# ML in PostgreSQL
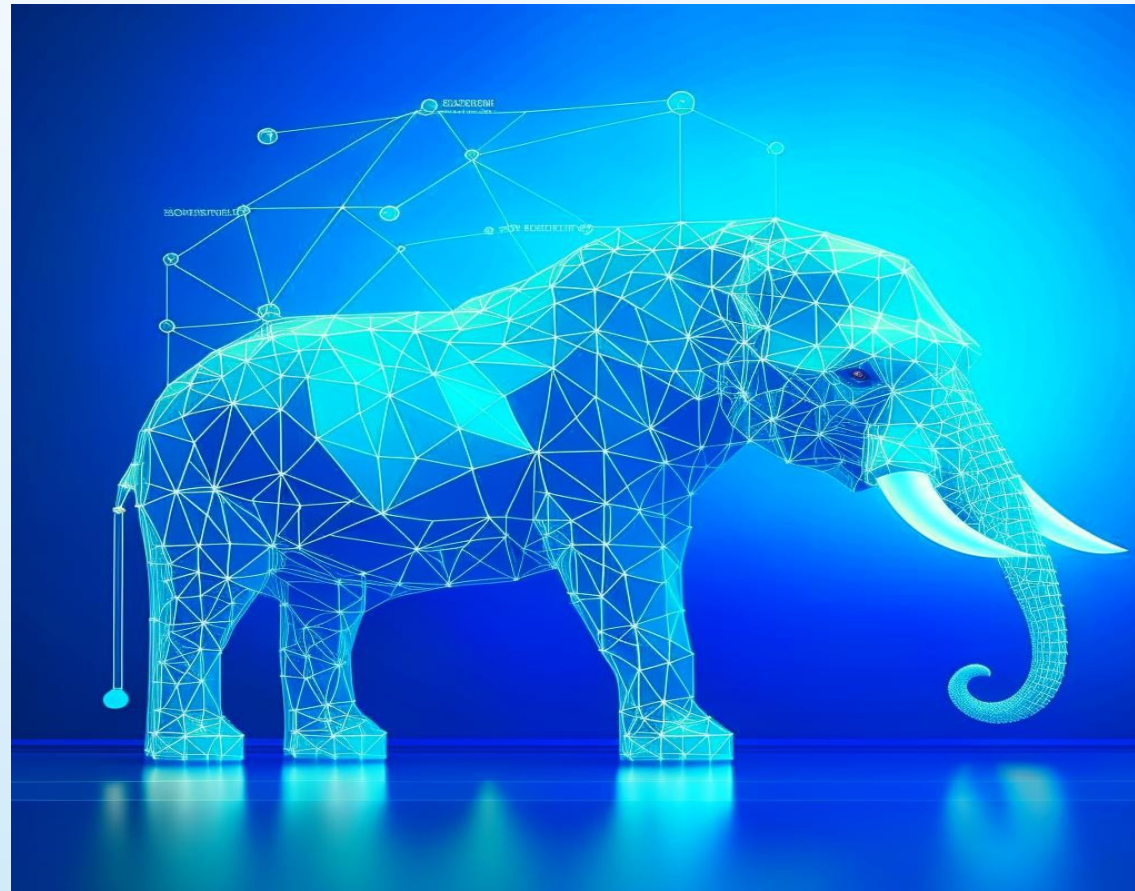
# Using models:

- Binary Classification

- Multi Classification

- Regression

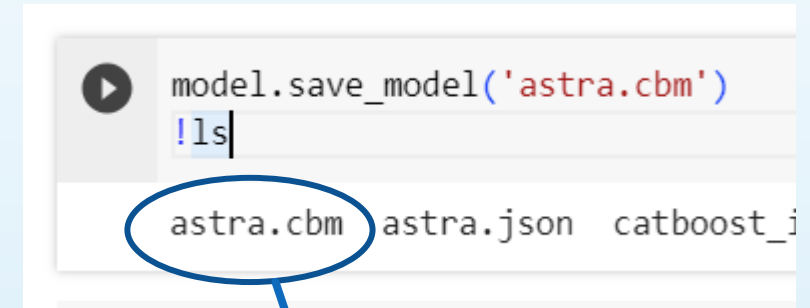# Using framework:

CatBoost: https://catboost.ai/

only prediction

# ML process:

1. Training model
2. Save model to database server

3. prediction



```
model.save_model('astra.cbm')
!ls
```

astra.cbm   astra.json   catboost_

```
adult=# SELECT ml_predict('astra3.cbm','astra3');
WARNING:   field run_id not used
WARNING:   field field_id not used
WARNING:   field spec_obj_id not used
WARNING:   field predict not used
      ml_predict
-----------------------
 public.astra3_predict
(1 row)
```

# Installation

- git clone https://github.com/akalend/pg_ml.git
- export PG_HOME=/usr/local/pgsql     //where is main postgres  folder
- wget https://github.com/catboost/catboost/releases/download/v1.2.2/libcatboostmodel.so
- mv libcatboostmodel.so $PG_HOME/lib
- cd pg_ml
- export PG_CONFIG=$PG_HOME/bin/pg_config
- export LD_LIBRARY_PATH=$PG_HOME/lib
- USE_PGXS=1 make
- sudo su
- export PATH=$PATH:$PG_HOME/bin
- USE_PGXS=1 make install
- chown postgres model.cbm
- [optional] cp model.cbm $PG_HOME/data

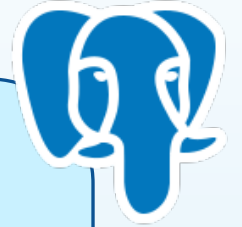# DataFrame columns convert to PostgreSQL fields

ID, Id, id

FieldID, Field_ID

Import-port

Port, PORT

<=>

id

field_id

import_port

port

# DataFrame columns and PostgreSQL fields

```
[14] df = pd.read_csv('star_classification.csv')

[16] for it in df.columns:
         print(it)

     obj_ID
     alpha
     delta
     u
     g
     r
     i
     z
     run_ID
     rerun_ID
     cam_col
     field_ID
     spec_obj_ID
     class
     redshift
     plate
     MJD
     fiber_ID
```

<=>

```
adult=# \d astra3
                        Table "public.astra
    Column      |        Type          | Collation
----------------+----------------------+----------
 alpha          | double precision     |
 delta          | double precision     |
 u              | double precision     |
 g              | double precision     |
 r              | double precision     |
 i              | double precision     |
 z              | double precision     |
 run_id         | bigint               |
 cam_col        | bigint               |
 field_id       | bigint               |
 spec_obj_id    | double precision     |
 redshift       | double precision     |
 plate          | bigint               |
 mjd            | bigint               |
 fiber_id       | bigint               |

adult=#
```

# DataFrame columns and PostgreSQL fields

```
df.head()
```

|   | obj_ID | alpha | delta | u | g | r | i | z | run_ID | rerun_ID | cam_col | field_ID | spec_obj_ID |
|---|--------|-------|-------|---|---|---|---|---|--------|----------|---------|----------|-------------|
| 0 | 1.237661e+18 | 135.689107 | 32.494632 | 23.87882 | 22.27530 | 20.39501 | 19.16573 | 18.79371 | 3606 | 301 | 2 | 79 | 6.543777e+18 |
| 1 | 1.237665e+18 | 144.826101 | 31.274185 | 24.77759 | 22.83188 | 22.58444 | 21.16812 | 21.61427 | 4518 | 301 | 5 | 119 | 1.176014e+19 |
| 2 | 1.237661e+18 | 142.188790 | 35.582444 | 25.26307 | 22.66389 | 20.60976 | 19.34857 | 18.94827 | 3606 | 301 | 2 | 120 | 5.152200e+18 |
| 3 | 1.237663e+18 | 338.741038 | -0.402828 | 22.13682 | 23.77656 | 21.61162 | 20.50454 | 19.25010 | 4192 | 301 | 3 | 214 | 1.030107e+19 |
| 4 | 1.237680e+18 | 345.282593 | 21.183866 | 19.43718 | 17.58028 | 16.49747 | 15.97711 | 15.54461 | 8102 | 301 | 3 | 137 | 6.891865e+18 |

```
adult=#  select * from astra3 limit 3;
     alpha       |      delta       |    u     |    g     |    r     |    i     |    z     | run_id | cam_col | field_id |     spec_obj_id
-----------------+------------------+----------+----------+----------+----------+----------+--------+---------+----------+-----------------------
 16.95688978450004 | 3.64613008870454 | 23.33542 | 21.95143 | 20.48149 |   19.603 | 19.13094 |   7712 |       6 |      442 | 4.855016555329904e+18 |
 240.063240247767 | 6.1341305981397 3 | 17.86033 | 16.79228 | 16.43001 | 16.30923 | 16.25873 |   3894 |       1 |      243 | 2.4489280322708705e+18 |
 30.887222067625 | 1.18870964120799 | 18.18911 | 16.89469 | 16.42161 | 16.24627 | 16.18549 |   7717 |       1 |      536 | 8.255357438959835e+18 |
(3 rows)
```

