

# Introduction to AutoML and Data Science using the Oracle Autonomous Database

#### Sandesh Rao

VP AIOps for the Autonomous Database

🕥 <u>@sandeshr</u>

in <a href="https://www.linkedin.com/in/raosandesh/">https://www.linkedin.com/in/raosandesh/</a>

https://www.slideshare.net/SandeshRao4

### Safe harbor statement

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, timing, and pricing of any features or functionality described for Oracle's products may change and remains at the sole discretion of Oracle Corporation.

Statements in this presentation relating to Oracle's future plans, expectations, beliefs, intentions and prospects are "forward-looking statements" and are subject to material risks and uncertainties. A detailed discussion of these factors and other risks that affect our business is contained in Oracle's Securities and Exchange Commission (SEC) filings, including our most recent reports on Form 10-K and Form 10-Q under the heading "Risk Factors." These filings are available on the SEC's website or on Oracle's website at <u>http://www.oracle.com/investor</u>. All information in this presentation is current as of September 2019 and Oracle undertakes no duty to update any statement in light of new information or future events.

### Agenda

- 1. Overview of ML and the Autonomous Database
- 2. Journey of the DBA to Data Scientist
- 3. OML Examples
- 4. AutoML and what's coming
- 5. Questions

### **Traditionally DBAs are Responsible for:**

### Tasks Specific to Business and Innovation

- Architecture, planning, data modeling
- Data security and lifecycle management
- Application related tuning
- End-to-End service level management

### **Maintenance Tasks**

- Configuration and tuning of systems, network, storage
- Database provisioning, patching
- Database backups, H/A, disaster recovery
- Database optimization



### **Autonomous Database Removes Generic Tasks**

Freedom from Drudgery for DBA: More Time to Innovate and Improve the Business

### **Tasks Specific to Business and Innovation**

- Architecture, planning, data modeling
- Data security and lifecycle management
- Application related tuning
- End-to-End service level management

#### Maintenance Tasks

- Configuration and tuning of systems, network, storage
- Database provisioning, patching
- Database backups, H/A, disaster recovery
- Database optimization



### The Evolution of the DBA/Database Developer Role

**Data Engineer** Architecture, "data wrangler"





**Data Security** Data classification, Data life-cycle mgmt

### Machine Learning

Solving data-driven problems Discovering insights Making predictions





**Application Tuning** SQL tuning, connection mgmt

### **Database Developer to Data Scientist Journey**

You Are Probably Already Doing Most of This Work!

Data extraction

Data wrangling Deriving new attributes ("feature engineering")

. . .

. . .

. . .

Import predictions & insights Translate and deploy ML models Automate

# Typically 80% of the work

Most data scientists spend only 20 percent of their time on actual data analysis and 80 percent of their time finding, cleaning, and reorganizing huge amounts of data, which is an inefficient data strategy<sup>1</sup>

### Eliminated or minimized with Oracle

Data Management platform becomes combine/hybrid DM + machine learning platform



# What is Machine Learning?

Algorithms *automatically* sift through large amounts of data to discover hidden patterns, new insights and make predictions



Copyright © 2020 Oracle and/or its affiliates

# **CRISP-DM** Methodology



#### DATA UNDERSTANDING

Assemble the "right data" Data profiling

- Data visualization
- Univariate statistics/group by
- Bi-variate statistics

#### DATA PREPARATION Sampling/Stratified Algorithm req'd transforms • Auto Data Preparation • Missing Values, Binning, Normalization, etc.

Normalization, etc.

- Unstructured data
- Aggregations

Domain specific transforms

• "Engineered Features" Features Selection

#### MODELING

- Algorithm settings/defaults
- Stratified sampling
- Feature selection
- Build model(s)

<sup>10</sup> Copyright © 2020 Oracle and/or its affiliates. https://en.wikipedia.org/wiki/Cross-industry\_standard\_process\_for\_data\_mining

# **Oracle Machine Learning**





# Database Developer to Data Scientist Journey Six Major Steps (Oracle Machine Learning POV)

- Business Understanding—Week 1
- Data Understanding—Week 2
- Data Preparation—Week 3
- Modeling (ML)—Week 4
- Evaluation—Week 5
- Deployment—Week 6



https://en.wikipedia.org/wiki/Cross-industry\_standard\_process\_for\_data\_mining

# Week 1—Business Understanding

Start with a Well-Defined Business Problem Statement

| Poorly Defined                     | Better  | Data Mining<br>Technique |
|------------------------------------|---|--------------------------|
| Predict employees that leave       | <ul> <li>Based on past employees that voluntarily left:</li> <li>Create New Attribute EmplTurnover → O/1</li> </ul>   |                          |
| Predict customers that churn       | <ul> <li>Based on past customers that have churned:</li> <li>Create New Attribute Churn → YES/NO</li> </ul>   |                          |
| Target "best" customers            | <ul> <li>Recency, Frequency Monetary (RFM) Analysis</li> <li>Specific Dollar Amount over Time Window:         <ul> <li>Who has spent \$500+ in most recent 18 months</li> </ul> </li> </ul> |                          |
| How can I make more \$\$?          | • What helps me sell soft drinks & coffee?  |                          |
| Which customers are likely to buy? | <ul> <li>How much is each customer likely to spend?</li> </ul>  | 0000                     |
| Who are my "best customers"?       | <ul> <li>What descriptive "rules" describe "best customers"?</li> </ul>   |                          |
| How can I combat fraud?            | <ul> <li>Which transactions are the most anomalous?</li> <li>Then roll-up to physician, claimant, employee</li> </ul>   | x2<br>x1                 |

# Week 1—Business Understanding Be Extremely Specific in your Problem Statement

#### Target "best" customers who have GOOD CREDIT and make payments



# Week 2—Data Understanding

📼 🔅 defa

### Review the Data; Does it Makes Sense?

▷洗餌

100

%sql

/\* This shows the credit scoring data we will use historical data to predict the likelihoo

Select \* from credit\_scoring\_100k\_v where rown

Credit Score Predictions

STEP 5: Review Credit Scoring Data

TOOK 4 Sec. Last updated by CHARLIE at February 07 2019, 3.29.45 PM

▦ Ш ІІ 🖬 🗠 🖄

Are AGEs all positive, 0-120?
Are INCOME values weekly or monthly?
Are the LOAN\_AMOUNTS reasonable?

