

# Introduction to Casual Inference (DAG Framework) using pgmpy

Ankur Ankan

Postdoctoral Researcher  
Radboud University, The Netherlands

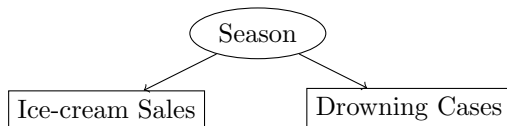
# Predictive vs Causal Modelling

## Predictive Modelling

- Exploits correlation among variables.
- Usually interested in predicting outcome variables.

## Causal Modelling

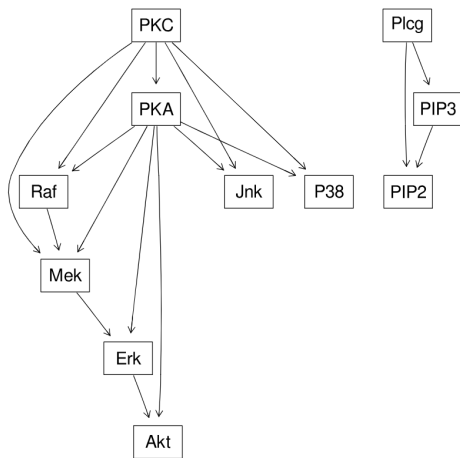
- Try to learn causal-effect relationships between variables.
- Interested in understanding the causal structure of the data generating mechanism and/or causal effect estimation.



Would an intervention on Ice-cream Sales affect Drowning Cases?

# Example: Protein Signalling Causal Network <sup>1</sup>

- Models the concentration of signalling proteins in cells.
- Understanding the mechanism such as how the signals bring cellular responses.
- Insights on how signalling pathways are altered in diseases.
- Can be used to identify potential targets for therapeutic intervention.



<sup>1</sup>Sachs, Karen, et al. "Causal protein-signaling networks derived from multiparameter single-cell data." Science 308.5721 (2005): 523-529.

# Examples

- **Epidemiology:** Analyzing data on how different treatment or exposures affect health outcomes.
- **Economics and Social Sciences:** Understanding impacts of policy interventions and help in designing potential policies.
- **Machine Learning:** For interpretability, feature selection, making models robust to out of distribution predictions.

# Two Main Frameworks

Mathematical frameworks for causal inference:

- **Potential Outcomes Framework** (a.k.a. Rubin's Casual Model)
  - ▶ Usually interested only in estimating causal effect.
  - ▶ Provides statistical methods for estimating the causal effect. Examples are Propensity Score based methods, Doubly robust estimators, etc.
- **Directed Acyclic Graphs (DAGs) / Structural Equation Models (SEMs)**
  - ▶ Causal diagram, i.e., DAG at the core.
  - ▶ The DAG can be used to define estimators for causal effects of interest.
  - ▶ DAGs make modelling assumptions explicit.

# Landscape of Causal Inference Python Packages

## Potential Outcomes Frameworks

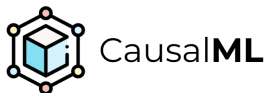
- Collection of estimation methods



- Doubly Robust Estimation



- Meta Learners: T/S/X-Learners.

