ORACLE

Art of Possible with AI & Data Science

Tax Fraud Detection powered by Generative AI, Graph Analytics and Machine Learning

EMEA Data Science Team



Tax Fraud Briefly

Income Tax Evasion

- Individuals and businesses underreporting their income or inflating deductions to reduce tax liabilities.
- Data Sources: Tax Returns, Wage and Salary Data

VAT/GST Fraud

- Fraudulent schemes involving the collection of value-added tax (VAT) or goods and services tax (GST) but not remitting it to the government.
- Data Sources: Transaction data, invoices, and sales records from businesses
- Taxpayer Records

Corporate Tax Fraud

- Large corporations
 engaging in aggressive tax
 planning or transfer pricing
 schemes to minimize their
 tax obligations.
- Data Sources: Corporate Financial Statements, Transfer Pricing Documentation

Employment Tax Fraud

- Employers misclassifying employees as independent contractors, thereby avoiding payroll taxes.
- Data Sources: Payroll data



Tax Fraud Briefly

Tax Shelter Abuse

- Exploiting offshore tax shelters or abusive tax shelters to hide income from taxation.
- Data sources: Financial institution data to reveal offshore accounts and financial holdings
- Legal documents

Identity Theft and Refund Fraud

- Criminals using stolen identities to file fraudulent tax returns and claim refunds.
- Data sources: Taxpayer Identification Data,

Sales Tax Fraud

- Retailers underreporting sales or manipulating sales records to evade sales tax.
- Data sources: POS data, inventory records, sales transactions

Cryptocurrency Tax Evasion

- Analyzing blockchain data to trace cryptocurrency transactions and identify tax evasion involving digital assets.
- Cryptocurrency blockchain data, Cryptocurrency Exchange Data





"It's not only about fraud. We can analyze our payments network, our device agent networks, we can feed graph topology and machine learning algorithms. They help us uncover how much more we can do"

Stanka Dalekova

Chief Technical Lead, Paysafe

Paysafe Group

- Online Payments provider, headquartered in Canada
- Annualized transactional volume \$100bn (2020)
- Processing up to 500000 payments per day

Business requirements

- Improve accuracy of fraud detection
- minimize chargebacks from Credit Card providers while not impacting customer experience
- Automate work of analysts in fraud department
- Integrate with operational payment systems (Oracle, MSFT)

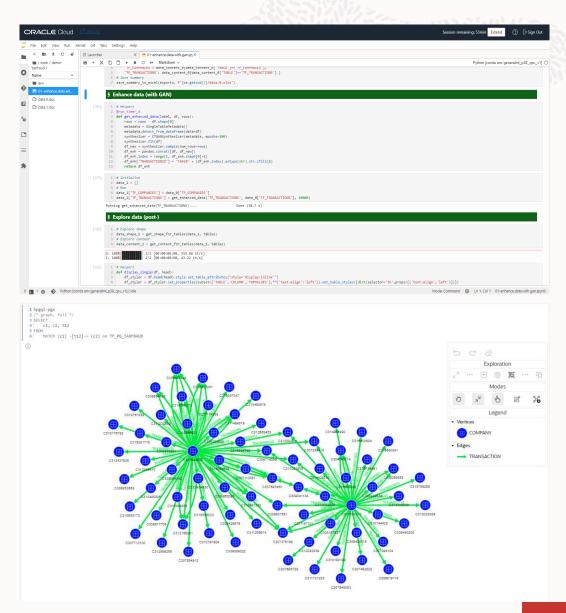
Fraud analytics platform based on Oracle Graph

- In-memory analytics engine with millisecond response time enables real-time transaction monitoring
- Integrated interactive visualization of transactions network eliminates manual work of fraud analysts
- Inclusion of account and device-related data in combination with ML improves prediction accuracy
- Seamless integration with payment systems through Oracle GoldenGate



Demo Flow

- 1. Demo Inspiration & Target Personas
- How have we achieved this?
 - Business Understanding
 - Data Enhancement with Generative Al
 - Data Exploration
 - Data Preparation (for Graph)
 - Graph Model
 - Graph Metrics
 - Data Preparation (for Machine Learning)
 - Machine Learning Model
 - Business Insights
- 3. Behind the Scenes
 - Oracle Data Platform
- 4. Next Steps





Demo Inspiration & Target Personas



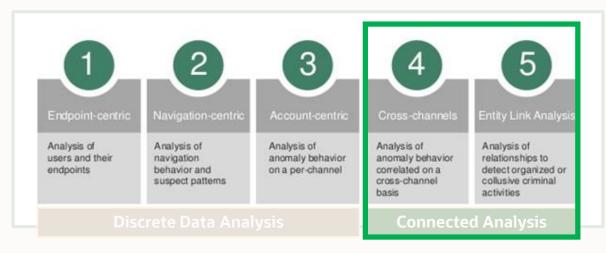
Demo Inspiration

This demo is going to showcase Oracle AI/ML platform capabilities to identify fraudulent behavior in companies through analysis of interactions.

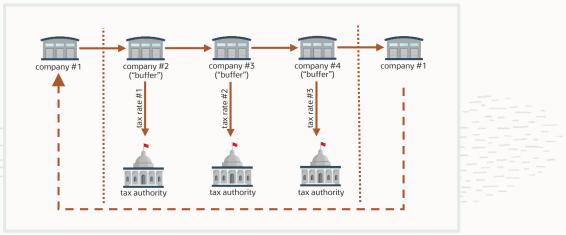
A particular type of behavior investigated includes circular money transfers, which can be indicative of abnormal interactions and can constitute sufficient cause for further analysis.