> ※ 目 焱

| ION_LEVEL TENURE LOAN_TYPE LOAN_AMOUNT LOAN_LENGTH GENDER  | EDUCATION_LEVEL TENU  | WEALTH                          | NUMBER_OF_LIABLES | MARITAL_STATUS                           | INCOME                       | AGE                  | CUSTOMER_ID                          |
|--|---|---------------------------------|-------------------|--|------------------------------|----------------------|--------------------------------------|
|  |   |                                 | ▼                 | ▼  |                              |                      |                                      |
| Degree 15 Education 15000 12 Male  | Master's Degree 15  | Average                         | 5                 | Single                                   | 3750                         | 27                   | 41908                                |
| Degree 5 Housing 55000 72 Male   | Master's Degree 5   | Average                         | 2                 | Married                                  | 2150                         | 24                   | 54255                                |
| Degree 24 Education 20000 12 Female  | Master's Degree 24  | Average                         | 2                 | Married                                  | 2150                         | 37                   | 46270                                |
| Degree 17 Need 35000 12 Male   | Master's Degree 17  | Average                         | 1                 | Divorced                                 | 2150                         | 47                   | 54244                                |
| Degree 24 Housing 95000 108 Female   | Master's Degree 24  | Average                         | 2                 | Married                                  | 2750                         | 21                   | 49074                                |
| s Degree 7 Need 30000 24 Female  | Bachelor's Degree 7   | Poor                            | 5                 | Married                                  | 2750                         | 63                   | 10647                                |
| s Degree 11 Housing 80000 60 Male  | Bachelor's Degree 11  | Poor                            | 3                 | Married                                  | 2750                         | 28                   | 9535                                 |
| s Degree 11 Need 15000 12 Male   | Bachelor's Degree 11  | Poor                            | 3                 | Married                                  | 2150                         | 67                   | 34680<br>∢                           |
| Degree24Housing95000108Females Degree7Need3000024Females Degree11Housing8000060Males Degree11Need1500012Male | Master's Degree24Bachelor's Degree7Bachelor's Degree11Bachelor's Degree11 | Average<br>Poor<br>Poor<br>Poor | 2<br>5<br>3<br>3  | Married<br>Married<br>Married<br>Married | 2750<br>2750<br>2750<br>2150 | 21<br>63<br>28<br>67 | 49074<br>10647<br>9535<br>34680<br>∢ |

Etc....

Took 5 sec. Last updated by CHARLIE at February 07 2019, 3:29:56 PM.

# Week 2—Data Understanding

### Review the Data; Does it Makes Sense?



# Week 2—Data Understanding

### Review the Data; Does it Makes Sense?



Copyright © 2020 Oracle and/or its affiliates.

18

# Week 3—Data Preparation

Prepare the Data, Create New Derived Attributes or "Engineered Features"

| Source Attribute           | New Attribute/"Engineered Feature"            |
|----------------------------|---|
| Date of Birth              | <br>AGE                                       |
| Address                    | DISTANCE_TO_DESTINATION<br>COMMUTE_TIME       |
| Call detail records (CDRs) | <br>#_DROPPED_CALLS                           |
|                            | <br>PERCENT_INTERNATIONAL                     |
| Salary                     | <br>PERCENT_VS_PEERS                          |
| Purchases                  | <br>TOTALS_PER_CATEGORY (e.g. Food, Clothing) |

# Week 3—Data Preparation

### Prepare the Data, Create New Derived Attributes or "Engineered Features"

Oracle Machine Learning's Auto Data Prep (ADP) and ML algorithms are designed with intelligent defaults and can automatically deal with:

- Missing values
- Outliers
- Binning
- Too many distinct values
- Too many constants
- Trans data/aggregations
- Unstructured data
- Correlated data

