# A comparison of causal discovery algorithms across different datasets

SUPERVISED BY:
V. CHRISTOPHIDES AND E. VAREILLE

PRESENTED BY:
ANIS AFLOU AND MARTIN GERVAIS

# team project ↘

Anis AFLOU

Martin GERVAIS

# content ↘

# context & motivation

# Families of Causal Discovery Algorithms



**OBSERVATIONAL DATA**

**CONSTRAINT BASED**

start with a fully connected undirected graph and then remove edges based on the result of conditional independence tests

- CAUSAL SUFFICIENCY
- CAUSAL FAITHFULNESS
- ACYCLICITY

**OUTPUT**

CPDAG

**SCORE BASED**

use a "goodness-of-fit" score of the model to the data while imposing a sparsity penalty to prevent overfitting

- CAUSAL SUFFICIENCY
- CAUSAL FAITHFULNESS
- ACYCLICITY

**OUTPUT**

DAG

**STRUCTURAL BASED**

assume an explicit functional causal model and identify causal directions by exploiting asymmetries in the noise distribution

- CAUSAL SUFFICIENCY
- CAUSAL FAITHFULNESS
- ACYCLICITY

**OUTPUT**

Fully Directed DAG

**FAMILY OF ALGORITHM**

**ASSUMPTIONS**

6

# Why are we comparing different **causal discovery algorithms?**

**Many causal discovery algorithms exist, but their practical behavior remains difficult to assess.**

**1**

Numerous causal discovery algorithms have been proposed

**2**

Algorithms are often tested on different datasets and with different metrics

**3**

However, evaluations and comparisons are still limited

**context ↗ & motivation**

# How do different causal discovery algorithms **behave** when **evaluated** under the same experimental conditions?

context ↗
& motivation

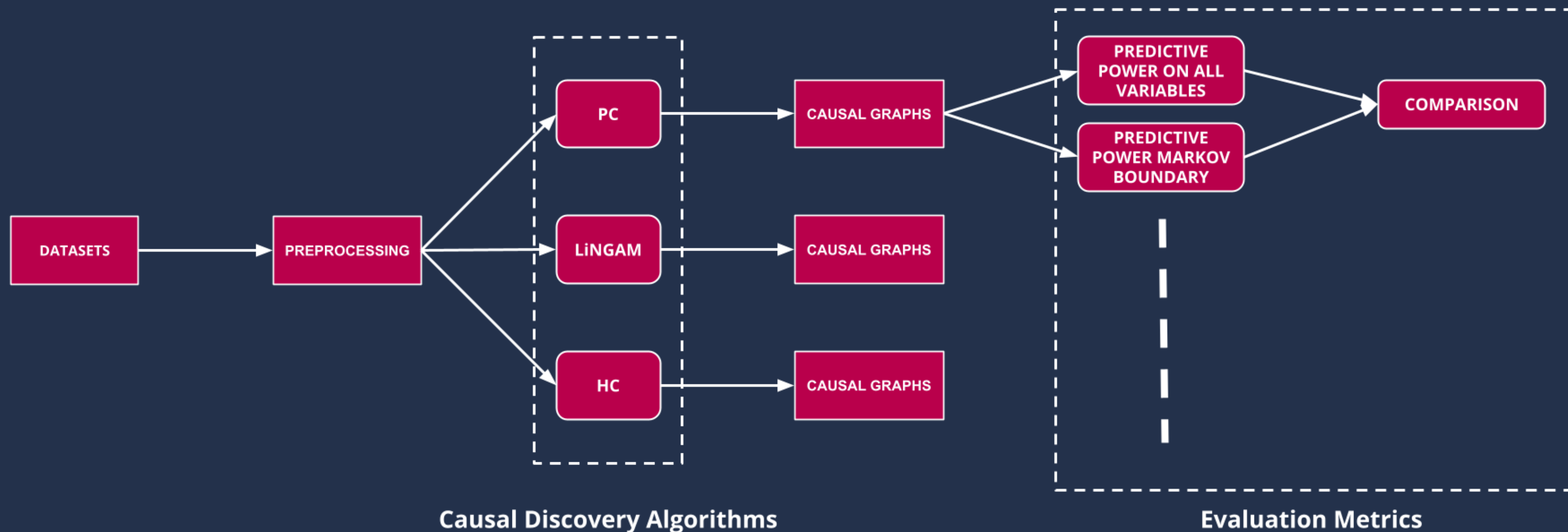# project pipeline

# project pipeline



figure 1: experimental pipeline

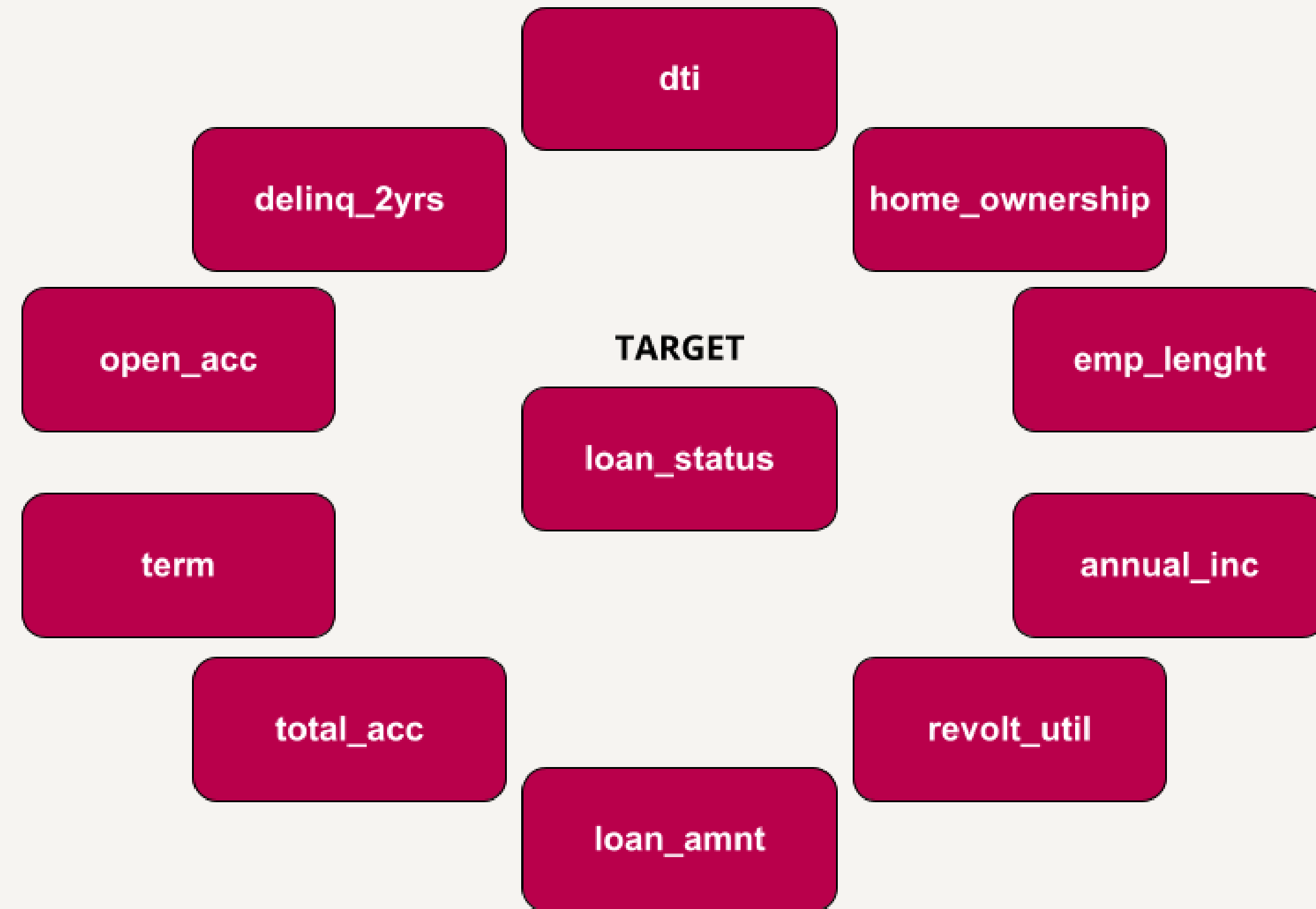# presentation ↘ of the datasets

DATASETS

ENSEA

# loan.csv

dti

delinq_2yrs

home_ownership

open_acc

TARGET

emp_lenght

loan_status

term

annual_inc

