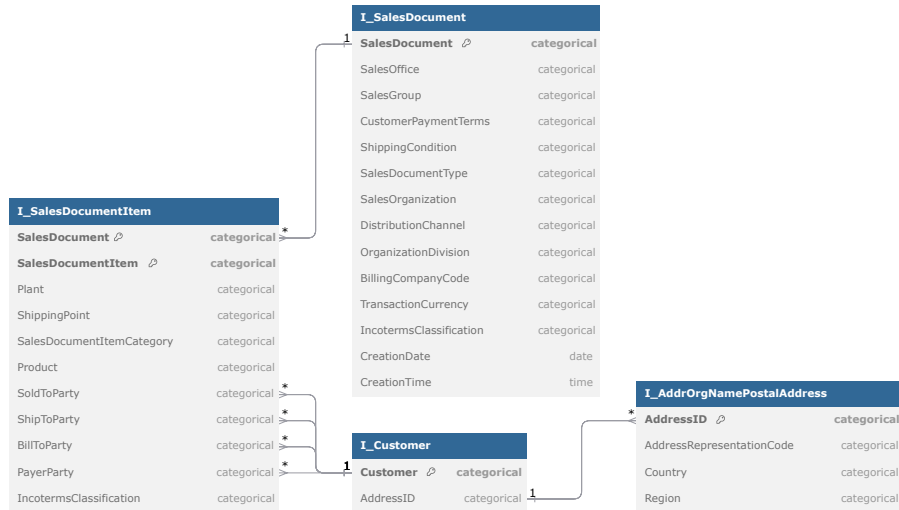


Figure 1: **Dataset Schemas**. Schemas for the four tables constituting the SALT dataset. Primary keys are highlighted in **bold** letters and with a key symbol. Foreign keys interconnecting tables.

38 To tackle these issues, we have curated the Sales Autocompletion Linked Business Tables (SALT)
 39 dataset, sourced from an Enterprise Resource Planning (ERP) system. ERP systems are comprehen-
 40 sive, multifunctional platforms essential for managing all core business operations including finance,
 41 human resources, production, and supply chains. As the backbone of organizational data manage-
 42 ment, ERP systems provide an excellent foundation for developing and evaluating data models that
 43 accurately reflect complex, real-world enterprise environments. The SALT dataset, which includes
 44 interconnected relational tables with a focus on sales, encompasses several million entries across
 45 various enterprise sales operations.

46 By sharing the SALT dataset with the research community, we aim to stimulate advancements in
 47 table representation learning and refine algorithm development to enhance applicability and perfor-
 48 mance in real-world settings. This initiative is crucial for evolving deep learning models that not
 49 only understand but also effectively function within the complexities of large-scale enterprise data
 50 landscapes, thereby promoting the development of enterprise-specific machine learning applications.

51 **Related Work:** The majority of existing table datasets originate from scraping the Web, notably
 52 extracted from HTML pages or CSV files from GitHub, which inadequately capture the complexity
 53 and dynamics typical of large database tables that are employed in operational enterprise systems.
 54 WebTables (Cafarella et al., 2008) corpus includes a massive collection of 233 million tables, sourced
 55 from HTML pages via the Common Crawl project. While WebTables offers an extensive quantity of
 56 tables, its diversity is constrained because it solely comprises HTML tables from web pages. TURL
 57 Deng et al. (2020) provides a cleaner corpus of 580 thousand tables extracted from Wikipedia. In
 58 contrast, GitTables (Hulsebos et al., 2023) contains over 10 million tables extracted from "comma-
 59 separated value" files (CSVs) found on GitHub. Tables from GitTables generally exhibit structural
 60 differences compared to those from WebTables, making GitTables an essential corpus despite its
 61 focus on a single file type. TabLib (Eggert et al., 2023) comprises 627 million tables across various
 62 file formats and totaling 69 TiB, sourced from GitHub and Common Crawl. Notably, it comprises
 63 exceptionally large tables of several million rows and columns. LakeBench (Deng et al., 2024) is
 64 a collection of benchmarks to resemble enterprise data lakes, containing tables from a variety of
 65 sources such as open government data for the purpose of unionability, joinability, and subset tasks.