| Build Partition Sampling In     | put Text         | Show          |             |  |                                   |
|---------------------------------|------------------|---------------|-------------|--|-----------------------------------|
| Columns: 28 included out of 31. | sing neuristics) | SHOW          |             | <b>∳</b>   → ⇒   | Mining Type Auto Prep Q Name      |
| Name                            | Data Type        | Input         | Mining Type | Auto P   | Rules                             |
| XYZI AGE                        | NUMBER           | ->            | N.          | <ul> <li>Image: A second s</li></ul> |                                   |
| XX BANK FUNDS                   | NUMBER           | ÷             | N.          | <ul> <li>Image: A start of the start of</li></ul>  |                                   |
| BUY_INSURANCE                   | VARCHAR2         | ->            | 52          | 1  |                                   |
| XXE CAR_OWNERSHIP               | NUMBER           | ->            | 5           | <ul> <li>Image: A start of the start of</li></ul>  | Change mining type to Categorical |
| XX CHECKING_AMOUNT              | NUMBER           | ->            | 15          | <ul> <li>Image: A start of the start of</li></ul>  |                                   |
| CREDIT_BALANCE                  | NUMBER           | ->            | 15          | <ul> <li>Image: A start of the start of</li></ul>  |                                   |
| XXE CREDIT_CARD_LIMITS          | NUMBER           | ⇒             | 1.5         | ~  |                                   |
| 🖙 CUST_ID                       | VARCHAR2         | -048          | 52          | 1  | Exclude for all models            |
| XXE FIRST                       | VARCHAR2         |               | 5           | ~  | Exclude for all models            |
| HAS_CHILDREN                    | NUMBER           | ⇒             | 5           | ~  | Change mining type to Categorical |
| HOUSE_OWNERSHIP                 | NUMBER           | ÷             | 5           | <b>~</b>   | Change mining type to Categorical |
| XVII LAST                       | VARCHAR2         | -28           | 2           | <b>~</b>   | Exclude for all models            |
| XYZ LTV                         | NUMBER           | ÷             | 15          | <b>~</b>   |                                   |
| XY2 LTV_BIN                     | VARCHAR2         | ÷             | 5           | <b>~</b>   |                                   |
| MARITAL_STATUS                  | VARCHAR2         | $\rightarrow$ |             | Image: A start of the start           |                                   |
| MONEY_MONTLY_OVERDRAWN          | NUMBER           | ÷             | <u>N</u>    | <ul> <li>Image: A set of the set of the</li></ul>  |                                   |
| MONTHLY_CHECKS_WRITTEN          | NUMBER           | ÷             | 11          | <b>~</b>   |                                   |
| MORTGAGE_AMOUNT                 | NUMBER           | ->            | 11          | <b>~</b>   |                                   |
| N_MORTGAGES                     | NUMBER           | ÷             | 2           | <b>~</b>   | Change mining type to Categorical |
| N_OF_DEPENDENTS                 | NUMBER           | ⇒             | 15          | <b>~</b>   |                                   |
| XYEN TRANS ATM                  |                  | ھ             | 1 et        |  | ,                                 |

# Week 4—Modeling (Machine Learning)

First, Identify the Key Attributes That Most Influence the Target Attribute

|  |  | Charlie Phase 2 New N   | lotebooks 🔻 😍 CHARLIE 🔻       |
|--|--|---|-------------------------------|
| MODEL BUILDING   |  |   | *                             |
| Took 0 sec. Last updated by CHARLIE at March 31 2020, 6:48:22 PM. (outdated)   |  |   |                               |
| In this notebook, we will build attribute importance and classification related topics.<br>https://docs.oracle.com/en/databas<br>Took 0 sec. Last updated by CHARLIE at April 01 most influence  | ation models. Please refer to the Oracle Machine<br>Importance to i<br>Jence the targe | Learning Documentation for additional info<br>dentify key va<br>t attribute | rmation on UFINISHED N ≥2 1 ® |
| Run attribute importance to identified the second s | Display the Top N Attributes for Good  | Credit Customers  | FINISHED D X 目                |
| "GOOD_CREDIT" and "OTHER_CREDIT"   | SELECT * FROM DM\$VAai_explain_output_credi  | t_score_bin WHERE ROWNUM <15;   | *<br>*                        |
| %script  |  | settings ▼  |                               |
| Find the importance of attributes that<br>independently impact the target attribute:<br>CREDIT_SCORE_BIN   | 0.428  |   | ATTRIBUTE_IMPORTANCE          |
| <pre>BEGIN DBMS_DATA_MINING.DROP_MODEL    ('ai_explain_output_credit_score_bin');</pre>  | 0.4  |   |                               |
| EXCEPTION WHEN OTHERS THEN NULL; END;  | 0.35   |   |                               |
| DECLARE<br>v_set1st DBMS_DATA_MINING.SETTING_LIST;<br>BEGIN  | 0.3  |   |                               |
| <pre>v_setlst('ALGO_NAME') := 'ALGO_AI_MDL';<br/>V setlst('PREP AUTO') := 'ON';</pre>  | 0.2  |   |                               |
| DBMS_DATA_MINING.CREATE_MODEL2(  | 0.15   |   | _                             |
| <pre>MUDEL_NAME =&gt;     'ai_explain_output_credit_score_bin',     MINING FUNCTION =&gt;</pre>  | 0.1  |   |                               |
| <pre>'ATTRIBUTE_IMPORTANCE',<br/>DATA_QUERY =&gt; 'select * from<br/>CREDIT_SCORING_100K_V',<br/>CREDIT_SCORING_100K_V',</pre>   | 0.05<br>0  |   |                               |
| CASE TO COLUMN NAME -> 'CUSTOMER TO'   | ш  | > 0   | ω                             |

21 Copyright © 2020 Oracle and/or its affiliates.

# Week 4—Modeling (Machine Learning) Training and Testing ML Models using 60/40% Random Samples



# Week 4—Modeling (Machine Learning)

Build multiple models with different algorithms and settings



# Week 5—Model Evaluation (ML)

### Next, test model accuracy

Randomly selected "hold out" sample of data that was used to train the ML model

Compute Cumulative Gains, Lift, Accuracy, etc.

Review the attributes used in the model and model coefficients

Make sure the model makes sense



# Week 6—Deployment

### Apply the Models to Predict "Best Customers"

Simple SQL Apply scripts run 100% inside the Database for immediate ML model deployment