Objective is uncovering hidden fraud patterns through **graph analysis**, **machine learning** and **generative Al** techniques.

Gartner's 5 of Fraud Prevention / Detection



Circular money transfers





Demo Inspiration

With Al driven fraud detection model:

- Tax investigators will be able to identify companies with high fraud risks more effectively
- Tax authorities will better allocate investigator resources
- Government will promote a fair tax system and minimizing revenue losses

Tax fraud poses a significant challenge globally, with estimates suggesting that it costs governments billions of dollars in lost revenue each year. It not only affects public finances but also erodes trust in the tax system.

HMRC (UK tax collecting authority) estimates that losses to tax fraud amount to £16 billion each year. This is nearly half of HMRC's estimate of the tax gap (£34 billion): the difference between the amount of tax HMRC should collect each year and the amount it actually collects.*

*https://www.nao.org.uk/reports/tackling-taxfraud-how-hmrc-responds-to-tax-evasion-thehidden-economy-and-criminal-attacks/



Target Personas



Fraud Investigator wants to gather evidence to identify tax fraud investigating suspicious activities



Data Scientist wants to use Generative AI, machine learning and graph analytics to develop models and algorithms that can detect fraud patterns in large datasets.



Business Analyst wants to build final dashboards to create visually informative and interactive dashboards that enable fraud investigator to monitor and analyze suspicious activities



How have we achieved this?



Business Understanding



Fraud Investigator wants to understand any cyclic money transfers among companies, as well as importance of said companies.

Cyclic money transactions can be a potential signal of tax fraud because they often indicate an attempt to hide the true nature of financial transactions and income. A measure of importance would also be very useful for prioritization.

Transactions

Transaction ID

Company From ID

Company To ID

Transaction Date

Transaction Amount

...

Companies

Company ID

Company Name

Other Potential Data Sources

Geolocation

VAT forms/tax returns

Whistleblower

Reports

Customs and

Import/Export Data

Social Media data

3rd party data

tax fraud data sources



Data Enhancement with Generative Al



Data Scientist assesses available data and sees that existing data volume is very low. She uses Generative Al to increase data volume that helps understanding fraud cases better.

Companies Company ID Company Name ...

Transactions Transaction ID Company From ID Company To ID Transaction Date Transaction Amount

Transactions Transaction ID Company From ID Company To ID Transaction Date Transaction Amount ...

200 companies

1K transactions

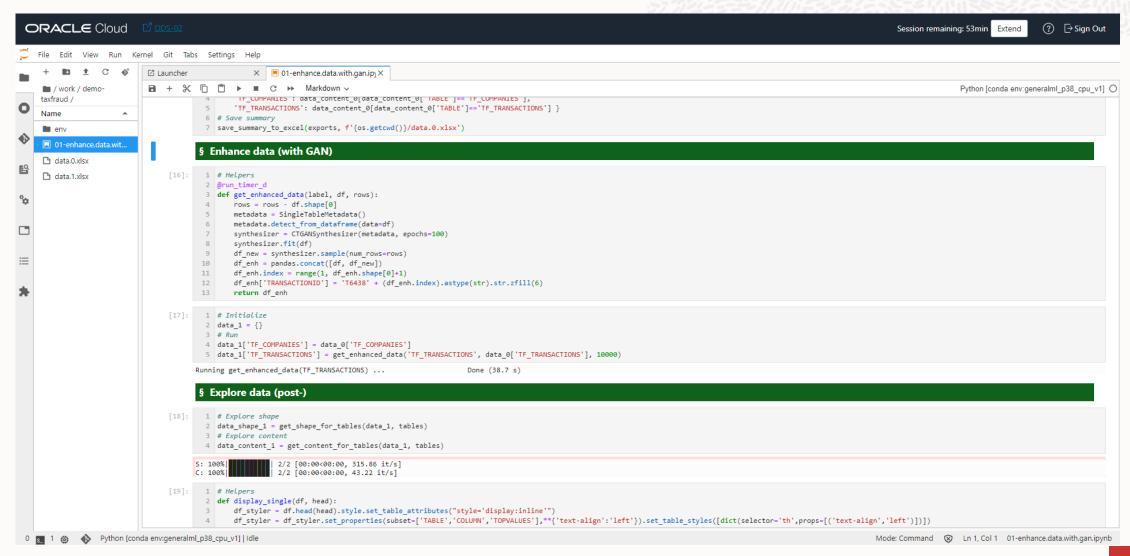


10K transactions

synthetic data generation



Data Enhancement with Generative Al



Data Exploration



Data Scientist explores the data to have a good understanding before proceeding with more advanced analysis and modeling tasks.

Companies Company ID Company Name ...