total_acc

revolt_util

loan_amnt

figure 2: list of the chosen features

**loan.csv**

**short description**
- 890000 loans
- 75 features

**short description**
- 20000 loans
- 11 features

Continuous 30

Discrete 25

Continuous 4

Categorical 4

Categorical 15

Temporal 5

Discrete 3

**presentation ↘
of the datasets**

figure 3: bubble charts of loan.csv

# GiveMeSomeCredit.csv



NumberOfTime[...]NotWorse

DebtRatio

NumberOfTime[...]NotWorse

NumberRealEstateLoansOrLines

RevolvingUtilizationOfUnsecuredLines

TARGET

Serious_Dlquin2yrs

NumberOfDependents

NumberOfTimes90DaysLate

MonthlyIncome

Age

NumberOfOpenCreditLinesAndLoans

figure 4: list of the chosen features

**short description**
- 150000 borrowers
- 11 features
- Mixed data

Discrete 7

Continuous 3

Categorical 1

figure 5: bubble chart for GiveMeSomeCredit.csv

**presentation ↘ of the datasets**

13

# OnlineShoppersIntention.csv



figure 6 : list of the chosen features



**short description**
- 12330 users sessions
- 18 features

figure 7 : bubble chart for OnlineShoppersIntention.csv

**presentation ↘
of the datasets**

# Algorithms used ⬊

PC

LiNGAM

HC

**Causal Discovery Algorithms**

# Algorithms used in this study ↘

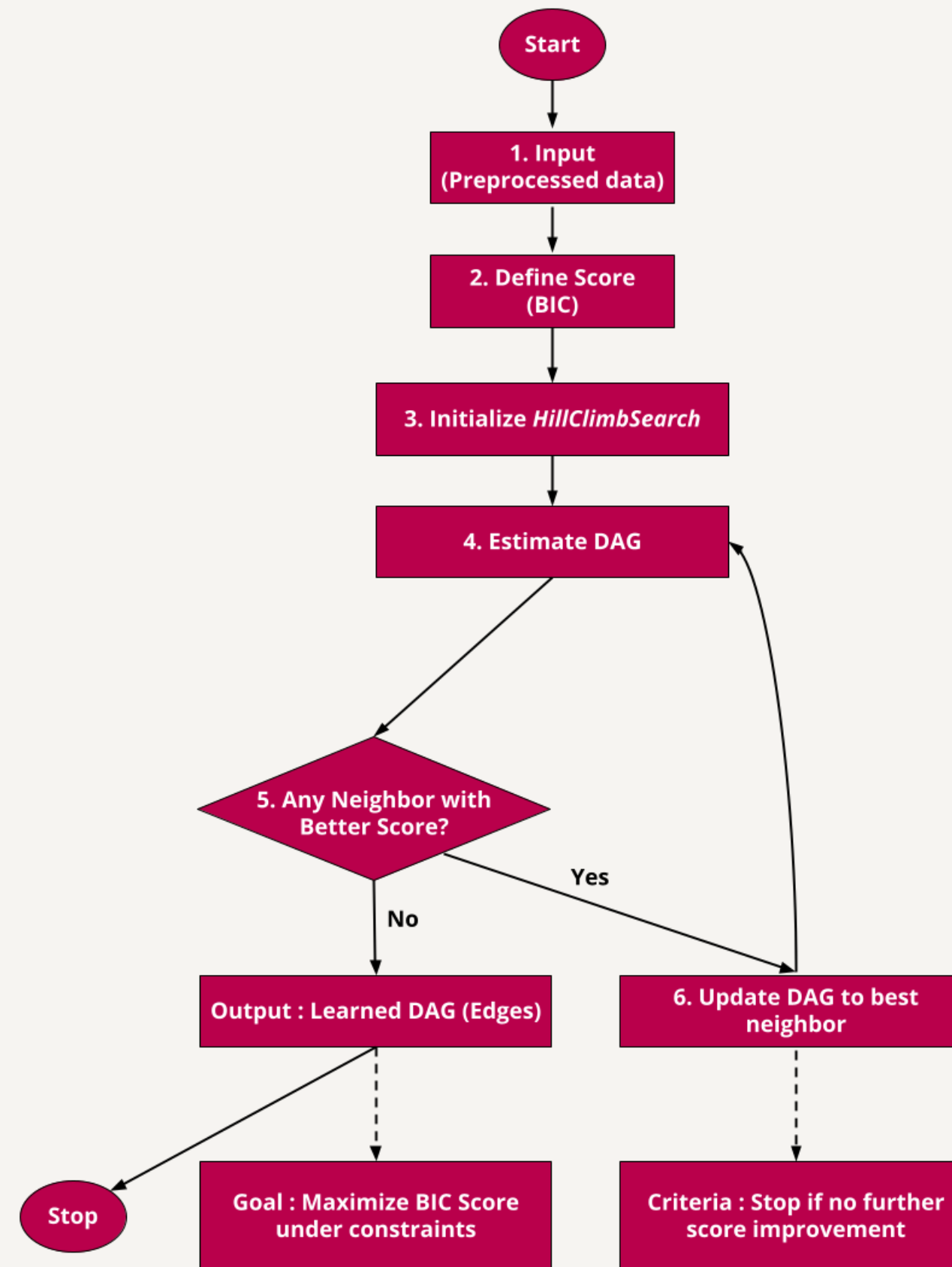| | Algorithm | Used on | Assumption | Parameters |
|---|---|---|---|---|
| 01 | PC Algorithm (Constraint Based) | Mixed datasets (continuous + categorical) | Causal Markov condition $(X \perp NonDesc(X) \mid Parents(X))$ | ci_test = 'pillai' significance_level=0.05 max_cond_vars = 3 njobs=–1 |
| 02 | Hill–Climbing (Score Based) | Discretized versions of the datasets | Causal Markov condition $(X \perp NonDesc(X) \mid Parents(X))$ | Scoring method = BIC |
| 03 | LiNGAM (Structural Based) | Continuous numerical variables | Independent and no gaussian noise | random_state=None, prior_knowledge=None, measure='pwling' |