| select a.customer_id<br>, a.prob_Credit_Score_Bin<br>, b.age, b.income, b.tenure, b.loan_type, b.loan_amount, b.occupation, b.education_level, b<br>marital_status<br>from (select * from (select Customer_id, round(prob_Credit_Score_Bin *100,2) prob_Credit_Score_Bin<br>(select Customer_ID, prediction_probability(N1_CLASS_MODEL, NULL using *) prob_Credit_Score_Bin<br>credit_scoring_100k_v b<br>where a.customer_id = b.customer_id<br>order by a.prob_Credit_Score_Bin desc<br><b>CUSTOMER_ID PROB_CREDIT_SCORE_BIN AGE INCOME TENURE LOAN_TYPE</b><br>34673 100 31 5250 8 Education<br>77936 100 37 6250 6 Need<br>56154 100 45 4250 10 Need<br>11610 100 63 4250 4 Housing<br>56733 100 54 4250 33 Housing   | //JQ1   |  |   |                                     | Mo   | del Apply                                       | /"Scorin   | g″             |
|---|---|--|---|-------------------------------------|--|---|--|----------------|
| Image: Second problemedit_Scone_Bin desc         Image: Scone_Bin desc  | <pre>select a.cus     , a.prob     , b.age,     .man from (select     (select     credit_sco where a.cust</pre> | <pre>tomer_id<br/>_Credit_Score_Bin<br/>b.income, b.tenur<br/>ital_status<br/>* from (select Cu<br/>Customer_ID, predi<br/>coring_new_cust_v)<br/>ring_100k_v b<br/>omer_id = b.custom</pre> | e, b.loan_type, b.lo<br>stomer_id, round(pro<br>ction_probability(N1<br>)) a<br>er_id | an_amount<br>b_Credit_<br>_CLASS_MO | , b.occupati<br>Score_Bin *1<br>DEL, NULL us   | on, b.educati<br>00,2) prob_Cr<br>ing *) prob_C | on_level, b<br>edit_Score_Bi<br>redit_Score_B                | n fro<br>in fo |
| Image: | order by a.p  | rob_Credit_Score_B   | in desc   |                                     |  |   |  |                |
| CUSTOMER_ID       PROB_CREDIT_SCORE_BIN       AGE       INCOME       TENURE       LOAN_TYPE         34673       100       31       5250       8       Education         77936       100       37       6250       6       Need         56154       100       45       4250       10       Need         11610       100       28       6250       3       Housing         56733       100       63       4250       4       Housing         57999       100       54       4250       33       Housing   |   |  | <u>+</u> -  |                                     |  |   |  |                |
| 34673       100       31       5250       8       Education         77936       100       37       6250       6       Need         56154       100       45       4250       10       Need         11610       100       28       6250       3       Housing         56733       100       63       4250       4       Housing         57999       100       54       4250       33       Housing   |   |  |   |                                     |  |   |  |                |
| 346731003152508Education779361003762506Need5615410045425010Need116101002862503Housing567331006342504Housing5799910054425033Housing  | CUSTOMER  |  |   | AGE                                 | INCOME   |   | LOAN TYPE  | =              |
| 779361003762506Need5615410045425010Need116101002862503Housing567331006342504Housing5799910054425033Housing  | CUSTOMER  | ID ▼ PROB_CRE  | DIT_SCORE_BIN   | AGE                                 |  | TENURE  | LOAN_TYPE  |                |
| 5615410045425010Need116101002862503Housing567331006342504Housing5799910054425033Housing   | CUSTOMER_<br>34673  | <b>ID PROB_CRE</b>   | EDIT_SCORE_BIN  | <b>AGE</b> 31                       | INCOME<br><b>5</b> 250                         | TENURE V  | LOAN_TYPE<br>Education                                       |                |
| 11610       100       28       6250       3       Housing         56733       100       63       4250       4       Housing         57999       100       54       4250       33       Housing  | CUSTOMER_<br>34673<br>77936   | <b>ID PROB_CRE</b><br>100<br>100   | DIT_SCORE_BIN   | <b>AGE</b><br>31<br>37              | INCOME<br>5250<br>6250                         | <b>TENURE *</b><br>8<br>6                       | LOAN_TYPE<br>Education<br>Need                               |                |
| 56733       100       63       4250       4       Housing         57999       100       54       4250       33       Housing  | CUSTOMER_<br>34673<br>77936<br>56154  | <b>_ID</b> ▼ <b>PROB_CRE</b><br>100<br>100<br>100  | DIT_SCORE_BIN   | AGE<br>31<br>37<br>45               | INCOME<br>5250<br>6250<br>4250                 | <b>TENURE</b><br>8<br>6<br>10                   | LOAN_TYPE<br>Education<br>Need<br>Need                       |                |
| 57999 100 54 4250 33 Housing  | CUSTOMER<br>34673<br>77936<br>56154<br>11610  | _ID ▼ PROB_CRE<br>100<br>100<br>100<br>100<br>100  | EDIT_SCORE_BIN  | AGE<br>31<br>37<br>45<br>28         | INCOME<br>5250<br>6250<br>4250<br>6250         | <b>TENURE</b><br>8<br>6<br>10<br>3              | LOAN_TYPE<br>Education<br>Need<br>Need<br>Housing            |                |
|   | CUSTOMER<br>34673<br>77936<br>56154<br>11610<br>56733   | _ID ▼ PROB_CRE<br>100<br>100<br>100<br>100<br>100<br>100   | EDIT_SCORE_BIN  | AGE<br>31<br>37<br>45<br>28<br>63   | INCOME<br>5250<br>6250<br>4250<br>6250<br>4250 | <b>TENURE</b> 8 6 10 3 4                        | LOAN_TYPE<br>Education<br>Need<br>Need<br>Housing<br>Housing |                |

FINISHED D 光 目 🕸

Apply the Oracle Machine Learning Model to New Customers to

Took 0 sec. Last updated by ADWC\_WS2 at October 17 2018, 2:51:03 PM. (outdated)

# Week 6—Deployment

Apply the Models to Predict "Best Customers"