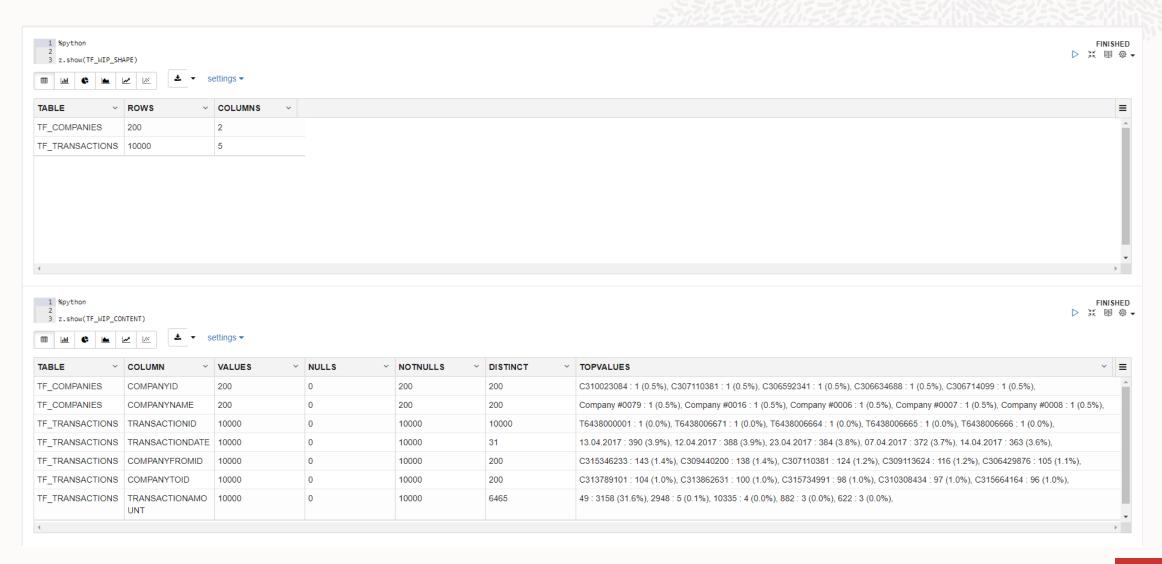
Transactions Transaction ID Company From ID Company To ID Transaction Date Transaction Amount

200 companies

10K transactions



Data Exploration



Data Preparation (for Graph)



Data Scientist transforms tabular data into a format suitable for representing companies as nodes and transaction amount as edges in a graph.

CompaniesCompany ID

Company Name

...

Transactions

Transaction ID

Company From ID

Company To ID

Transaction Date

Transaction Amount

...

Companies



200 companies

10K transactions



200 nodes 10K edges



Data Preparation (for Graph)

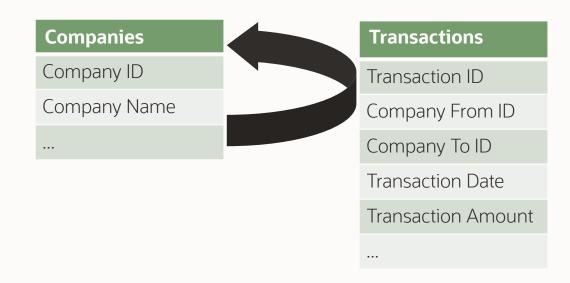
```
1 %python
                                                                                                                                FINISHED
                                                                                                                                                                                                                                                                                FINISHED
                                                                                                                         ▷ ※ 園 ◎ •
                                                                                                                                                                                                                                                                         ○ ※ 国 ※ 
    3 # imports
                                                                                                                                                     3 # sync proxy objects
    4 import warnings
                                                                                                                                                    4 TF_COMPANIES = oml.sync(table="TF_COMPANIES")
                                                                                                                                                     5 TF_TRANSACTIONS = oml.sync(table="TF_TRANSACTIONS")
    5 import oml
    6 import pandas
                                                                                                                                                    7 # pull data
    8 # settings
                                                                                                                                                    8 data = {}
    9 warnings.filterwarnings('ignore')
                                                                                                                                                    9 data['TF_COMPANIES'] = TF_COMPANIES.pull()
   10 pandas.set_option('display.width', 2000)
                                                                                                                                                   10 data['TF_TRANSACTIONS'] = TF_TRANSACTIONS.pull()
    1 %python
                                                                                                                                                                                                                                                                                FINISHED
                                                                                                                                                                                                                                                                         ▷無圓廢▼
    3 # init
    4 data_pg = {}
    6 # transform data 'COMPANIES'
    7 data_pg['TF_PGV_COMPANIES'] = data['TF_COMPANIES']
    9 # transform data 'TRANSACTIONS'
   10 data_pg['TF_PGE_TRANSACTIONS'] = data['TF_TRANSACTIONS']
   12 # check
   13 print(f'COMPANIES data now has {data_pg["TF_PGV_COMPANIES"].shape[0]} rows and {data_pg["TF_PGV_COMPANIES"].shape[1]} columns.')
   14 print(f'TRANSACTIONS data now has {data_pg["TF_PGE_TRANSACTIONS"].shape[0]} rows and {data_pg["TF_PGE_TRANSACTIONS"].shape[1]} columns.')
COMPANIES data now has 200 rows and 2 columns.
TRANSACTIONS data now has 10000 rows and 5 columns.
    1 %python
                                                                                                                                                                                                                                                                                FINISHED
                                                                                                                                                                                                                                                                         ▷ ※ 園 ◎ ▼
    3 # init
    4 data_s = data_pg
    6 # push data
          oml.drop(table='TF_PGV_COMPANIES')
           oml.drop(table='TF_PGE_TRANSACTIONS')
   10 except:
   12 TF_PGV_COMPANIES = oml.create(data_s['TF_PGV_COMPANIES'], table='TF_PGV_COMPANIES',
         dbtypes={'COMPANYID':'VARCHAR2(100)','COMPANYNAME':'VARCHAR2(100)'})
   14 TF_PGE_TRANSACTIONS = oml.create(data_s['TF_PGE_TRANSACTIONS'], table='TF_PGE_TRANSACTIONS',
         dbtypes={'TRANSACTIONID':'VARCHAR2(100)','TRANSACTIONDATE':'VARCHAR2(100)','COMPANYFROMID':'VARCHAR2(100)','COMPANYTOID':'VARCHAR2(100)','TRANSACTIONAMOUNT':'NUMBER(10)'})
   1 %script
                                                                                                                          1 %script
                                                                                                        FINISHED
                                                                                                                                                                                                                                                                               FINISHED
                                                                                                 ▷※圓碌▼
                                                                                                                                                                                                                                                                         ▷※圓碌▼
   3 alter table TF_PGE_TRANSACTIONS drop constraint FK_TF_PGE_TRANSACTIONS_FROM;
                                                                                                                          3 alter table TF_PGV_COMPANIES add constraint PK_TF_PGV_COMPANIES PRIMARY KEY (COMPANYID);
   4 alter table TF_PGE_TRANSACTIONS drop constraint FK_TF_PGE_TRANSACTIONS_TO;
                                                                                                                          4 alter table TF_PGE_TRANSACTIONS add constraint PK_TF_PGE_TRANSACTIONS PRIMARY KEY (TRANSACTIONID);
                                                                                                                          5 alter table TF_PGE_TRANSACTIONS add constraint FK_TF_PGE_TRANSACTIONS_FROM FOREIGN KEY (COMPANYFROMID) references TF_PGV_COMPANIES(COMPANYID);
6 alter table TF_PGE_TRANSACTIONS add constraint FK_TF_PGE_TRANSACTIONS_TO FOREIGN KEY (COMPANYTOID) references TF_PGV_COMPANIES(COMPANYID)
   5 alter table TF_PGE_TRANSACTIONS drop constraint PK_TF_PGE_TRANSACTIONS;
   6 alter table TF_PGV_COMPANIES drop constraint PK_TF_PGV_COMPANIES
```