66 2 SALT Dataset

67 **Background:** The SALT dataset is specifically curated to mirror customer interactions within an
 68 Enterprise Resource Planning (ERP) system and is designed to train models that assist users by
 69 predicting fields typically missing in sales orders. This dataset is crucial to the sales and distribution
 70 process, especially for creating the "Sales Order Document." Each of these documents records a
 71 single transaction that includes various items intended for specific customers, marking a distinct

Table 1: **Breakdown of atomic data composition of datasets:** Datasets SALT, *GitTables* (Hulsebos et al., 2023), *WebTables* (Lehmberg et al., 2016) and *TabLib* (Eggert et al., 2023)

Atomic data type	SALT	GitTables	WebTables	TabLib
Numeric	38.7%	57.9%	51.4%	33.6%
String	58.1%	41.6%	47.4%	61.8%
Other	3.2%	0.5%	1.2%	4.6%

72 phase in the sales cycle.

73 Structured around four principal tables—sales documents, sales document items, customers, and
 74 addresses—the dataset consolidates data from a single enterprise that underwent anonymization. The
 75 sales documents table logs vital details such as sales office, sales group, payment conditions, and
 76 shipping arrangements, limiting its entries to those specifically categorized as sales orders. The sales
 77 document items table captures detailed information for each line item in these documents, including
 78 the product sold, shipping point, and the parties involved in the transaction. Concurrently, the
 79 customer table holds comprehensive master data about customers, further elaborated in the addresses
 80 table with specifics like country and region. The input variables in the dataset include a mix of fields
 81 typically populated by users during the creation of a sales order, augmented by master data fields
 82 like material number and customer details. The target variables are not always maintained; they
 83 are optional and may not be filled out for certain transactions depending on particular scenarios or
 84 requirements. This intricate structure of SALT not only enhances model training for missing field
 85 predictions but also effectively replicates complex ERP interactions within sales order automation.

86 **Task:** In the dataset, 21 fields are categorized as input variables, serving as features for predictive
 87 modeling applications, while 8 fields are designated as target variables, intended for prediction
 88 based on the input data analysis. The predictive model, which will be trained using this dataset,
 89 is specifically tasked with performing multiclass classification on seven critical variables. These
 90 variables are essential for ensuring the seamless execution of sales orders:

- 91 • `I_SalesDocument.SalesOffice` - Sales activities for specific products and regions
- 92 • `I_SalesDocument.SalesGroup` - Sales representatives managing responsibilities
- 93 • `I_SalesDocument.CustomerPaymentTerms` - Payment conditions, i.e., deadlines and
 94 early payment discounts
- 95 • `I_SalesDocument.ShippingCondition` - Logistics terms
- 96 • `I_SalesDocumentItem.ShippingPoint` - Dispatch location
- 97 • `I_SalesDocumentItem.Plant` - Production/ storage facility, critical for inventory control
- 98 • `I_SalesDocument.IncotermsClassification` and
 99 `I_SalesDocumentItem.IncotermsClassification` - International commercial terms,
 100 outline transaction responsibilities like shipping and insurance¹

101 **Structure:** The dataset is structured into four primary tables encompassing a total
 102 of 573,810 sales orders (`I_SalesDocument`), which include 2,706,491 sales order items
 103 (`I_SalesDocumentItem`) associated with 136,317 unique business partners (`I_Customer`) and
 104 1,650,641 (`I_AddrOrgNamePostalAddress`) addresses - see Fig. 1 for the table schemas. The table
 105 fields are filtered to include only the data relevant to the specific use case described above. After
 106 filtering, the tables are merged to form a single flat dataset containing 2,706,491 rows, such that
 107 each row in the dataset represents a single sales order item (for more details see Appendix Sec. A.3).
 108 The entries cover transactions conducted between January 1, 2018, and December 31, 2020. To
 109 assess the dataset’s predictive modeling utility, data was divided with temporal splits, with validation
 110 segments starting from February 1, 2020, and test segments from July 1, 2020. For on analysis on the
 111 distribution values see Tab. 1.