Simple SQL Apply scripts run 100% inside the Database for model build, model apply and immediate ML model deployment

|  | Claims and  | omaly detection script  | sq/ ×   |
|--|---|---|---|
| SQL Worksheet Histo  | iry   |   |   |
| 🕨 📃 🕲 🗸 🎉 🗟  | l 🕼 🛃 I 🔏   | 🗄 🥔 🐻 🗛 i   |   |
| Worksheet Query  | Builder   |   |   |
| drop table (<br>exec dbms_da<br>create table<br>insert into<br>insert into<br>commit;                              | CLAIMS_SET;<br>ata_mining.dr<br>• CLAIMS_SET<br>CLAIMS_SET V<br>CLAIMS_SET V                    | cop_model('CLAIMSMC<br>(setting_name vard<br>ralues ('ALGO_NAME'<br>ralues ('PREP_AUTO'   | <pre>DEL');<br/>:har2(30), setting_value varchar2(4000));<br/>,'ALGO_SUPPORT_VECTOR_MACHINES');<br/>,'ON');</pre> |
| <pre>BEGIN DBMS_DAT: model_ mining data_t: case_it target settin END;</pre>  | A_MINING.CREA<br>name<br>function<br>able_name<br>d_column_name<br>column_name<br>ys_table_name | <pre>TTE_MODEL( =&gt; 'CLAIMSMODEL', =&gt; dbms_data_mini =&gt; 'CLAIMS', :=&gt; 'POLICYNUMBER' =&gt; null, :=&gt; 'CLAIMS_SET');</pre> | Model Build   |
|  |   |   |   |
| Select * fro<br>(select POL)<br>rank (<br>(select POL)<br>from CLAIMS<br>where PASTN<br>where rnk <<br>order by pe | om<br>ICYNUMBER, ro<br>) over (order<br>ICYNUMBER, pr<br>UMBEROFCLAIMS<br>= 5<br>rcent_fraud d  | <pre>pund(prob_fraud*100 by prob_fraud des ediction_probabili in ('2to4', 'more lesc;</pre>   | <pre>IVIOUELADDIY ',2) percent_fraud, ic) rnk from .ty(CLAIMSMODEL, '0' using *) prob_fraud :than4')))</pre>      |
|  |   |   |   |
| Script Output ×  |   |   |   |
| 📌 🥔 🖯 📇 📃 I  | Task complete   | ed in 3.825 seconds   |   |
| Commit complete.   |   |   |   |
| PL/SQL procedure :   | successfully  | completed.  |   |
| PL/SQL procedure :   | successfully<br>ENT_FRAUD   | completed.<br>RNK   | Roculto   |

### **Congratulations!** Almost there ©



# Ainlona THIS CERTIFICATE IS PRESENTED TO

# Data Scientist

#### LOREM IPSUM DOLOR SIT AMET

Obtain a signed certificate. Obtaining a signed certificate involves creating a certificate signing request (CSR) and sending it to a CA in accordance with the CA's enrollment process. After conducting some checks on your company, the CA signs your request, encrypts it with a private key, and sends you a validated certificate. See the instructions provided by the CA for more information.



DATE

------SIGNATURE



### **Statistical Functions**

#### Simple SQL Syntax—Statistical Comparisons (t-tests)

Compare AVE Purchase Amounts Men vs. Women Grouped\_By INCOME\_LEVEL

| SELECT SUBSTR(cust_income_level, 1, 22) income_level,                                       |
|---|
| AVG(DECODE(cust_gender, 'M', amount_sold, null))                               sold_to_men, |
| AVG(DECODE(cust_gender, 'F', amount_sold, null))  |
| STATS_T_TEST_INDEPU(cust_gender, amount_sold, 'STATISTIC', 'F') t_observed                  |
| STATS_T_TEST_INDEPU(cust_gender, amount_sold)    two_sided_p_value                          |
| FROM customers c, sales s   |
| WHERE c.cust_id = s.cust_id   |
| GROUP BY ROLLUP(cust_income_level)  |
| ORDER BY income_level, sold_to_men, sold_to_women, t_observed;                              |

#### Query Result ×

| ړ 🖈 |  | 62 | <b>*</b> | SQL | All Rows Fetched: 14 in 1.523 seconds |  |
|-----|--|----|----------|-----|---------------------------------------|--|
|-----|--|----|----------|-----|---------------------------------------|--|

|    | INCOME_LEVEL         | SOLD_TO_MEN    | SOLD_TO_WOMEN     | T_OBSERVED     | TWO_SIDED_P_VALUE     | 1 |
|----|----------------------|----------------|-------------------|----------------|-----------------------|---|
| 1  | A: Below 30,000      | 105.2834897729 | 99.42814466653473 | -2.05425922984 | 0.039964704379552678  | 4 |
| 2  | B: 30,000 - 49,999   | 102.5965095067 | 109.8296418272003 | 2.969223321889 | 0.0029877419365879512 | 4 |
| 3  | C: 50,000 - 69,999   | 105.6275880730 | 110.1279310121247 | 2.349685400926 | 0.018792276771129993  |   |
| 4  | D: 70,000 - 89,999   | 106.6302994897 | 110.4728699326023 | 2.268392806338 | 0.023307831257217089  | K |
| 5  | E: 90,000 - 109,999  | 103.3967414937 | 101.6104162583700 | -1.26035091954 | 0.20754566236328209   |   |
| 6  | F: 110,000 - 129,999 | 106.7647596205 | 105.9813119482142 | -0.60580010770 | 0.54464855287037528   |   |
| 7  | G: 130,000 - 149,999 | 108.8775321810 | 107.3137698570293 | -0.85219780969 | 0.39410775484348759   |   |
| 8  | H: 150,000 - 169,999 | 110.9872579252 | 107.1521911799573 | -1.94514858879 | 0.051762623899376248  |   |
| 9  | I: 170,000 - 189,999 | 102.8082379709 | 107.4355601412162 | 2.149669205899 | 0.031587875078399455  |   |
| 10 | J: 190,000 - 249,999 | 108.0405638372 | 115.3433560297627 | 2.547498669040 | 0.010854966021230945  | K |
| 11 | K: 250,000 - 299,999 | 112.3779929260 | 108.1960973300511 | -1.41155136806 | 0.15809167565415438   |   |
| 12 | L: 300,000 and above | 120.9702345758 | 112.2163421398336 | -2.07261936500 | 0.038225611271820475  |   |
| 13 | (null)               | 106.6637691587 | 107.2763858209415 | 1.078537818864 | 0.28079420737509053   |   |
| 14 | (null)               | 107.1218447412 | 113.8044098205854 | 0.689462437157 | 0.4905957646106357    |   |
|    |                      |                |                   |                |                       |   |