Graph Model



Data Scientist builds a
Graph model to investigate
cyclic money transfers and
companies who are at the
center of these
transactions.

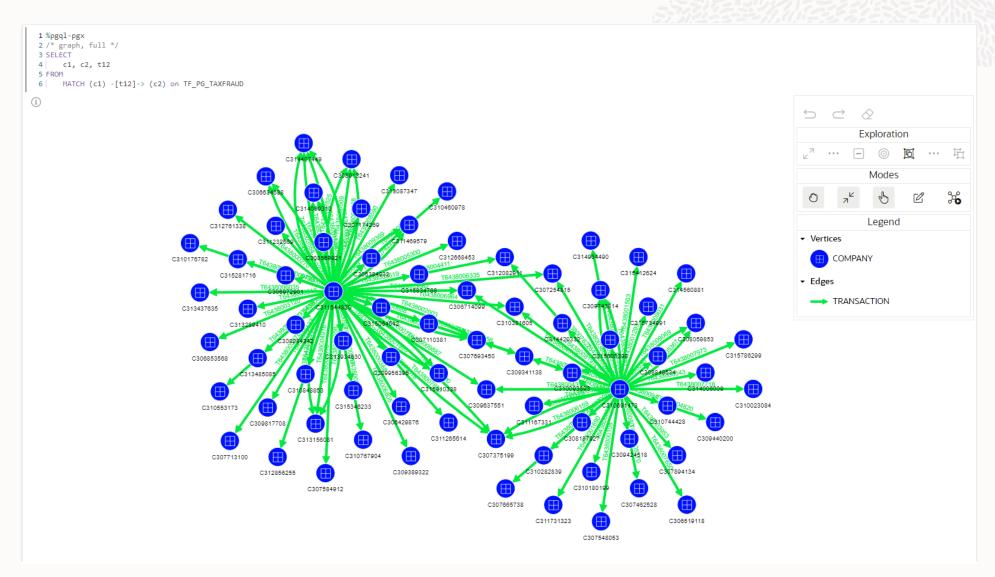


10K edges

200 nodes



Graph Model

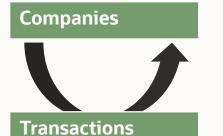


Graph Metrics



Data Scientist identifies 1-, 2- and 3-step cycles of transactions as well as calculates importance scores of the companies.







fraud patterns

importance

1-step cycles

2-step cycles

3-step cycles

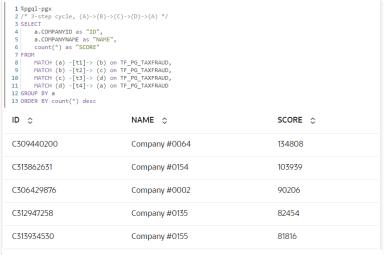
betweenness centrality closeness centrality