112 **Data Insights:** The dataset employed in this study is derived from authentic industry data captured by
 113 an Enterprise Resource Planning (ERP) system, documenting sales orders. This dataset has undergone
 114 minimal pre-processing primarily aimed at addressing privacy concerns. Several challenges arise
 115 from the nature and quality of the dataset, which need careful consideration:

- 116 • **Diversity:** There is a substantial diversity in certain data fields due to the wide range of
 117 unique values they contain. For instance, the field `I_SalesDocumentItem.ShipToParty`

¹This field can be defined independently on item and header level, which is why both are included.

118 includes 21,997 distinct customer IDs, while `I_SalesDocumentItem.Product` comprises
 119 209,823 unique product identifiers.

- 120 • **Class imbalance:** The dataset demonstrates a pronounced class imbalance. The distribution
 121 of sales offices across sales orders is highly skewed; the most frequently occurring sales office
 122 is associated with 75% of the orders, and the two most common sales offices collectively
 123 account for 98% of the data. Despite this, there are 33 distinct sales offices represented in
 124 more than one order, suggesting a long-tail distribution.
- 125 • **Noise:** A considerable amount of input noise is evident within the dataset. Since data entry
 126 is frequently manual, discrepancies may arise as different employees might handle identical
 127 business scenarios differently or make inadvertent errors. Moreover, certain fields may be
 128 occasionally left blank, potentially leading to gaps in the data.
- 129 • **Data drift:** Technically, the dataset is prone to data drift, a phenomenon where the cate-
 130 gorizations such as sales groups within the ERP system evolve over time. This drift may
 131 particularly impact analyses involving temporal splits of the data, as category definitions
 132 may shift across the time periods considered. Notably, the target categories are not subject
 133 to such data drift.

Table 2: **Classification performance of baseline models:** Evaluation of baseline models on the eight different tasks on SALT. **Top:** Simple baselines **Middle:** Gradient-boosted decision tree models **Bottom:** Deep learning methods. **Performance metric:** Mean Reciprocal Rank.

Performance Baseline - MRR (\uparrow)					
Method \ Target Variable	Plant	Shipping Point	Item Incoterm Cls.	Header Incoterm Cls.	
Random Classifier	0.32	0.22	0.33	0.33	
Majority Class Baseline	0.51	0.41	0.49	0.49	
XGBoost (Chen and Guestrin, 2016)	0.99	0.96	0.75	0.75	
LightGBM (Ke et al., 2017)	0.99	0.86	0.82	0.82	
CatBoost (Prokhorenkova et al., 2018)	0.99	0.99	0.82	0.82	
CARTe (Kim et al., 2024)	0.97	0.96	0.77	0.77	
AutoGluon (Erickson et al., 2020)	0.99	0.98	0.79	0.79	
(continued)	Sales Office	Sales Group	Pay. Terms	Ship. Condition	Avg.
Random Classifier	0.97	0.01	0.13	0.16	0.31
Majority Class Baseline	0.99	0.03	0.23	0.29	0.43
XGBoost (Chen and Guestrin, 2016)	0.98	0.62	0.75	0.71	0.81
LightGBM (Ke et al., 2017)	0.99	0.21	0.83	0.80	0.79
CatBoost (Prokhorenkova et al., 2018)	0.99	0.17	0.56	0.74	0.76
CARTe (Kim et al., 2024)	0.99	0.40	0.70	0.77	0.79
AutoGluon (Erickson et al., 2020)	0.99	0.66	0.84	0.82	0.86