 STATS\_T\_TEST\_INDEPU (SQL) Example;
 P\_Values < 05 show statistically significantly differences in the amounts purchased by men vs. women



### **OAA Model Build and Real-time SQL Apply Prediction**

Simple SQL Syntax—Attribute Importance - ML Model Build (PL/SQL)



# OML for R Model Build Simple R Language Syntax—Attribute Importance

### ML Model Build (R)

```
> ore.odmAI (BUY_INSURANCE ~ ., CUST_INSUR_LTV)
```

```
Call:
ore.odmAI(formula = BUY INSURANCE ~ ., data = CUST INSUR LTV)
```

### **Model Results (R)**

| importance   | rank   |
|--------------|--|
| 0.2161187797 | 1  |
| 0.1489347141 | 2  |
| 0.1463026512 | 3  |
| 0.1155879786 | 4  |
| 0.1095178647 | 5  |
|              | importance<br>0.2161187797<br>0.1489347141<br>0.1463026512<br>0.1155879786<br>0.1095178647 |





# OML for Python Model Build—*Coming soon!* Simple Python Language Syntax—Attribute Importance

### **ML Model Build (Python)**

```
> ai_mod = ai(**setting) # Create AI model object
> ai_mod = ai_mod.fit(train_x, train_y)
```

### **Model Results (Python)**

Tmnortanco.

| impor cance.            |              |      |
|-------------------------|--------------|------|
|                         |              |      |
| variable                | importance   | rank |
| BANK_FUNDS              | 0.2161187797 | 1    |
| MONEY_MONTLY_OVERDRAWN  | 0.1489347141 | 2    |
| N_TRANS_ATM             | 0.1463026512 | 3    |
| N TRANS TELLER          | 0.1155879786 | 4    |
| T_AMOUNT_AUTOM_PAYMENTS | 0.1095178647 | 5    |





# **Oracle Machine Learning**

### Machine Learning Notebooks included in Autonomous Databases

Key Features:

- Collaborative UI for data scientist and analysts
- Packaged with Autonomous Databases
- Quick start Example notebooks
- Easy access to shared notebooks, templates, permissions, scheduler, etc.
- OML4SQL
- OML4Py coming soon
- Supports deployment of OML models

|      |  | _ <b>E</b> * Machine  | Learning  |   | III NEWEST OML NOTEBOOKS [Charlie 🔻 😍 CHARLIE 💌  |   |
|------|--|---|---|---|--|---|
|      |  |   |   |   |  | • Connected   |
|      | Targeting To   | op Custo  | mers 10K_1  |   | 🚍 🏠 default 🕶  |   |
| ases | Finding "BES<br>Heather has spent m<br>of her time waiting for<br>processed there. The<br>simple SQL comman<br>incoming data on the<br>no change to analysis<br>Updated August 2019<br>Copyright (c) 2019 O  | ST" Lifetim<br>ost of her time ow<br>r jobs to finish and<br>a liternative cloud<br>ds in ADW are far<br>fly, and allow end<br>s/reporting Data V<br>9 By: Siddesh Ujjr<br>racle and/or its aff | te Value (LTV)<br>er the past couple of year<br>d very little of her time an<br>solution is more complex<br>d user analysts to immedii<br>//sualization toolset that u<br>i, Dhvani Shehr, Charlie<br>filiates. All rights reserved | Customers using C<br>s extracting and preparing data<br>hydring the data. Demands from,<br>and has no direct out of the by<br>hely fast, leveraging all the perfor-<br>tably see mining results. This dr<br>sers are familiar with.<br>Barger |  |   |
|      | Took 2 sec. Last updated by 0  | CHARLIE at February 13 2  | 2020, 5:18:38 PM. (outdated)  |   | Took 0 sec. Last updated by CHARLIE at February 06 2020, 5:40.45 PM. (outdated)  |   |
| etc. | View 1000 rows of data     FINISHED ▷ X ()) ●       [sal     This shows the first 1000 rows of customer's credit scoring data we will use to predict the lifetime value of a customer.       SELECT * FROM CREDIT_SCORING_IANK WHERE ROWNW < 1000; |   |   |   | Create Bar Chart of Customer's MAX_CC_SPENT_AMOL    Finished D X    @<br>Saal<br>Visualize the Number of Customers with Max Credit Card Amount<br>SELECT * FROM CREDIT_SCORING_NEW FROM ROMAWH < 1000; | Compare the Classification models FINISHED D 2010 00 00 00 00 00 00 00 00 00 00 00 00 |
|      | CUSTOMER_ID  | ~ AGE   | ~ INCOME  | ✓ MARITAL ST =  | Grouped      Stacked     Good Credit     Other Credit  |   |
| lels | 85336  | 19  | 2150  | Single 🔺  | 35   | 0.8   |
|      | 55140  | 36  | 3750  | Single  | 30   |   |
|      | 8446   | 51  | 6250  | Married   | 25   |   |
|      | 32627  | 38  | 6250  | Married   | 20   | 0.4   |
|      | 23786  | 47  | 2150  | Single  |  |   |
|      | 18646  | 36  | 5250  | Married   | in the second  | 0.2   |
|      | 13646  | 49  | 4250  | Married   |  | 0.051   |