# Flowchart of HillClimb Search Model ⬊



figure 7: flowchart of hillclimb search model
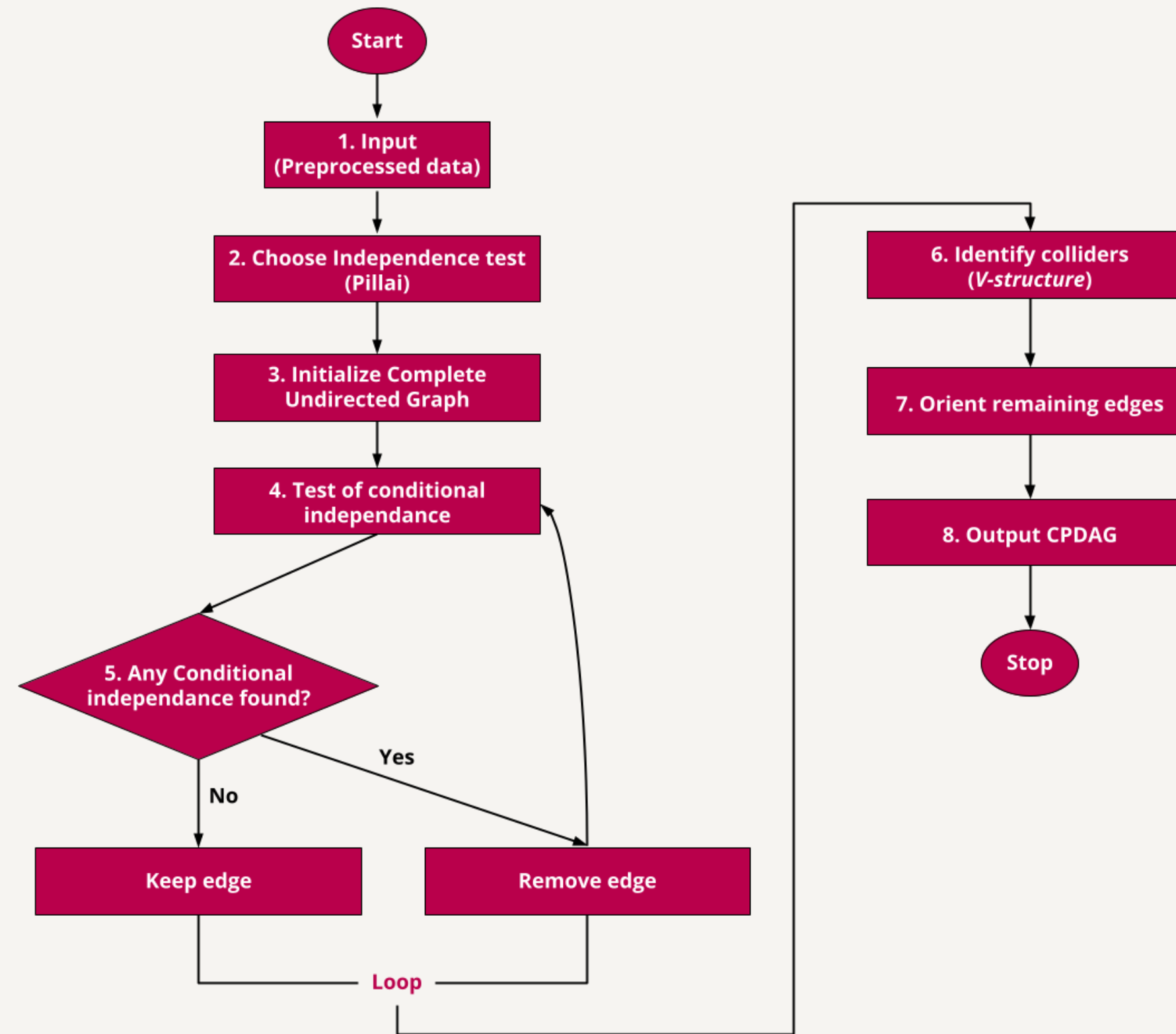
# Flowchart of PC Model ↘