134 3 Experiments & Results

135 We evaluate the SALT dataset using several baselines for tabular data on the joined table. The only
 136 preprocessing applied is filling in missing values with either a constant value (for the categorical
 137 features) or the mean value (for numerical features). The fields related to creation date and time were
 138 only used to split the data and then discarded. The validation set was used for early stopping and
 139 no hyperparameter tuning was performed. See Tab. 2 for the detailed breakdown of performance
 140 evaluation tasks of each task. As can be seen the AutoGluon Erickson et al. (2020) shows best perfor-
 141 mance on SALT with a significant margin of (+0.05 p .) w.r.t. the next best approach XGBoost (Chen
 142 and Guestrin, 2016). The analysis reveals several noteworthy insights: **i)** Certain target variables
 143 demonstrate substantial predictability, achieving prediction scores near 0.99, indicating a high degree
 144 of accuracy. **ii)** The dataset exhibits significant class imbalance, which is particularly evident from
 145 the performance of the majority class baseline. This imbalance is most pronounced when predicting
 146 variables such as the Sales Office. **iii)** The predictive performance of the model is adversely affected
 147 when tasked with predicting fields like the Sales Group, which suffers from high cardinality issues,
 148 complicating accurate classification.

149 4 Conclusion

150 We introduce a novel dataset focused on linked business data, demonstrating the characteristics of
 151 data within actual enterprise systems. We further assessed the performance of current tabular models
 152 against tree-based and cutting-edge models. The empirical data reveal that most tabular models
 153 effectively manage the prediction tasks in SALT. To augment the dataset’s complexity and utility
 154 in future work, we plan to include additional tables from a broader range of scenarios, data from
 155 multiple companies, and enhance the semantic richness of the dataset to present greater challenges.

References

- 156
157 Jan-Micha Bodensohn, Ulf Brackmann, Liane Vogel, Matthias Urban, Anupam Sanghi, and Carsten
158 Binnig. 2024. Llms for data engineering on enterprise data. In *Joint proceedings of workshops at*
159 *the 50th International Conference on Very Large Data Bases (VLDB 2024)*, Guangzhou, China,
160 *August 26 - August 30, 2023, VLDBW 2024*, Tabular Data Analysis Workshop Proceedings.
- 161 Michael J. Cafarella, Alon Halevy, Daisy Zhe Wang, Eugene Wu, and Yang Zhang. 2008. Webtables:
162 exploring the power of tables on the web. *Proc. VLDB Endow.*, 1(1):538–549.
- 163 Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A scalable tree boosting system. In *Proceedings*
164 *of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,
165 *KDD '16*, pages 785–794, New York, NY, USA. ACM.
- 166 Xiang Deng, Huan Sun, Alyssa Lees, You Wu, and Cong Yu. 2020. Turl: table understanding through
167 representation learning. *Proc. VLDB Endow.*, 14(3):307–319.
- 168 Yuhao Deng, Chengliang Chai, Lei Cao, Qin Yuan, Siyuan Chen, Yanrui Yu, Zhaoze Sun, Junyi
169 Wang, Jiajun Li, Ziqi Cao, Kaisen Jin, Chi Zhang, Yuqing Jiang, Yuanfang Zhang, Yuping Wang,
170 Ye Yuan, Guoren Wang, and Nan Tang. 2024. Lakebench: A benchmark for discovering joinable
171 and unionable tables in data lakes. *Proc. VLDB Endow.*, 17(8):1925–1938.
- 172 Gus Eggert, Kevin Huo, Mike Biven, and Justin Waugh. 2023. Tablib: A dataset of 627m tables with
173 context.
- 174 Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander
175 Smola. 2020. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint*
176 *arXiv:2003.06505*.
- 177 Matthias Fey, Weihua Hu, Kexin Huang, Jan Eric Lenssen, Rishabh Ranjan, Joshua Robinson,
178 Rex Ying, Jiaxuan You, and Jure Leskovec. 2024. Position: Relational deep learning - graph
179 representation learning on relational databases. In *Proceedings of the 41st International Conference*
180 *on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 13592–
181 13607. PMLR.
- 182 Leo Grinsztajn, Edouard Oyallon, and Gael Varoquaux. 2022. Why do tree-based models still
183 outperform deep learning on typical tabular data? In *Thirty-sixth Conference on Neural Information*
184 *Processing Systems Datasets and Benchmarks Track*.
- 185 Madelon Hulsebos, Çagatay Demiralp, and Paul Groth. 2023. Gittables: A large-scale corpus of
186 relational tables. *Proceedings of the ACM on Management of Data*, 1(1):1–17.
- 187 Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan
188 Liu. 2017. Lightgbm: a highly efficient gradient boosting decision tree. In *Proceedings of the 31st*
189 *International Conference on Neural Information Processing Systems*, NIPS'17, page 3149–3157,
190 Red Hook, NY, USA. Curran Associates Inc.
- 191 Myung Jun Kim, Leo Grinsztajn, and Gael Varoquaux. 2024. Carte: Pretraining and transfer for
192 tabular learning. In *Forty-first International Conference on Machine Learning*.
- 193 Oliver Lehmberg, Dominique Ritze, Robert Meusel, and Christian Bizer. 2016. A large public corpus
194 of web tables containing time and context metadata. In *Proceedings of the 25th International*
195 *Conference on World Wide Web, WWW 2016, Montreal, Canada, April 11-15, 2016, Companion*
196 *Volume*, pages 75–76. ACM.
- 197 Liudmila Ostroumova Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush,
198 and Andrey Gulin. 2018. Catboost: unbiased boosting with categorical features. In *NeurIPS*, pages
199 6639–6649.
- 200 Boris Van Breugel and Mihaela Van Der Schaar. 2024. Position: Why tabular foundation models
201 should be a research priority. In *Proceedings of the 41st International Conference on Machine*
202 *Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 48976–48993. PMLR.