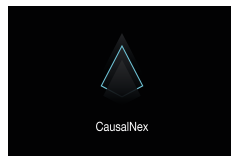


## DAG Frameworks

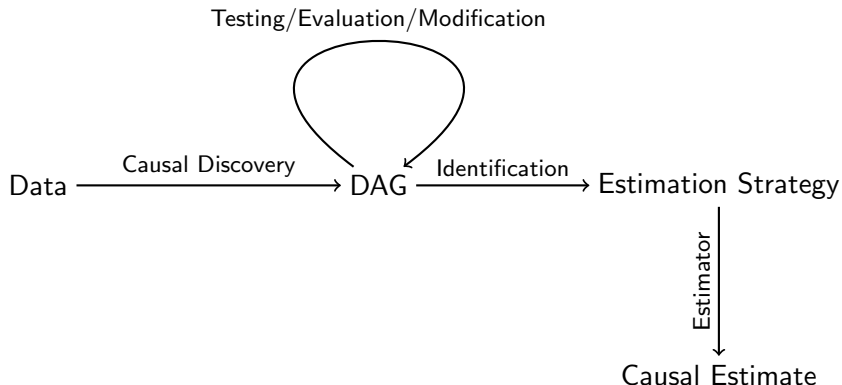
- pgmpy



- casualnex

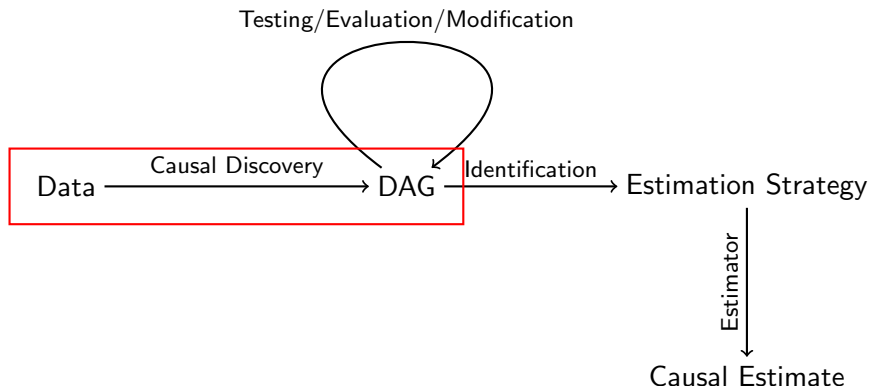


# A General Workflow in the DAG Framework



pgmpy provides functionality to perform each of these steps.

# Causal Discovery



**Causal Discovery:** Learn the causal DAG from data.

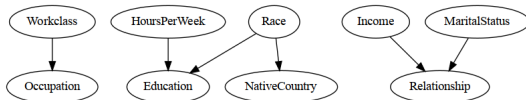


# Causal Discovery: Automated Algorithms

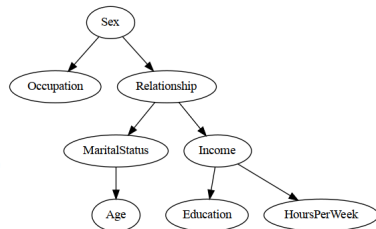
- Many automated algorithms with nice asymptotic properties.
  - ▶ Constraint-based: PC, Fast Causal Inference.
  - ▶ Score-based: Greedy Equivalence Search, Hill-Climb Search.
  - ▶ Optimization-based: NoTears, DAGMA.
- In practice, output varies significantly depending on sample size, algorithm used, and their hyperparameters.
- With no standard evaluation method, difficult to decide the correct model.

# Causal Discovery: Automated Algorithms

```
1. from pgmpy.estimators import PC
2.
3. est = PC(df)
4. dag = est.estimate(
5.     ci_test='chi_square')
```



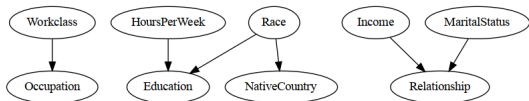
```
1. from pgmpy.estimators import HillClimbSearch
2.
3. est = HillClimbSearch(df)
4. dag = est.estimate(scoring_method='bicscore')
```



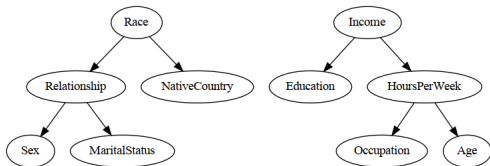
<sup>1</sup>Becker, Barry and Kohavi, Ronny. (1996). Adult. UCI Machine Learning Repository.  
<https://doi.org/10.24432/C5XW20>.

# Causal Discovery: Automated Algorithms

```
1. from pgmpy.estimators import PC
2.
3. est = PC(df)
4. dag = est.estimate(
5.     ci_test='chi_square')
```



```
1. from pgmpy.estimators import PC
2.
3. est = PC(df)
4. dag = est.estimate(
5.     ci_test='pillai')
```



# Causal Discovery: Problems Automated Algorithms

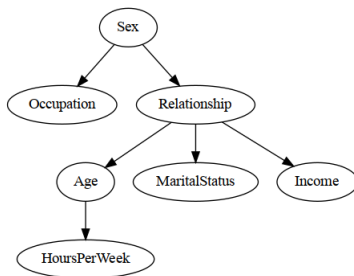
- Difficult to choose the best algorithm/model.
- In finite sample scenario, all algorithms make mistakes.
- Usually requires expert knowledge input.

pgmpy implements option to incorporate expert knowledge to these algorithms.