# **Oracle Machine Learning**

### Machine Learning Notebooks included in Autonomous Databases

Key Features:

- Collaborative UI for data scientist and analysts
- Packaged with Autonomous Databases
- Quick start Example notebooks
- Easy access to shared notebooks, templates, permissions, scheduler, etc.
- OML4SQL
- OML4Py coming soon
- Supports deployment of OML models





# Oracle Machine Learning for R / Python

## Multiple Components/APIs of Oracle Machine Learning

# Transparency layer

- Leverage proxy objects so data remain in database
- Overload native functions translating functionality to SQL
- Use familiar R/Python syntax to manipulate database data

# Parallel, distributed algorithms

- Scalability and performance
- Exposes in-database algorithms available from OML4SQL

# Embedded execution

- Manage and invoke R or Python scripts in Oracle Database
- Data-parallel, task-parallel, and non-parallel execution
- Use open source packages to augment functionality

# OML4Py, Automated Machine Learning - AutoML

- Feature selection, model selection, hyper-parameter tuning



\* Coming soon

### **Oracle Machine Leaning**

Multiple Languages UIs Supported for End Users & Apps Development







"Citizen" Data Scientists







# *Coming Soon!* | AutoML – *new* with OML4Py

Increase data scientist productivity – reduce overall compute time



Auto Algorithm Selection

- Identify in-database algorithm that achieves highest model quality
- Find best algorithm faster than with exhaustive search

# Auto Feature Selection Auto Tune Hyperparameters

- Reduce # of features by identifying most predictive
- Improve performance and accuracy

- Significantly improve model accuracy
- Avoid manual or exhaustive search techniques

Enables non-expert users to leverage Machine Learning

# Coming Soon! | OML AutoML User Interface

### "Code-free" user interface supporting automated end-to-end machine learning

Automate production and deployment of ML models

- Enhance Data Scientist productivity and user-experience
- Enable non-expert users to leverage ML
- Unify model deployment and monitoring
- Support model management

### Features

- Minimal user input: data, target
- Model leaderboard
- Model deployment via REST
- Model monitoring
- Cognitive features for image and text



# Coming Soon! | OML AutoML User Interface

# "Code-free" user interface supporting automated end-to-end machine learning

Automate production and deployment of ML models

- Enhance Data Scientist productivity and user-experience
- Enable non-expert users to leverage ML
- Unify model deployment and monitoring
- Support model management

### Features

- Minimal user input: data, target
- Model leaderboard
- Model deployment via REST
- Model monitoring
- Cognitive features for image and text



# *Coming Soon!* | Algorithms for Database 20c

## **Two major new ML algorithms**

## Gradient Boosted Trees (XGBoost)

- Highly popular and powerful algorithm Kaggle winners
- Classification, regression, ranking, survival analysis MSET-SPRT
  - Multivariate State Estimation Technique Sequential Probability Ratio Test (MSET-SPRT)
  - Nonlinear, nonparametric anomaly detection algorithm designed to monitor critical processes.
  - Detects subtle anomalies while also producing minimal false alarms.
  - Calibrates expected behavior from historical normal operational sequence of monitored signals.
  - Re-implemented and sped up in-DB and based on original Oracle Labs algorithm





# Oracle Data Miner UI

# Drag and Drop, Workflows, Easy to Use UI for "Citizen Data Scientist"

Easy to use to define analytical methodologies that can be shared

SQL Developer Extension

Workflow API and generates SQL code for immediate deployment





C Secure https://adwc.uscom-west-1.oraclecloud.com/oml/tenants/OCID1.TENANCY.OC1..AAAAAAAAAFCUE47PQMRF4VIGNEEBGBCMMOY5R7X... Q 🕁 +

#### 0

# Example Templates

New Notebook

 $\leftarrow \rightarrow$ 

| nomaly Detection   | Association Rules   | Attribute Importance  | Classification Prediction M   |
|--|---|---|---|
| his notebook shows how to detect                               | Notebook to show the use of Assoc                                   | Notebook to identify key attributes                                 | Example notebook to predict custo                                     |
| uthor:   | Author:   | Author:   | Author:   |
| ate Added: 5/4/18 6:59 AM<br>ags: 'Anomaly Detection' 'Machine | Date Added: 5/4/18 6:59 AM<br>Tags: 'SQL' 'Associations' 'Rules' 'M | Date Added: 5/4/18 6:59 AM<br>Tags: 'SQL' 'Attribute Importance' 'K | Date Added: 5/4/18 6:59 AM<br>Tags: 'Classification' 'Prediction' 'De |
| ★ 1 Likes 🔍 24 🙀 0   | ★ 0 Likes 🔍 2 📭 0   | ★ 0 Likes 🔍 10 📭 0  | ★ 0 Likes 🔍 6 📑 1   |
|  |   |   |   |

#### Clustering

This notebook shows how to identi...

Author:

Date Added: 5/4/18 6:59 AM Tags: 'Clustering' 'K-Means' 'Expect...

★ 0 Likes 🔍 9 📑 0

### My First Notebook

Oracle Machine Learning example ...

Author:

Date Added: 5/4/18 6:59 AM Tags: 'SQL' 'Data' 'Graph'

★ 0 Likes 🔍 15 📑 0

#### Regression

This notebook shows how to predic...

Author:

Date Added: 5/4/18 6:59 AM Tags: 'Regression' 'SVM' 'GLM' 'Logi...

★ 0 Likes ⊕ 6 📭 0

#### **Statistical Function**

Oracle Machine Learning example ...

#### Author:

Date Added: 5/4/18 6:59 AM Tags: 'Statistics' 'ANOVA' 'T-test' 'F-...