200 nodes 10K edges

Graph Metrics



1 %1		
1 %pgql-pgx 2 /* 2-step cycle, (A)->(B)->	(C)->(A) */	
3 SELECT		
4 a.COMPANYID as "ID", 5 a.COMPANYNAME as "NAME"		
6 count(*) as "SCORE"	,	
7 FROM	TE DE TAVERAUD	
8 MATCH (a) -[t1]-> (b) o 9 MATCH (b) -[t2]-> (c) o		
10 MATCH (c) -[t3]-> (a) o		
11 GROUP BY a 12 ORDER BY count(*) desc		
12 ONDER BY COUNTY () desc		
ID \$	NAME ≎	SCORE ≎
C309440200	Company #0064	2733
	. ,	
C313862631	Company #0154	2088
	,	
C306429876	Company #0002	1721
2300427070	company #0002	1721
C312947258	Company #0135	1665
C312747230	Company #0133	1005
C313934530	Company #01EE	1582
C313934330	Company #0155	1302



```
2 /* betweenness centrality */
                                                                                   2 /* closeness centrality */
3 SELECT
                                                                                  3 SELECT
4 a.COMPANYID as "ID",
                                                                                   4 a.COMPANYID as "ID",
5 a.COMPANYNAME as "NAME",
                                                                                  5 a.COMPANYNAME as "NAME",
                                                                                   6 a.CLOSENESS as "SCORE"
6 a.BETWEENNESS as "SCORE"
8 MATCH (a) -[t]-> (b) on TF_PG_TAXFRAUD
9 GROUP BY a
10 ORDER BY a.BETWEENNESS desc
Type to search
                      NAME ≎
ID 🗘
                                                 SCORE ♦
C309440200
                                                712.1449272244015
                     Company #0064
C313862631
                      Company #0154
                                                 537.4907683695352
C306429876
                      Company #0002
                                                 444.29669081790325
C315346233
                                                 423.10032538396234
                      Company #0186
C312947258
                      Company #0135
                                                 401.4715513258675
```

```
8 MATCH (a) -[t]-> (b) on TF_PG_TAXFRAUD
9 GROUP BY a
10 ORDER BY a.CLOSENESS desc
Type to search
ID 💠
                  NAME ≎
                                        SCORE 0
C309440200
                                        0.0038461538461538464
                  Company #0064
C313862631
                  Company #0154
                                        0.0035587188612099642
C306429876
                  Company #0002
                                        0.0035211267605633804
C315346233
                                        0.0034965034965034965
                  Company #0186
C307110381
                  Company #0016
                                        0.003484320557491289
```

1 %pgql-pgx

1 %pgql-pgx

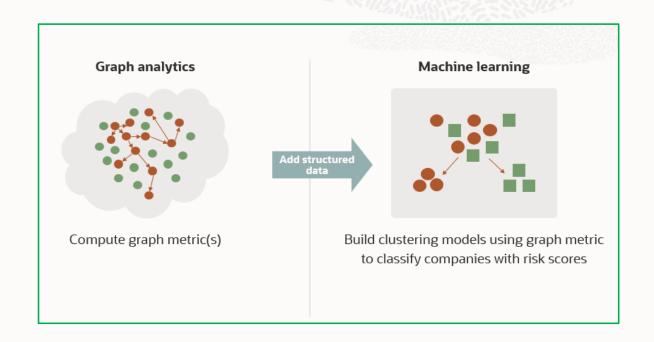
Detecting Companies with High Fraud Risk



Fraud Investigator would like to understand the list of companies with high fraud potential.



Data Scientist uses Graph metrics to segment companies to provide higher risk companies to fraud investigator.



Data Preparation (for Machine Learning)



Data Scientist prepares data for segmenting companies in terms of fraud risks.

Companies



1-step cycles2-step cycles3-step cycles

betweenness centrality closeness centrality

Companies

1-step cycles score

2-step cycles score

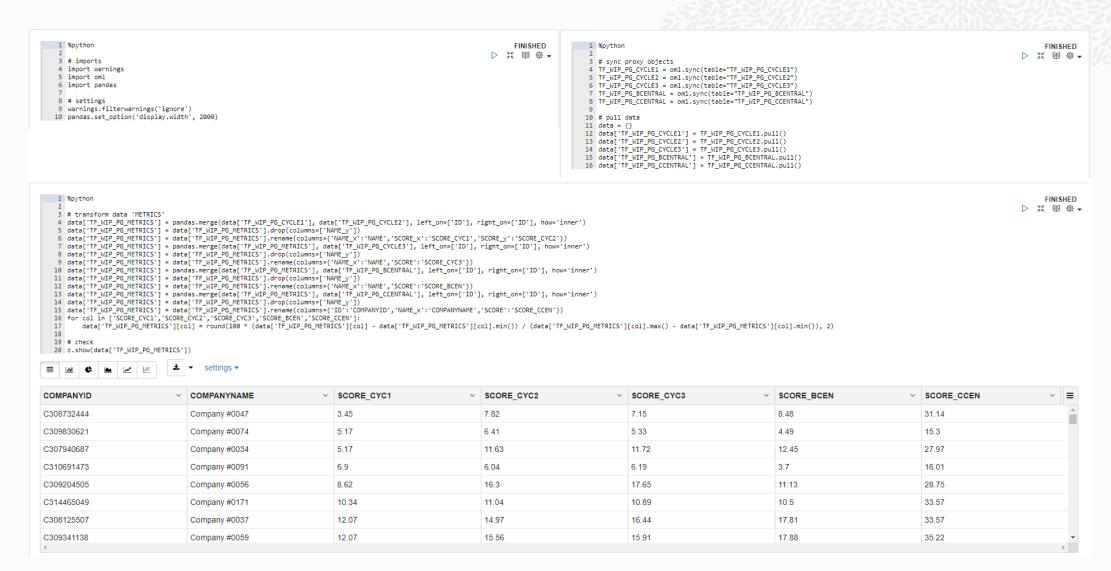
3-step cycles score

betweenness centrality score

closeness centrality score



Data Preparation (for Machine Learning)