figure 8: flowchart of PC model
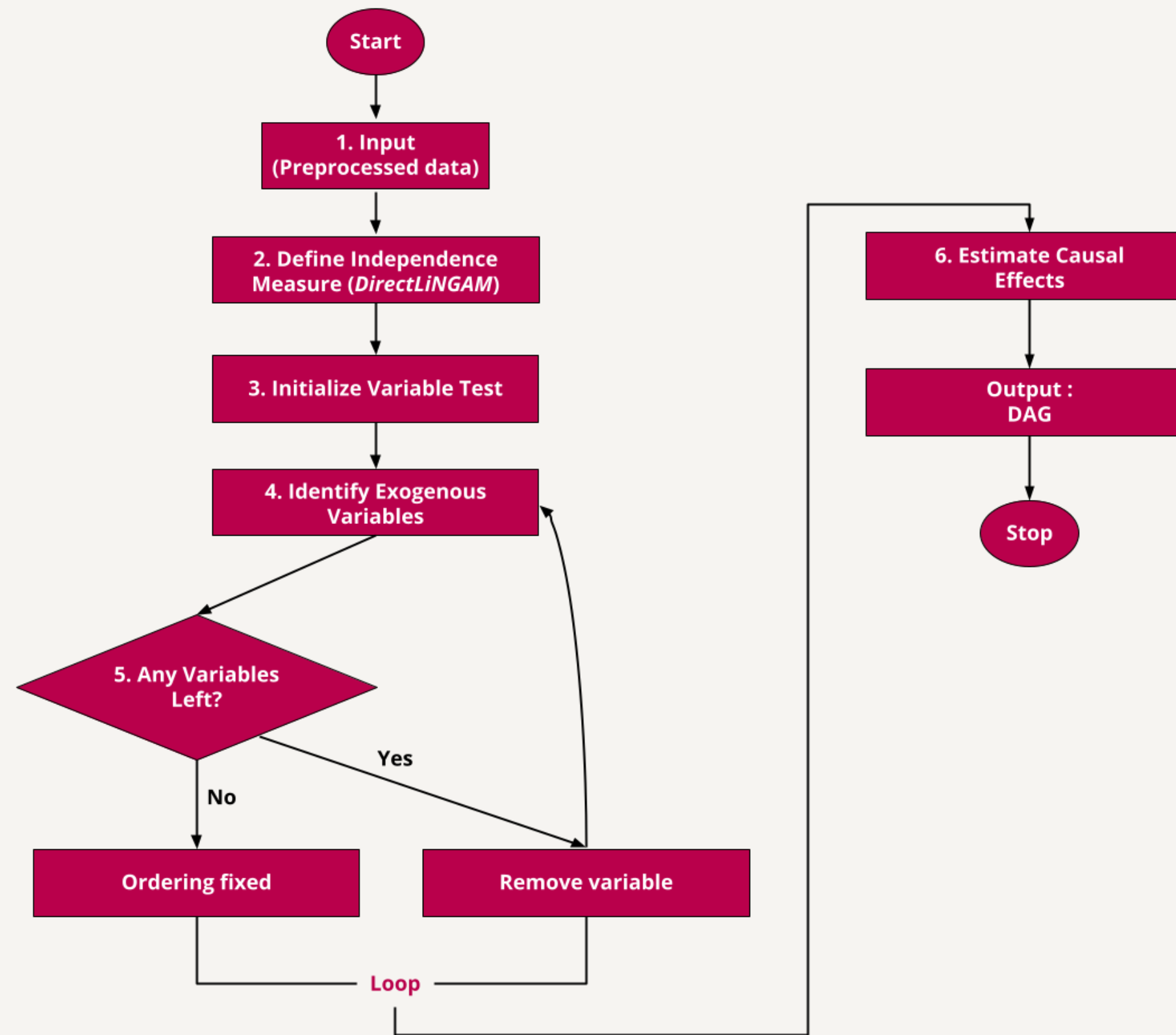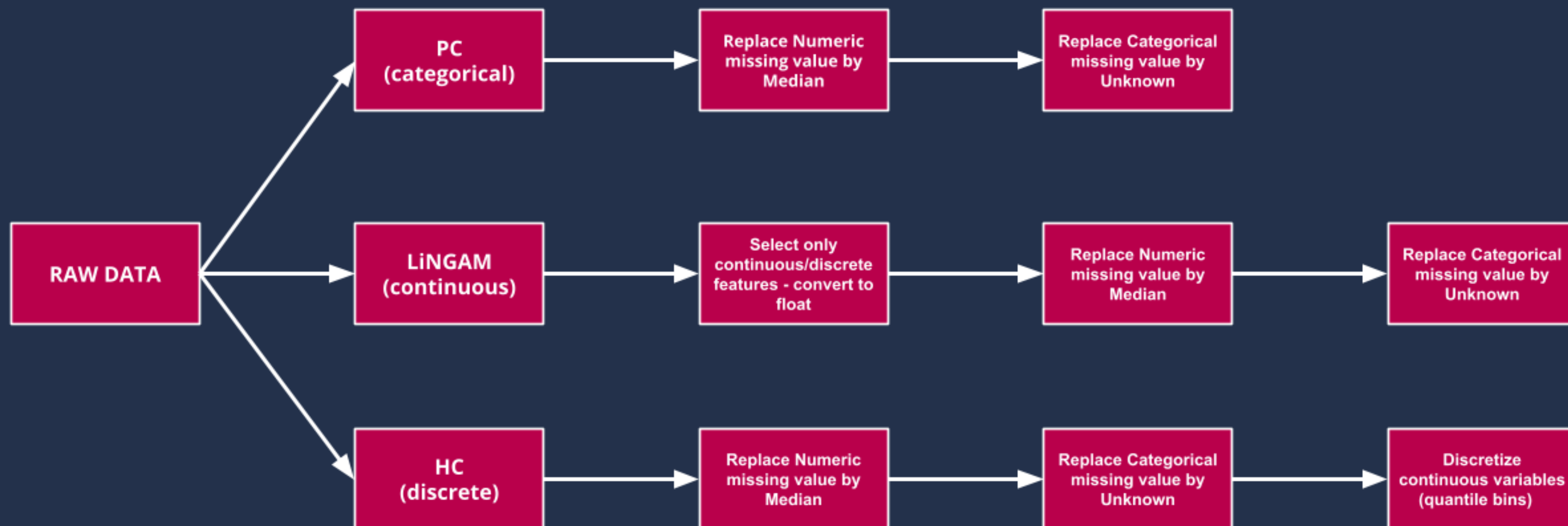
# Flowchart of LiNGAM Model



figure 9: flowchart of LiNGAM model

# data preprocessing

PREPROCESSING

# Preprocessing is adapted to each algorithm's assumptions

**RAW DATA** →

**PC (categorical)** → Replace Numeric missing value by Median → Replace Categorical missing value by Unknown

**LiNGAM (continuous)** → Select only continuous/discrete features - convert to float → Replace Numeric missing value by Median → Replace Categorical missing value by Unknown