203 **A Appendix**

204 **A.1 Table Detailed Information**

205 This section describes the schemas of the four tables that constitute the SALT dataset. Figure 1
 206 shows the schema of the tables. See Tab. 5 for an detailed overview of the *SalesDocumentItem*, Tab. 3
 207 for *SalesOrder* details, Tab. 6 for customer data detail and Tab. 4 for the customer address details.
 208 The column *Is Target Field* indicates the target fields which should be predicted based on the other
 209 fields in the tables, since they are populated in a later stage of the sales order creation.

210 **A.2 Additional Data Statistics**

211 This section provides additional statistics on data composition. Table 8 provides statistics on the
 212 target fields. Table 8 provides a detailed statistics of the table fields.

Table 3: **SalesDocument in Detail**

Field Name	Data Type	Description	Is Target Field
SalesDocument	Categorical (integer)	ID of the sales document	
SalesDocumentType	Categorical (string)	Type of sales document	
SalesOrganization	Categorical (string)	ID of the sales organization	
DistributionChannel	Categorical (string)	ID of the distribution channel	
OrganizationDivision	Categorical (string)	ID of the organization division	
BillingCompanyCode	Categorical (string)	Company code to be billed	
TransactionCurrency	Categorical (string)	Currency code (EUR, USD, ...)	
CreationDate	Date (string)	Date of sales document creation	
CreationTime	Time (string)	Time of day of sales document creation	
SalesOffice	Categorical (string)	ID of the sales office	✓
SalesGroup	Categorical (string)	ID of the sales group	✓
CustomerPaymentTerms	Categorical (string)	ID of the payment terms	✓
ShippingCondition	Categorical (string)	ID of the shipping condition	✓
IncotermsClassification	Categorical (string)	ID of the incoterms	✓

Table 4: **AddrOrgNamePostalAddress in Detail**

Field Name	Data Type	Description	Is Target Field
AddressID	Categorical (integer)	ID of the address	
AddressRepresentationCode	Categorical (integer)	System internal address code	
Country	Categorical (string)	Country name	
Region	Categorical (string)	Region name	