Machine Learning Model

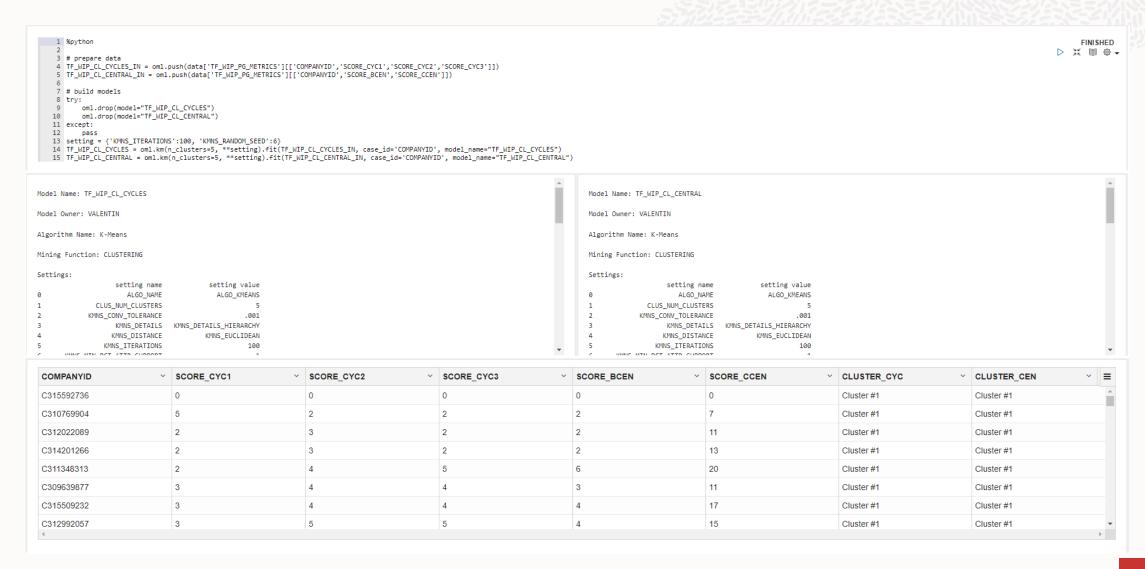


Data Scientist builds two machine learning models to detect fraud risk and propensity for fraud of the companies.

Companies 1-step cycles score 2-step cycles score 3-step cycles score betweenness centrality score closeness centrality score importance model 'fraud propensity' model

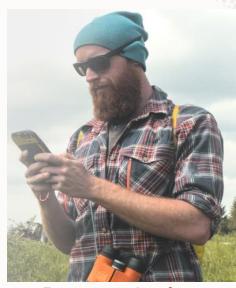


Machine Learning Model





Data Scientist shares tax fraud segments with Business Analyst.



Business Analyst creates dashboards for Fraud Investigator.

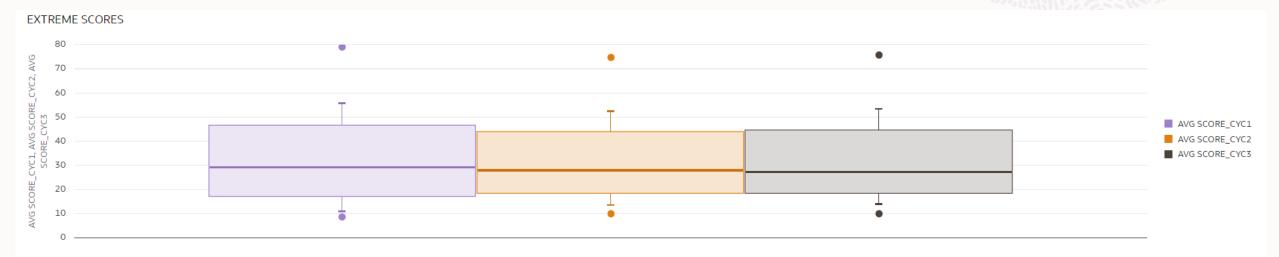


Fraud Investigator accesses all high-risk companies at the click of a button.







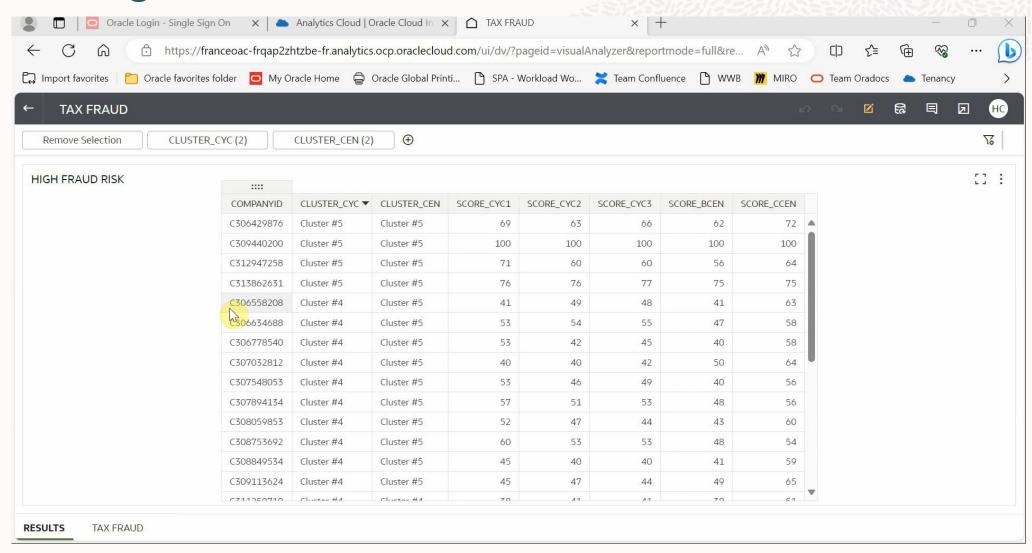


HIGH RISK COMPNIES

CLUSTER_CYC: Cluster #5 CLUSTER_CEN: Cluster #5

COMPANYID	CLUSTER_CYC	CLUSTER_CEN
C306429876	Cluster #5	Cluster #5
C309440200	Cluster #5	Cluster #5
C312947258	Cluster #5	Cluster #5
C313862631	Cluster #5	Cluster #5