★ 0 Likes 🔍 1 📭 0

### Manage and Analyze All Your Data



### **In-Database Machine Learning**



More Models **Better Models** 

Faster, More Secur

Less Cost

|  | No Need To Extract and<br>Move Data     | Data stays in Database                                  |
|--|---|---|
|  | Data Preparation and<br>Transformation  | Accelerated with<br>Automatic Data Prep                 |
|  | Data Mining and<br>Model Building       | SQL, R, Python<br>Oracle Data Miner UI<br>OML Notebooks |
|  | No Need to Transform<br>Production Data | Embedded Data<br>Preparation                            |
|  | Model Scoring                           | Accelerated Via<br>Exadata Database Machine             |
|  | Ready to Deploy!                        | Easy to repeat model by                                 |

Zero time required. No production impact.

No separate environment required. Much faster data prep. Data stays protected and secured.

**Oracle Data Miner and AutoML** greatly speed model building. Less skill required. No coding.

No need for second production instance.

Faster model validation

ding as often as needed

### **Oracle's Machine Learning & Adv. Analytics Algorithms**

## CLASSIFICATION

- Naïve Bayes
- Logistic Regression (GLM)
- Decision Tree
- Random Forest
- Neural Network
- Support Vector Machine
- Explicit Semantic Analysis

#### CLUSTERING

- Hierarchical K-Means
- Hierarchical O-Cluster
- Expectation Maximization (EM)

#### **ANOMALY DETECTION**

One-Class SVM

#### TIME SERIES

- State of the art forecasting using **Exponential Smoothing**
- Includes all popular models e.g. Holt-Winters with trends, seasons, irregularity, missing data

OAA includes support for Partitioned Models, Transactional, Unstructured, Geo-spatial, Graph data. etc,

#### REGRESSION



- Linear Model Generalized Linear Model
- Support Vector Machine (SVM)
- Stepwise Linear regression
- Neural Network
- LASSO \*

#### **ATTRIBUTE IMPORTANCE**

- Minimum Description Length
   <sup>AAAAAAA</sup>
   <sup>AAAAAAA</sup>
- Principal Comp Analysis (PCA)
- Unsupervised Pair-wise KL Div
- CUR decomposition for row & AI

#### **ASSOCIATION RULES**

• A priori/ market basket

#### **PREDICTIVE OUERIES**

• Predict, cluster, detect, features

#### **SOL ANALYTICS**

- SQL Windows, SQL Patterns,
- SQL Aggregates

#### **FEATURE EXTRACTION**

- Principal Comp Analysis (PCA)
- Non-negative Matrix Factorization
- Singular Value Decomposition (SVD)
- Explicit Semantic Analysis (ESA)

#### **TEXT MINING SUPPORT**



• Algorithms support text



- Tokenization and theme extraction
- Explicit Semantic Analysis (ESA) for document similarity

#### **STATISTICAL FUNCTIONS**

• Basic statistics: min, max, median, stdev, t-test, F-test, Pearson's, Chi-Sq, ANOVA, etc.

#### **R** PACKAGES

- CRAN R Algorithm Packages through Embedded R Execution
- Spark MLlib algorithm integration

#### **EXPORTABLE ML MODELS**

 REST APIs for deployment  $\bigcirc$ 



<sup>•</sup> OAA (Oracle Data Mining + Oracle R Enterprise) and ORAAH combined

### **Oracle's Machine Learning & Adv. Analytics Algorithms**

### **STATISTICAL FUNCTIONS**

- Descriptive statistics (e.g. <u>median</u>, <u>stdev</u>, <u>mode</u>, <u>sum</u>, etc.)
- Hypothesis testing (<u>t-test</u>, <u>F-test</u>, <u>Kolmogorov-</u> <u>Smirnov test</u>, <u>Mann Whitney</u> <u>test</u>, <u>Wilcoxon Signed Ranks test</u>

 <u>Correlations analysis</u> (parametric and nonparametric e.g. <u>Pearson's test for</u> <u>correlation, Spearman's rho</u> <u>coefficient, Kendall's tau-b</u> <u>correlation coefficient</u>)

- <u>Ranking functions</u>
- <u>Cross Tabulations with Chi-square</u> <u>statistics</u>



- Linear regression
- <u>ANOVA</u> (Analysis of variance)
- Test Distribution fit (e.g. <u>Normal distribution</u> <u>test</u>, <u>Binomial test</u>, <u>Weibull</u> <u>test</u>, <u>Uniform</u> <u>test</u>, <u>Exponential</u> <u>test</u>, <u>Poisson test</u>, etc.)
- <u>Statistical Aggregates</u> (min, max, mean, median, stdev, mode, quantiles, plus x sigma, minus x sigma, top n outliers, bottom n outliers)

### **ANALYTICAL SQL**

- <u>SQL Windows</u>
- <u>SQL Aggregate functions</u>
- LAG/LEAD functions
- SQL for Pattern Matching
- Additional approximate query processing: APPROX\_COUNT , APPROX\_SUM, APPROX\_RANK
- <u>Regular Expressions</u>

### ML and Al are just "Algorithms"

Algorithms Operate on Data







Posterior Probability

Predictor Prior Probability

 $P(c \mid \mathbf{X}) = P(x_1 \mid c) \times P(x_2 \mid c) \times \dots \times P(x_n \mid c) \times P(c)$ 

Move the Algorithms; Not the Data!; It Changes *Everything*!

### **Thank You**

## **Any Questions ?**

Sandesh Rao VP AIOps for the Autonomous Database



<u>@sandeshr</u> in <a href="https://www.linkedin.com/in/raosandesh/">https://www.linkedin.com/in/raosandesh/</a> https://www.slideshare.net/SandeshRao4