# Causal Discovery: Expert Knowledge Integration

- Users can specify edges to blacklist/whitelist.
- The algorithm never adds blacklisted edges and only searches over whitelisted edges.

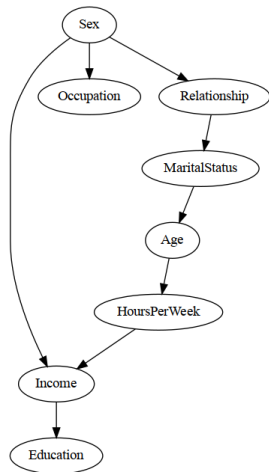
```
1. est = HillClimbSearch(df)
2. blist = [('Sex', 'MaritalStatus'),
3.          ('MaritalStatus', 'Income'),
4.          ('Income', 'Education'),
5.          ('MaritalStatus', 'Age'),
6.          ('Income', 'HoursPerWeek')]
7. dag = est.estimate(
8.     scoring_method="bicscore",
9.     black_list=blist)
```



# Causal Discovery: Expert Knowledge Integration

- Score based methods perform local optimization.

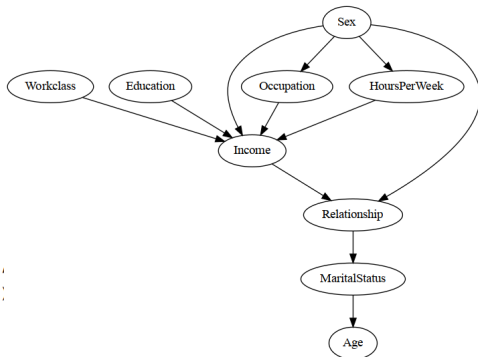
```
1. from pgmpy.base import DAG
2. est = HillClimbSearch(df)
3. start_dag=DAG(
4.     [('Age', 'Income'),
5.      ('Age', 'Education'),
6.      ('Age', 'MaritalStatus'),
7.      ('Age', 'HoursPerWeek'),
8.      ('NativeCountry', 'Race'),
9.      ('HoursPerWeek', 'Income'),
10.     ('Education', 'Occupation'),
11.     ('Education', 'Income'),
12.     ('Education', 'Workclass'),
13.     ('Occupation', 'Income'),
14.     ('Sex', 'Income'),
15.     ('Relationship', 'MaritalStatus'),
16. ])
17. dag_bic_start = est.estimate(
18.     scoring_method="bicscore",
19.     start_dag=start_dag)
```



# Causal Discovery: Expert Knowledge Integration

- Allows to specify fixed edges that will be present in the final model.

```
1. est = HillClimbSearch(df)
2. dag_bic_fixed = est.estimate(
3.     scoring_method="bicscore",
4.     fixed_edges=[
5.         ('HoursPerWeek', 'Income'),
6.         ('Education', 'Income'),
7.         ('Sex', 'Income'),
8.         ('Occupation', 'Income'),
9.         ('Workclass', 'Income')]
```



# Causal Discovery: Expert-In-The-Loop

Since, causal discovery in practice is an iterative manual process, we designed an algorithm that works with interactive input from the user.

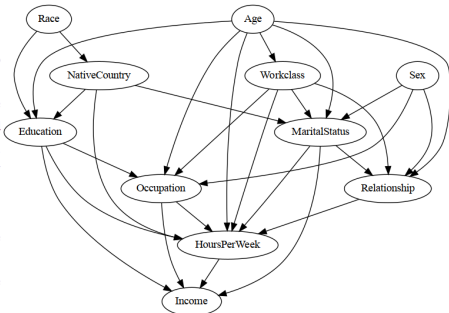
- Finds pair of variables whose observed correlation isn't explained by the model.
- Employs a ranking scheme to rank these violations.
- Iteratively adds and removes edges.
- Uses an expert to decide the direction of the edge.
- Greatly reduces the amount of manual intervention required.