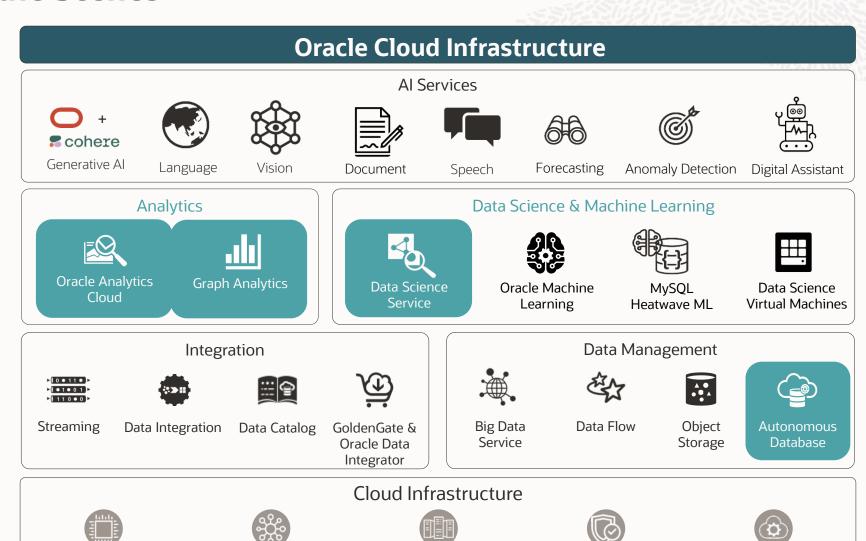






Behind the Scenes

Behind the Scenes



Storage

Security

Cloud Native

Compute

Networking

Autonomous Database



Self-Driving

- Scale-out database with fault-tolerance and DR
- Runs on enterpriseproven Exadata platform
- Full compatibility with existing enterprise databases



Self-Securing

- Automatically applies security updates online
- Secure configuration with full database encryption
- Sensitive data hidden from Oracle or customer admins



Self-Repairing

- Recovers automatically from any failure
- 99.995% uptime including maintenance
- Elastically scales compute or storage as needed



Oracle Graph

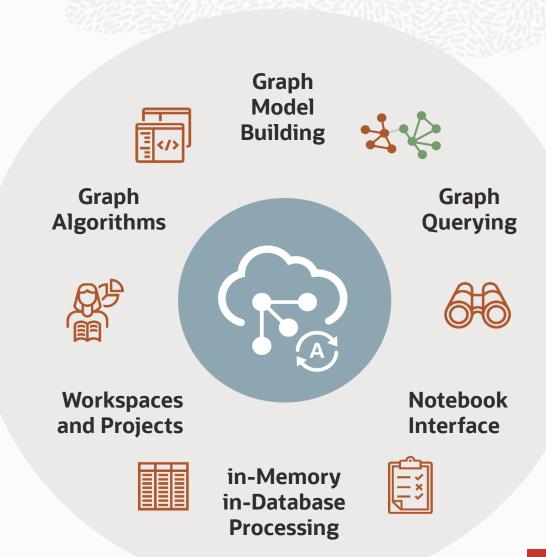
Model entities and relationships using a graph approach

Leverage comprehensive tooling and automation for graph modelling

Benefit from 50+ pre-built graph algorithms, graph querying language and visualization layer

Collaborate with teammates on shareable and reproducible graph assets

Run scalable and highly performant workloads on a managed environment





Oracle Machine Learning

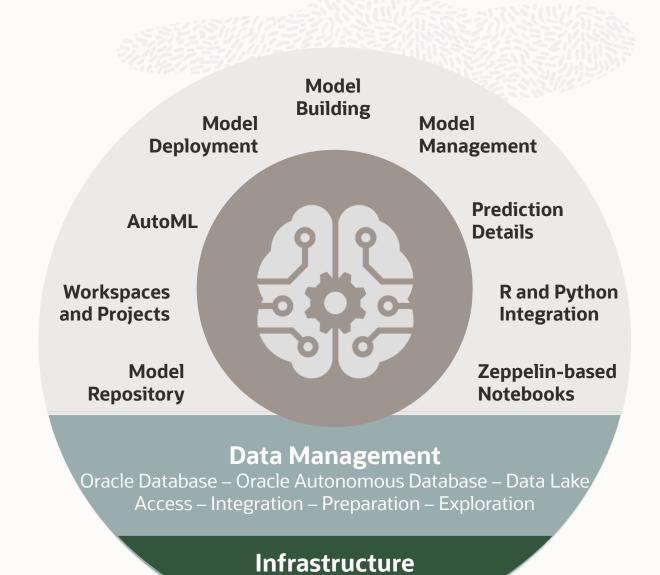
Support data scientist productivity and empower non-experts with AutoML

Gain algorithm-specific data preparation, integrated text mining, and partitioned models