# Information about model

```
adult=# SELECT ml_info ('astra3.cbm');
                                 ml_info
------------------------------------------------------------------
 dimension:3 numeric features:12 categorial features:0 modelType "MultiClass"+
 fieldName:alpha,delta,u,g,r,i,z,cam_col,redshift,plate,MJD,fiber_ID
(1 row)

adult=#
```

- Dimension result (How many classes)
- Feature count (categorical and float)
- Type of model
- Fields name

# Information about model

```
adult=# SELECT ml_info ('astra3.cbm');
                                ml_info
-------------------------------------------------------------------------
 dimension:3 numeric features:12 categorial features:0 modelType "MultiClass"+
 fieldName:alpha,delta,u,g,r,i,z,cam_col,redshift,plate,MJD,fiber_ID
(1 row)

adult=#
```

```
adult=# SELECT ml_info ('titanic.cbm');
                                ml_info
-------------------------------------------------------------------------
 dimension:1 numeric features:2 categorial features:9 modelType "Accuracy"     +
 fieldName:PassengerId,Pclass,Name,Sex,Age,SibSp,Parch,Ticket,Fare,Cabin,Embarked
(1 row)
```

# Model information

# More information

```
model.save_model('astra.json', format='json')
!ls

astra.cbm  astra.json  catboost_info  sample_data
```

```
adult=# SELECT * from  ml_json_info('astra.json');
                       ml_json_info
---------------------------------------------------------------------
 float feature:alpha,delta,u,g,r,i,z,cam_col,redshift,plate,MJD,fiber_ID,+
 categorial feature:
(1 row)
```

# Categorical Feature list
# Float Feature count list

# Prediction of model (recordset)

**Path to model file**

**Table name**

**Categorical field list**

```
adult=#    SELECT * from ml_cat_predict ('titanic.cbm',
'titanic','{name,passenger_id,pclass,sex,sibsp,parch,ticke
t,cabin,embarked }');
 row_num |        predict         | class
---------+------------------------+------
       0 |     -1.7937342449233795 |   0
       1 |     -0.7958399022225136 |   0
       2 |     -2.3928732160132470 |   0
       3 |     -1.9429766248990040 |   0
       4 |     -0.41747860726736713 |   0
       5 |     -2.0608914711097546 |   0
       6 |      0.5914467057444344 |   1
       7 |     -1.0786526230973736 |   0
       8 |      0.6757411102494171 |   1
       9 |     -3.2509569289807160 |   0
      10 |     -2.2747255881045620 |   0
      11 |     -1.3228896775643357 |   0
      12 |      2.7093190924641700 |   1
      13 |     -2.4233542239140187 |   0
```

# Prediction of model (result table)

**Path to model file**

**Table name**

```
adult=# SELECT * from  ml_predict_table('astra3.cbm','astra3');
   ml_predict_table
--------------------------
 public.astra3_predict
(1 row)
```

**Create the new table**

# Prediction of model (table)

**Path to model file**

**Table name**

**List of categorical fields**

```
adult-# ^C
adult=# SELECT ml_predict ('adult.cbm',   'adult2',
adult(# '{workclass,education,marital_status, occupation,relationship,race,sex,native_country}');
      ml_predict
----------------------
 public.adult2_predict
(1 row)
```

**Create the new table**

# Prediction results

## SELECT * FROM {table}_predict;



## SELECT * FROM ml_predict(...);

# Binary classification

# Binary classification

# Regression

# Multi classification

# Inner data model