**HC (discrete)** → Replace Numeric missing value by Median → Replace Categorical missing value by Unknown → Discretize continuous variables (quantile bins)

data processing

# Evaluation Metrics

# 1/ Predictive Power of Markov Blanket

- Extract the Markov Blanket of a target variable
- Train a predictive model using MB variables only
- Measure predictive performance

# 2/ Run time of compilation

**\*No SHD because we do not have the ground truth**
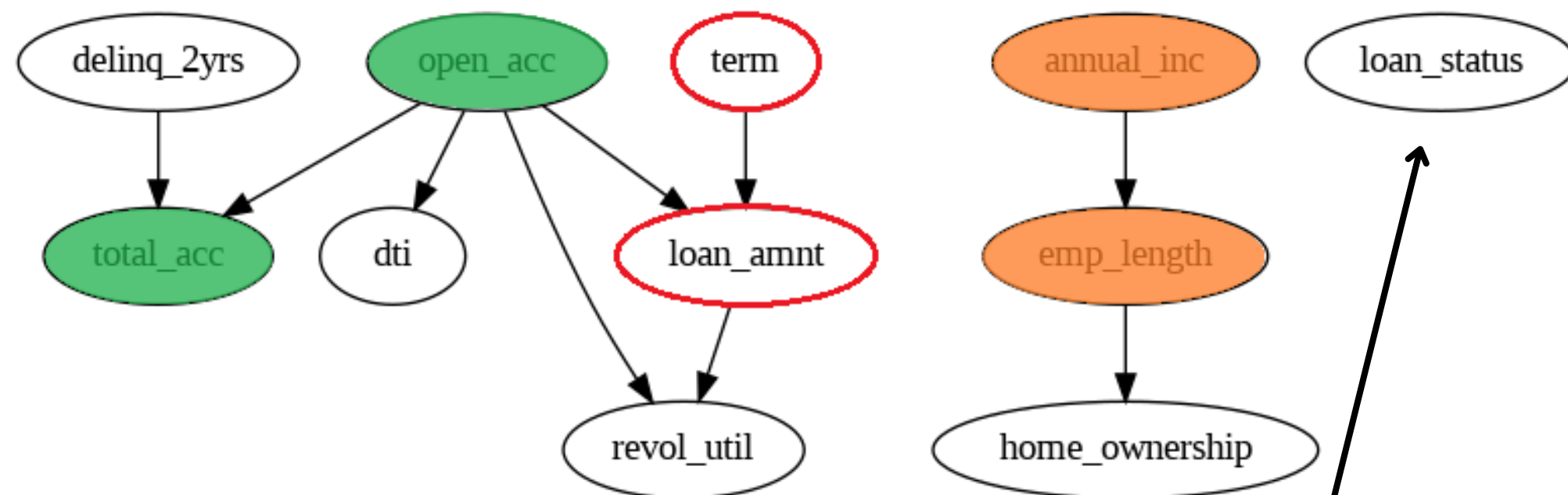
# Evaluation Metrics ↗

results 71

# 1st dataset : loan.csv
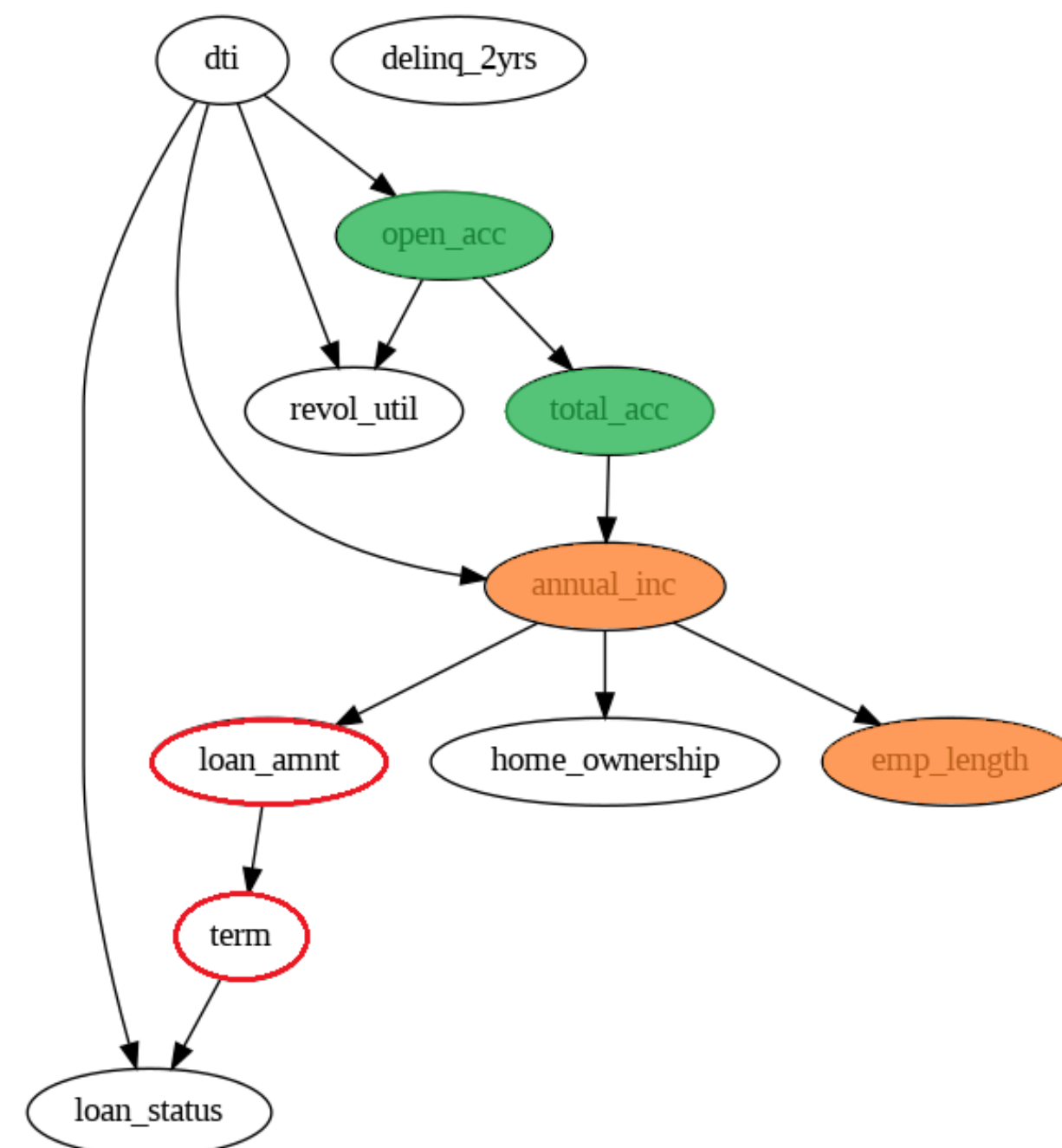
figure 10: causal graph on loan.csv with HC algorithm

# 2nd dataset : [GiveMeSomeCredit.csv](GiveMeSomeCredit.csv)

# GiveMeSomeCredit.csv PC/HC comparison

run time → <1s, 26.89it/s
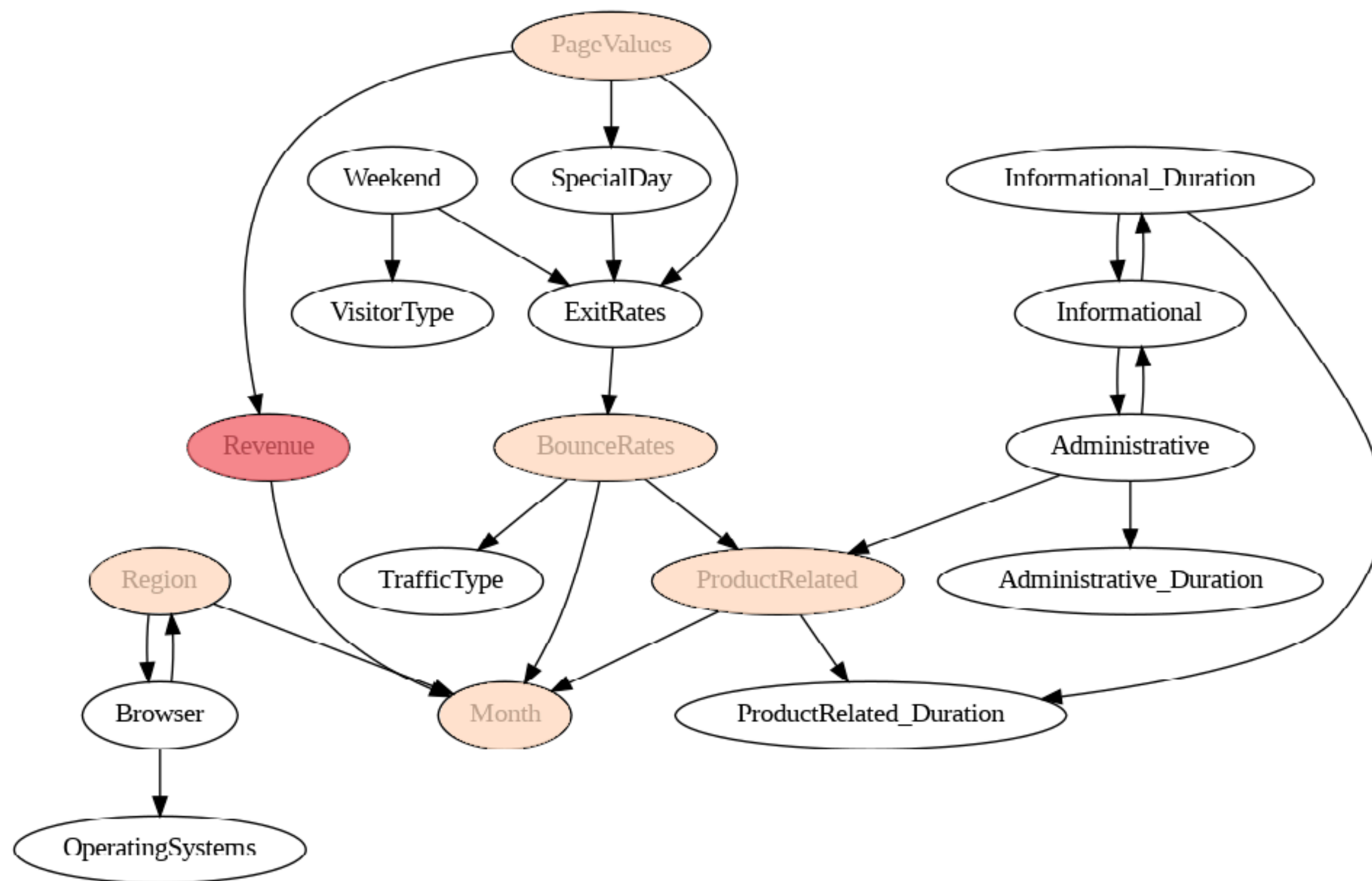
**HC algorithm**