Benefit from over 30+ high performance in-database machine learning algorithms

Deploy and update machine learning models in production via SQL and REST APIs

Deploy R and Python user-defined functions using managed processes



CPU – Storage – Network

Data Science Service

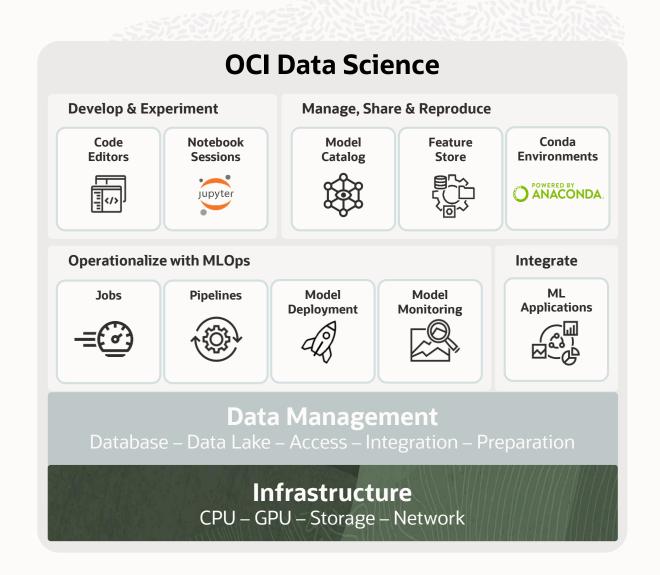
Accelerate and automate the entire end-to-end data science lifecycle

Use favorite open-source Python tools and frameworks

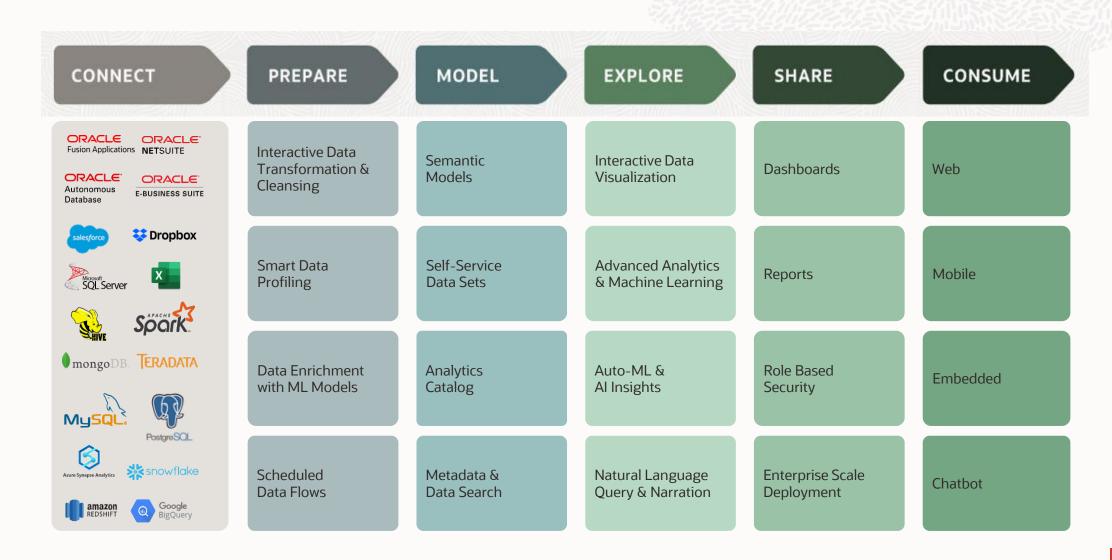
Gain enterprise-grade MLOps with flexible interfaces and unlimited scale

Collaborate with teammates on shareable and reproducible data science assets

Run large-scale workloads with access to bare metal GPUs and distributed data processing and model training

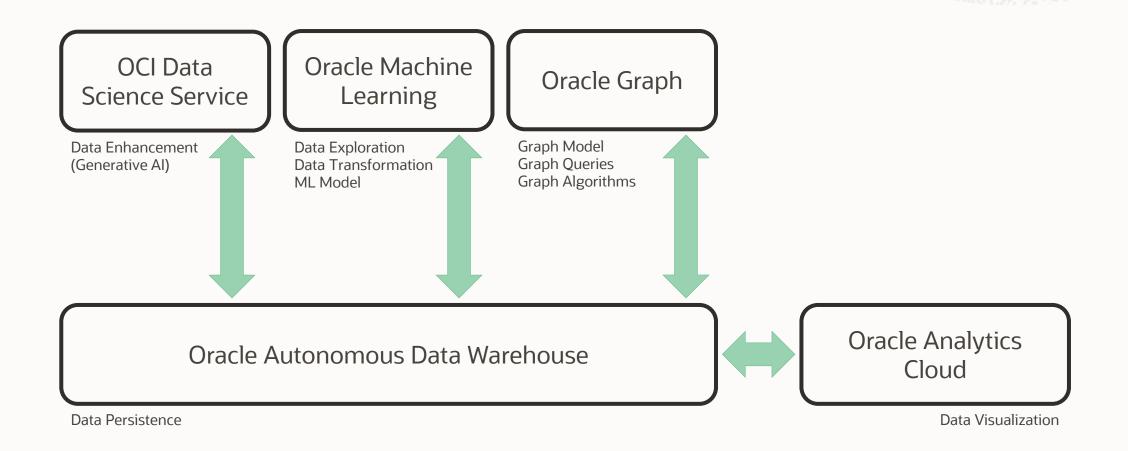


Oracle Analytics Cloud





High Level Architecture



Other Potential Contributions



ORACLE