Table 5: **SalesDocumentItem in Detail**

Field Name	Data Type	Description	Is Target Field
SalesDocument	Categorical (integer)	ID of the sales document	
SalesDocumentItem	Categorical (integer)	ID of the sales document item	
SalesDocumentItemCategory	Categorical (string)	ID of the item category	
Product	Categorical (integer)	ID of the product sold	
SoldToParty	Categorical (integer)	ID of the customer sold to	
ShipToParty	Categorical (integer)	ID of the customer shipped to	
BillToParty	Categorical (integer)	ID of the customer billed to	
PayerParty	Categorical (integer)	ID of the payer	
Plant	Categorical (string)	ID of the plant	✓
ShippingPoint	Categorical (string)	ID of the shipping point	✓
IncotermsClassification	Categorical (string)	ID of the incoterms	✓

213 **A.3 Data Extraction and Processing**

214 The extracted tables are filtered to contain only the data relevant to this use case.

Table 6: **Customer in Detail**

Field Name	Data Type	Description	Is Target Field
Customer	Categorical (integer)	ID of the customer	
AddressID	Categorical (integer)	ID of customer's address	

Table 7: **Target field statistics**

Target field	Unique values (#)	Missing values (%)	Normalized entropy
I_SalesDocument.CustomerPaymentTerms	157	0.96	0.52
I_SalesDocument.IncotermsClassification	13	0.58	0.56
I_SalesDocument.SalesGroup	636	0.77	0.85
I_SalesDocument.SalesOffice	37	0.01	0.18
I_SalesDocument.ShippingCondition	55	0.01	0.55
I_SalesDocumentItem.IncotermsClassification	13	0.21	0.54
I_SalesDocumentItem.Plant	41	0.02	0.49
I_SalesDocumentItem.ShippingPoint	99	0.03	0.51

Table 8: **Field statistics:** Number of unique values, percentage of missing values, normalized entropy.

Table	Field	Unique values (#)	Missing values (%)	Normalized entropy
I_SalesDocument	BillingCompanyCode	33	0.0	0.61
I_SalesDocument	CreationDate	1081	0.0	0.96
I_SalesDocument	CreationTime	63247	0.0	0.97
I_SalesDocument	CustomerPaymentTerms	158	0.96	0.52
I_SalesDocument	DistributionChannel	4	0.0	0.08
I_SalesDocument	IncotermsClassification	14	0.58	0.56
I_SalesDocument	OrganizationDivision	1	0.0	0
I_SalesDocument	SalesDocument	573810	0.0	1.0
I_SalesDocument	SalesDocumentType	16	0.0	0.37
I_SalesDocument	SalesGroup	637	0.77	0.85
I_SalesDocument	SalesOffice	38	0.01	0.18
I_SalesDocument	SalesOrganization	36	0.0	0.6
I_SalesDocument	ShippingCondition	56	0.01	0.55
I_SalesDocument	TransactionCurrency	28	0.0	0.29
I_SalesDocumentItem	BillToParty	19483	0.0	0.78
I_SalesDocumentItem	IncotermsClassification	14	0.21	0.54
I_SalesDocumentItem	PayerParty	19441	0.0	0.78
I_SalesDocumentItem	Plant	42	0.02	0.49
I_SalesDocumentItem	Product	209823	0.0	0.75
I_SalesDocumentItem	SalesDocument	573699	0.0	0.92
I_SalesDocumentItem	SalesDocumentItem	948	0.0	0.57
I_SalesDocumentItem	SalesDocumentItemCategory	24	0.1	0.4
I_SalesDocumentItem	ShipToParty	21998	0.0	0.76
I_SalesDocumentItem	ShippingPoint	100	0.03	0.51
I_SalesDocumentItem	SoldToParty	19454	0.0	0.78
I_Customer	AddressID	136287	0.0	1.0
I_Customer	Customer	136317	0.0	1.0
I_Address	AddressID	1625958	0.0	1.0
I_Address	AddressRepresentationCode	1	100.0	0
I_Address	Country	244	13.17	0.57
I_Address	Region	615	80.54	0.74