run time → 05:01s, 60.25s/it

**PC algorithm**

●, ○ : same blocks

**results** ↗

# 3rd Dataset : OnlineShoppers.csv

run time → 30:47s, 397.29s/it



figure 12: causal graph on OnlineShopperscsv with PC algorithm
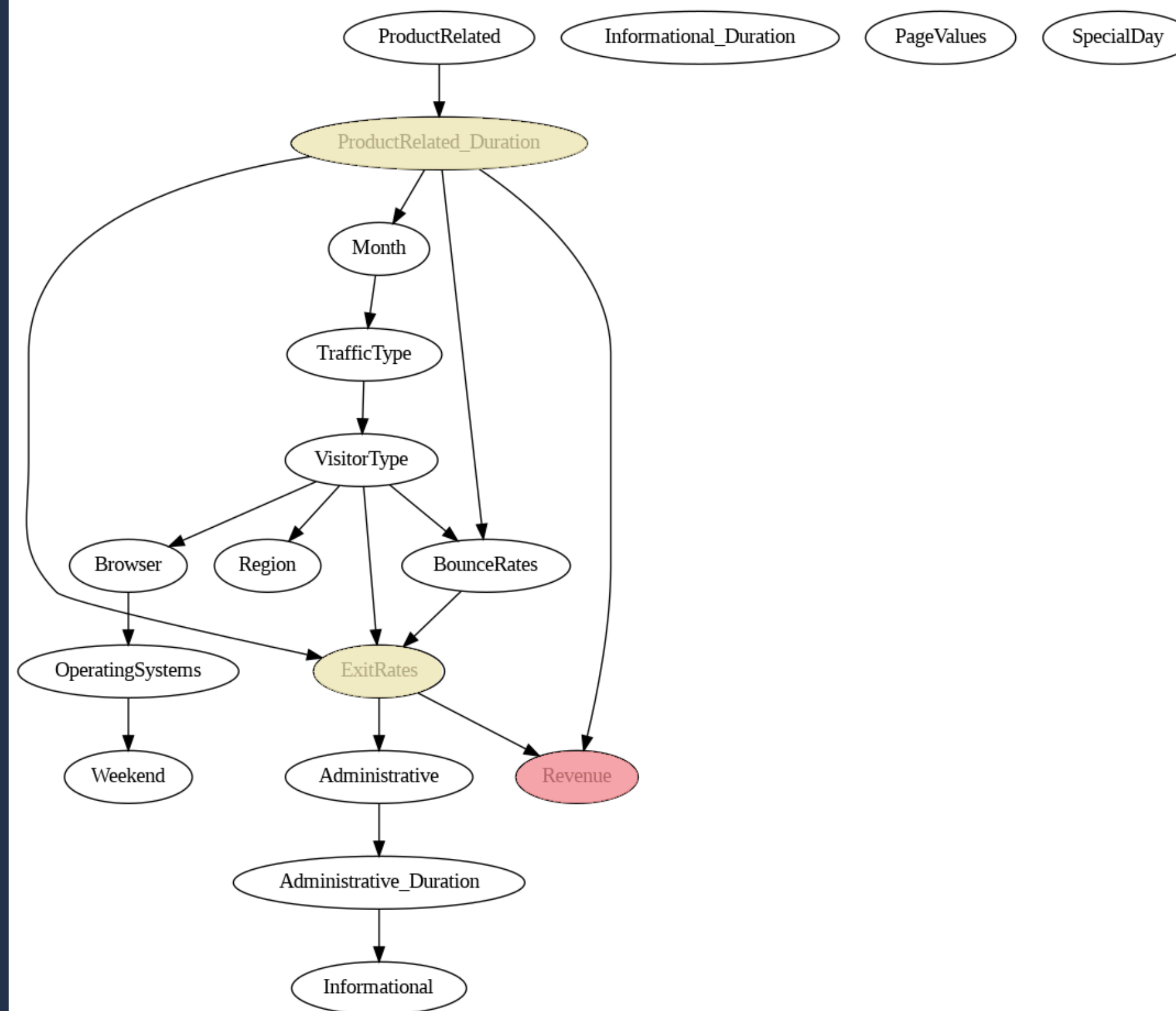
Markov Blanket

Target Feature

**5 features** on **Markov Blanket**
[PageValues, BounceRates,
ProductRelated, Month, Region]
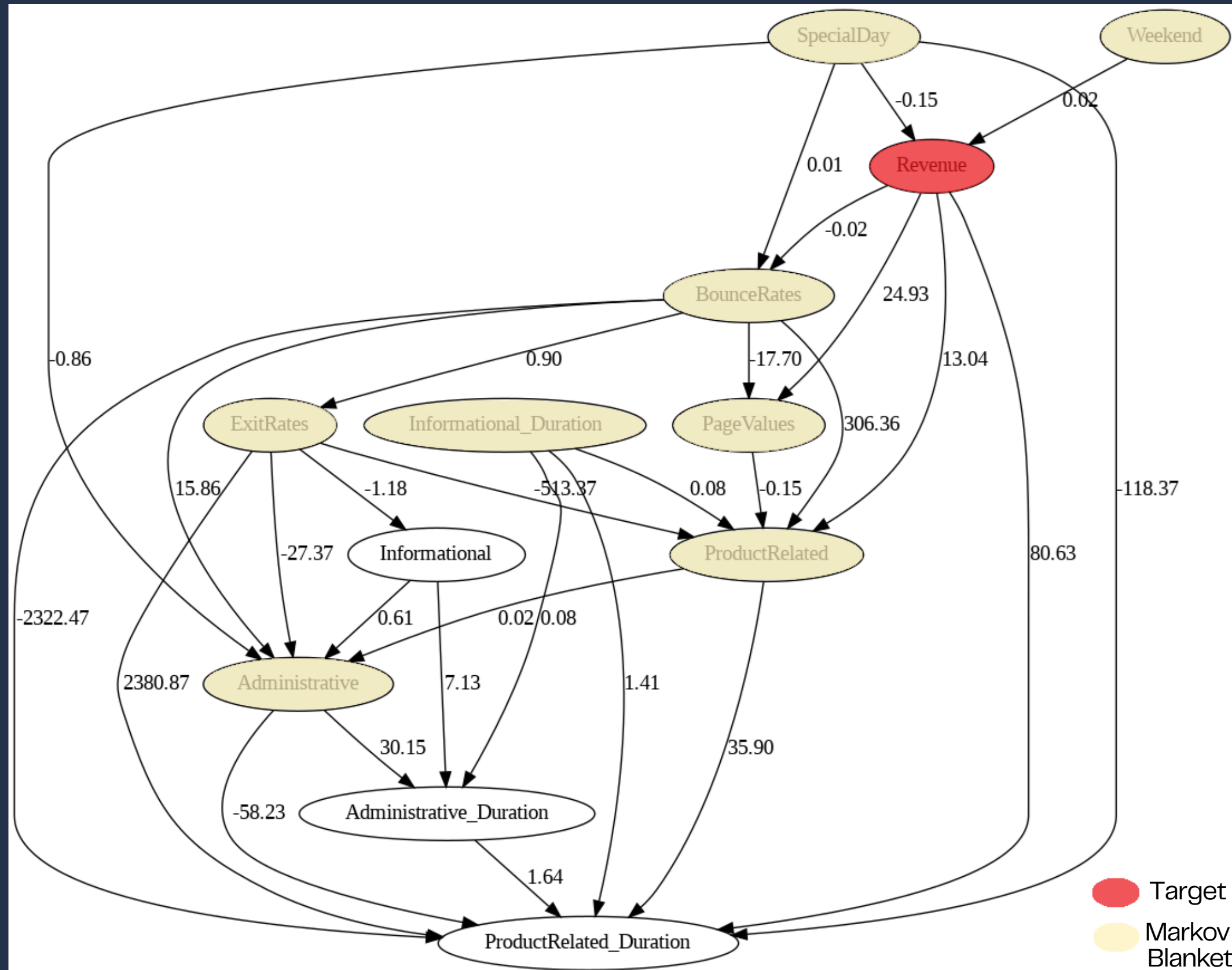
# OnlineShoppers.csv **HC** algorithm



