02 Statistical modeling

Author: Miao Cai miao.cai@slu.edu

Statistical modeling

We then use four different models to model the risk during the trip:

- · Logistic regression
- Poisson regression
- XGBoost
- · Deep learning (Neural networks)

import packages and read data

```
1 # !pip install h2o
2 import sys
3 import numpy as np
4 import h2o
5 from h2o.estimators.glm import H2OGeneralizedLinearEstimator
6 h2o.init(nthreads = -1, max_mem_size = 8)
7
8 print("Python version: " + sys.version)
9 print("numpy version:", np.__version__)
10 print("h2o version:", h2o.__version__)
```

```
Checking whether there is an H2O instance running at <a href="http://localhost:54321">http://localhost:54321</a> ..... not fo
    Attempting to start a local H2O server...
      Java Version: openjdk version "11.0.4" 2019-07-16; OpenJDK Runtime Environment (build
      Starting server from /usr/local/lib/python3.6/dist-packages/h2o/backend/bin/h2o.jar
      Ice root: /tmp/tmp294 9azi
      JVM stdout: /tmp/tmp294_9azi/h2o_unknownUser_started_from_python.out
      JVM stderr: /tmp/tmp294 9azi/h2o unknownUser started from python.err
      Server is running at <a href="http://127.0.0.1:54321">http://127.0.0.1:54321</a>
    Connecting to H2O server at <a href="http://127.0.0.1:54321">http://127.0.0.1:54321</a> ... successful.
    H20
    cluster
                02 secs
    uptime:
    H2O
    cluster
                Etc/UTC
    timezone:
    H2O data
    parsing
                UTC
    timezone:
    H20
    cluster
                3.26.0.11
    version:
    H20
    cluster
                11 days
    version
    age:
    H20
    cluster
                H2O from python unknownUser rtskve
    name:
    H20
    cluster total 1
    nodes:
    H20
    cluster free 8 Gb
    memory.
1 df = h2o.import file('https://raw.githubusercontent.com/caimiao0714/optimization stats cas
2 df[df['y'] > 0,'y binary'] = 1
3 df[df['y'] == 0,'y_binary'] = 0
4 df['y binary'] = df['y binary'].asfactor()
5 df['log Distance'] = df['Distance'].log()
6 df.head(5)
    Parse progress:
                                                                                           100%
    C1 y Distance Precipitation
                                   Traffic y binary log Distance
      0 0
              1018
                              0 0.299886
                                                 0
                                                          6.9256
      10
               973
                              0 0.565617
                                                 0
                                                        6.88038
      20
              1021
                              0 0.414564
                                                 0
                                                        6.92854
      30
               998
                              0 0.559767
                                                 0
                                                        6.90575
      40
               985
                              0 0.777217
                                                 0
                                                        6.89264
```

¹ lk = h2o.import_file('https://raw.githubusercontent.com/caimiao0714/optimization_stats_cas

^{2 14} hoad(5)

Ľ→

4 TV.IICau()

Parse progress:						
C1 # Node A	Node Z	Distance Prec	ipitation Traffic			
0 Ann_Arbor Ithaca		800	0 0.254345			
1 Ann_Arbor Princeton		800	0 0.243435			
2 Ann_Arbor Salt_Lake_City		2400	0 0.254188			
3 Atlanta	Houston	1200	0 0.424037			
4 Atlanta	Pittsburgh	900	0 0.573477			

Split into train and test sets

Logistic regression

Coefficients: glm coefficients

	names	coefficients	standardized_coefficients
0	Intercept	-3.604438	-2.144072
1	Distance	0.001008	0.032266
2	Precipitation	0.256380	0.092045
3	Traffic	0.830543	0.187609

```
1 print("Logistic regression model evaluation:")
 2 print("train AUC: " + str(fit_logit.auc()))
 3 print("test AUC: " + str(logit_test_fit.auc()))
 4 print("---")
 5 print("train Accuracy" + str(fit_logit.accuracy()))
 6 print("test Accuracy" + str(logit_test_fit.accuracy()))
 7 print("---")
 8 print("train MSE" + str(fit_logit.mse()))
 9 print("test MSE" + str(logit_test_fit.mse()))
10 print("---")
11 print("train R-square: " + str(fit_logit.r2()))
12 print("test R-square: " + str(logit_test_fit.r2()))
    Logistic regression model evaluation:
    train AUC: 0.5596289078650459
    test AUC: 0.5638801871833545
    train Accuracy[[0.18502530292639074, 0.8936048995869534]]
    test Accuracy[[0.17768208465318328, 0.8940620782726046]]
    train MSE0.09478002969662627
    test MSE0.09376565320196198
    train R-square: 0.004278462347631851
    test R-square: 0.004429388061521045
```

Poisson regression

names coefficients standardized coefficients

Coefficients: glm coefficients

glm Model Build progress:

100%

	names	COETTICIENTS	Standardized_Coefficients
0	Intercept	-4.371852	-2.102051
1	Distance	0.001747	0.055937
2	Precipitation	0.334264	0.120008
3	Traffic	0.947354	0.213995

→ XGBoost

```
1 from h2o.estimators import H2OXGBoostEstimator
 2 xgboost params = {
         "ntrees" : 50,
 3
         "max_depth" : 5,
 4
 5
         "learn_rate" : 0.001,
         "sample_rate" : 0.7,
 6
 7
         "col_sample_rate_per_tree" : 0.9,
         "min rows" : 5,
 8
         "seed": 4241,
 9
10
         "score tree interval": 10
11 }
12 fit_xgboost = H2OXGBoostEstimator(**xgboost_params)
13 fit xgboost.train(x = ['Precipitation', 'Traffic', 'Distance'],
14
                     y = 'y binary',
15
                     training frame = df train,
16
                     validation_frame = df_valid)
     xgboost Model Build progress: |
                                                                                  100%
 1 xgboost test fit = fit xgboost.model performance(df test)
 2 print("XGBoost regression model evaluation:")
 3 print("train AUC: " + str(fit_xgboost.auc()))
 4 print("test AUC: " + str(xgboost_test_fit.auc()))
 5 print("---")
 6 print("train Accuracy" + str(fit_xgboost.accuracy()))
 7 print("test Accuracy" + str(xgboost_test_fit.accuracy()))
 8 print("---")
 9 print("train MSE" + str(fit_xgboost.mse()))
10 print("test MSE" + str(xgboost_test_fit.mse()))
11 print("---")
12 print("train R-square: " + str(fit_xgboost.r2()))
13 print("test R-square: " + str(xgboost_test_fit.r2()))
```

```
XGBoost regression model evaluation:
    train AUC: 0.6058716756560456
    test AUC: 0.552575221410063
    ...
    train Accuracy[[0.48714303970336914, 0.8937473294402507]]
    test Accuracy[[0.4882924258708954, 0.8940620782726046]]
    ...
    train MSE0.2352391414490155
    test MSE0.2352131798112998
    ...
    train R-square: -1.4713294603237945
    test R-square: -1.4974105268199778
```

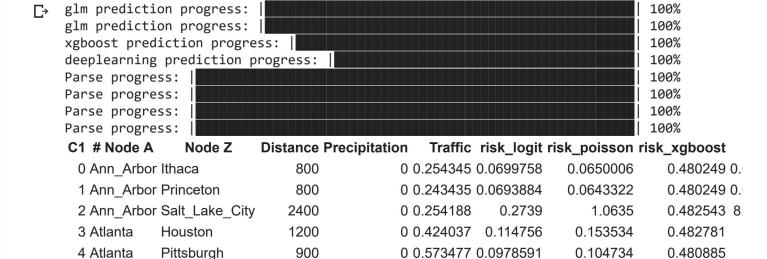
▼ Neural networks

C→

```
1 from h2o.estimators.deeplearning import H2OAutoEncoderEstimator, H2ODeepLearningEstimator
 1 fit_DL = H2ODeepLearningEstimator(epochs = 1000,
 2
                                     # hidden = [10, 10],
 3
                                     model_id = 'Deep learning',
 4
                                     seed = 1)
 5 fit DL.train(x = ['Precipitation', 'Traffic', 'Distance'],
               y = 'y binary',
 6
 7
               training_frame = df_train,
 8
                validation frame = df valid)
    deeplearning Model Build progress:
                                                                                  100%
 1 DL_test_fit = fit_DL.model_performance(df_test)
 2 print("Deep learning model evaluation:")
 3 print("train AUC: " + str(fit DL.auc()))
 4 print("test AUC: " + str(DL_test_fit.auc()))
 5 print("---")
 6 print("train Accuracy" + str(fit_DL.accuracy()))
 7 print("test Accuracy" + str(DL_test_fit.accuracy()))
 8 print("---")
 9 print("train MSE" + str(fit DL.mse()))
10 print("test MSE" + str(DL_test_fit.mse()))
11 print("---")
12 print("train R-square: " + str(fit_DL.r2()))
13 print("test R-square: " + str(DL_test_fit.r2()))
```

```
Deep learning model evaluation:
train AUC: 0.5614672763588284
test AUC: 0.5351742274819198
---
train Accuracy[[0.21823377000378666, 0.8933200398803589]]
test Accuracy[[0.2164417325947588, 0.8940620782726046]]
---
train MSE0.09477659936866129
test MSE0.09432099305838623
---
train R-square: 0.0043145000176650905
test R-square: -0.0014670145316892924
```

Prediction for links data



```
1 from google.colab import drive
2 from google.colab import files
3 drive.mount('/content/drive/', force_remount=True)
4 lk_risks.as_data_frame().to_csv('lk_risks.csv')
5 files.download('lk_risks.csv')
```

By Mounted at /content/drive/