- 215 • I_SalesDocument and I_SalesDocumentItem are filtered to contain only documents of
- 216 category sales order
- 217 • Only orders which have been fully processed are included, that is, sales orders which have
- 218 gone through the entire business process. The goal is to ensure that we are working with the
- 219 data in its most complete sense and no further changes would be expected.
- 220 • I_Customer and I_AddrOrgNamePostalAddress are filtered to contain only the business
- 221 partners which appear in the aforementioned sales orders and their respective addresses
- 222 • Additionally, the table fields are also filtered to include only those relevant to the use case,
- 223 which were listed in the previous section

224 After extraction, the tables were joined together to create a flat structure, as described in List. 1, such
 225 that each row in the final table represents one sales order item. Note that, due to this flat structure, the

226 prediction targets which are defined on the sales order level will be repeated across multiple items /
 227 rows. This choice was made to allow for the possibility of training a single model to predict all target
 228 fields.

Listing 1: **Table join SQL Query:** SQL query used to join the 4 tables together

```

1  SELECT
2      SalesDocumentItem . SalesDocument ,
3      SalesDocumentItem . SalesDocumentItem ,
4      SalesDocument . SalesOffice ,
5      SalesDocument . SalesGroup ,
6      SalesDocument . CustomerPaymentTerms ,
7      SalesDocument . ShippingCondition ,
8      SalesDocumentItem . Plant ,
9      SalesDocumentItem . ShippingPoint ,
10     SalesDocument . SalesDocumentType ,
11     SalesDocument . SalesOrganization ,
12     SalesDocument . DistributionChannel ,
13     SalesDocument . OrganizationDivision ,
14     SalesDocument . BillingCompanyCode ,
15     SalesDocument . TransactionCurrency ,
16     SalesDocumentItem . SalesDocumentItemCategory ,
17     SalesDocumentItem . Product ,
18     SalesDocumentItem . SoldToParty ,
19     SoldToPartyAddress . Country as SoldToPartyCountry ,
20     SoldToPartyAddress . Region as SoldToPartyRegion ,
21     SalesDocumentItem . ShipToParty ,
22     ShipToPartyAddress . Country as ShipToCountry ,
23     ShipToPartyAddress . Region as ShipToPartyRegion ,
24     SalesDocumentItem . BillToParty ,
25     BillToPartyAddress . Country as BillToPartyCountry ,
26     BillToPartyAddress . Region as BillToPartyRegion ,
27     SalesDocumentItem . PayerParty ,
28     PayerPartyAddress . Country as PayerCountry ,
29     PayerPartyAddress . Region as PayerRegion ,
30     SalesDocument . CreationDate ,
31     SalesDocument . CreationTime
32  FROM   I_SalesDocumentItem AS SalesDocumentItem
33         INNER JOIN I_SalesDocument AS SalesDocument
34             ON SalesDocument . SalesDocument = SalesDocumentItem . SalesDocument
35         LEFT JOIN I_Customer AS SoldToPartyTable
36             ON SalesDocumentItem . SoldToParty = SoldToPartyTable . Customer
37         LEFT JOIN I_Customer AS ShipToPartyTable
38             ON SalesDocumentItem . ShipToParty = ShipToPartyTable . Customer
39         LEFT JOIN I_Customer AS BillToPartyTable
40             ON SalesDocumentItem . BillToParty = BillToPartyTable . Customer
41         LEFT JOIN I_Customer AS PayerPartyTable
42             ON SalesDocumentItem . PayerParty = PayerPartyTable . Customer
43         LEFT JOIN I_AddrOrgNamePostalAddress AS SoldToPartyAddress
44             ON SoldToPartyAddress . AddressID = SoldToPartyTable . AddressID
45         LEFT JOIN I_AddrOrgNamePostalAddress AS ShipToPartyAddress
46             ON ShipToPartyAddress . AddressID = ShipToPartyTable . AddressID
47         LEFT JOIN I_AddrOrgNamePostalAddress AS BillToPartyAddress
48             ON BillToPartyAddress . AddressID = BillToPartyTable . AddressID
49         LEFT JOIN I_AddrOrgNamePostalAddress AS PayerPartyAddress
50             ON PayerPartyAddress . AddressID = PayerPartyTable . AddressID

```