figure 13 causal graph on OnlineShopperscsv with HC algorithm

run time → <1s, 26.21it/s

Markov Blanket

Target Feature

**2 features** on **Markov Blanket**
(too few)
[ExitRates,
ProductRelated_Duration]
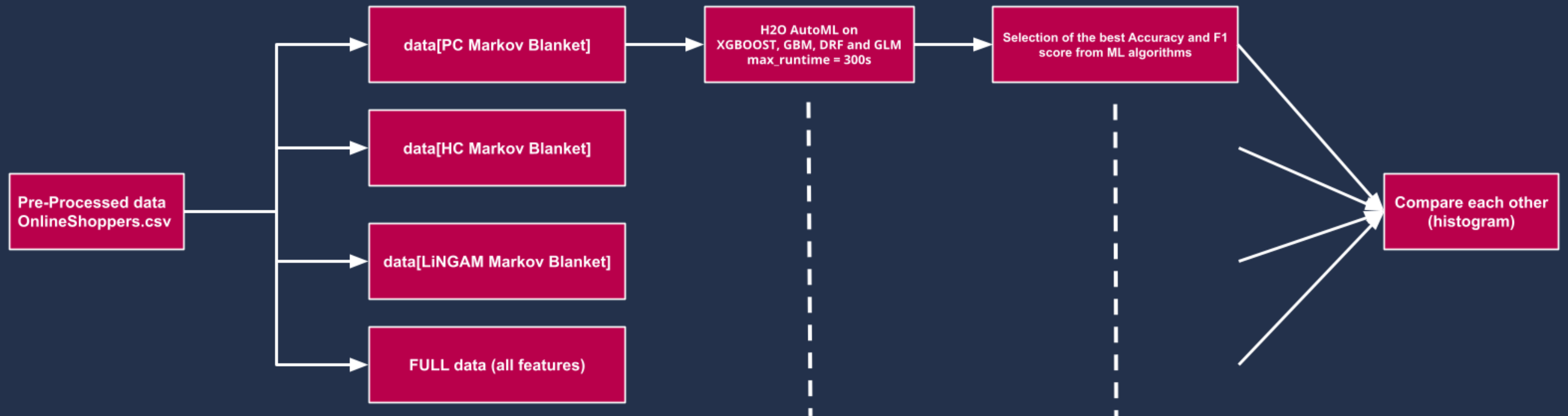
# OnlineShoppers.csv LiNGAM algorithm

run time → 00:15s , 24.11it/s

if X increases of 1 → Y increase by the value

Only on continuous/discrete variables

Settings set by default, because of low flexibility with LiNGAM

**8 features** on **Markov Blanket**
[SpecialDay, Weekend, BounceRates, PageValues, ProductRelated, Informational_Duration, ExitRates, Administrative]

figure 14 causal graph on OnlineShopperscsv with LiNGAM algorithm

# Comparaison des algorithmes PC/HC/LiNGAM

AutoML

**Pourquoi ?**

- Même condition pour évaluer la puissance prédicitive inter–algorithme

- Gain de temps



Pre-Processed data
OnlineShoppers.csv

data[PC Markov Blanket]

data[HC Markov Blanket]

data[LiNGAM Markov Blanket]

FULL data (all features)

H2O AutoML on
XGBOOST, GBM, DRF and GLM
max_runtime = 300s

Selection of the best Accuracy and F1
score from ML algorithms

Compare each other
(histogram)

results

Accuracy comparison by ML algorithm
FULL vs PC vs HC vs LiNGAM

| dataset | Accuracy |
|---|---|---|
| 0 | FULL | 0.900990 |
| 1 | MB_PC | 0.897140 |
| 2 | MB_HC | 0.847635 |
| 3 | MB_LiNGAM | 0.903740 |

figure 15 : AutoML performance comparison PC/HC

LiNGAM donne une meilleur accuracy max avec 8 features contre 18 dans le dataset original !

34

# Comparaison des algorithmes **PC et HC**



F1-score comparison by ML algorithm
FULL vs PC vs HC vs LiNGAM

On observe globalement une légère différence entre les F1–Scores du dataset entier par rapport aux autres

| | dataset | F1 |
|---|---|---|
| 0 | FULL | 0.668977 |
| 1 | MB_PC | 0.666667 |
| 2 | MB_HC | 0.380497 |
| 3 | MB_LiNGAM | 0.652241 |

figure 15 : AutoML performance comparison PC/HC

results

MB_UNION = MB_PC U MB_HC



AutoML performance comparison
FULL vs Markov Boundary UNION

From 18 to 7 features !

Résultat très
intéressant

figure 16 : AutoML performance comparison PC U HC

# Conclusion

## 1) Choix des algorithmes de découverte causale

| Algorithm | Family | Suitable when | Limitations |
|---|---|---|---|
| **Peter Clark** (PC) | Constraint-Based | • Moderate number of data<br>• Mixed data types | • High computational complexity |
| **Hill-Climbing** (HC) | Score-Based | • Trade-off between performance and computation time<br>• Mixed or discretized data | • Possible convergence to local optima |
| **LiNGAM** (L) | Structural-Based | • Mostly continuous variables<br>• Interest in causal effect strengh | • Non Gaussianity required |

# Perspectives : Découverte causale comme levier de réduction des coût ML

Sur des jeux de données massifs et bruités, l'entraînement de modèles ML devient extrêmement coûteux.

**Nous émettons l'hypothèse que :**

En identifiant des sous-ensembles  (ex. Markov Blanket), on réduirait drastiquement la complexité des modèles ML tout en maintenant, voire en améliorant, leurs performances prédictives.

**Cette approche pourrait :**

- Réduire le nombre de variables utilisées lors de l'entraînement
- Diminuer les coûts de calcul et de déploiement
- Renforcer l'interprétabilité des modèles

À plus long terme, cela ouvre la voie à des pipelines hybrides causaux–prédictifs, où la découverte causale agit comme un filtre, en amont de modèles ML complexes, combinant ainsi performance et interprétabilité.

# Thank you for your attention! ↘