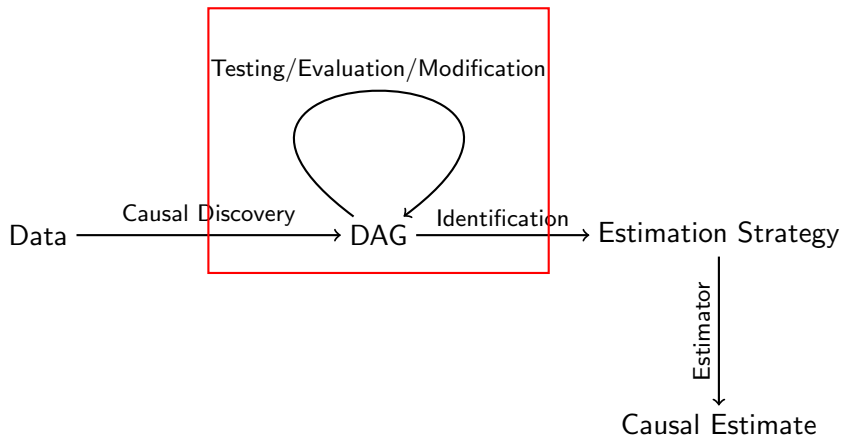
# Causal Discovery: Expert-In-The-Loop

```
1. from pgmpy.estimators import ExpertInLoop
2.
3. descriptions = {
4.     "Age": "The age of a person",
5.     "Workclass": "The workplace where the person is "
6.         "employed such as Private industry, or
7.         "self employed",
8.     "Education": "The highest level of education
9.         the person has finished",
10.    "MaritalStatus": "The marital status of the
11.        person",
12.    "Occupation": "The kind of job the person does
13.        example, sales, craft repair, clerical",
14.    "Relationship": "The relationship status of
15.        the person",
16.    "Race": "The ethnicity of the person",
17.    "Sex": "The sex or gender of the person",
18.    "HoursPerWeek": "The number of hours per week
19.        the person works",
20.    "NativeCountry": "The native country of the
21.        person",
22.    "Income": "The income i.e. amount of money
23.        the person makes",
24. }
```

```
22. est_llm = ExpertInLoop(df)
23. dag_llm = est_llm.estimate(
24.     variable_descriptions=descriptions)
```



# Model Evaluation

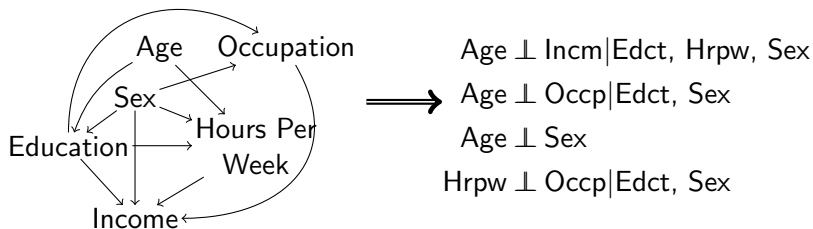


- As causal discovery algorithms make mistakes, important to test and modify learned models.

Similar to unsupervised learning, we do not have any ground truth data, so there is no straightforward way to evaluate models

pgmpy implements a few methods to test and compare models.

# Model Evaluation: Implied Conditional Independences (CIs)



- Each missing edge implies a Conditional Independence (CI).
- Statistical tests can be used to check whether they hold in data.
- Models can be modified based on these tests.

# Model Evaluation: Implied Conditional Independences (CIs)

```
1. >>> from pgmpy.metrics import implied_cis
2. >>> from pgmpy.estimators.CITests import ci_pillai
3. >>> implied_cis(model=dag, data=df, ci_test=ci_pillai)
```

```
4.
5.      u          v          cond_vars          p-value
6. Relationship  Race      [Sex, Workclass, MaritalStatus, Age]  5.516662e-07
7. Relationship  NativeCountry  [Age, Sex, Workclass, MaritalStatus]  2.326022e-03
8. Workclass     Race              []  2.469383e-01
9. Workclass     Income      [HoursPerWeek, Occupation, MaritalS...  9.043659e-01
10. Workclass     Education              [Age]  4.298111e-01
```

# Model Evaluation: Fisher's C Test

- Combines the implied CI tests to summarize it into a single p-value.
- Using some significance threshold, we can decide whether the model fits well to the data.

```
1. >>> from pgmpy.metrics import fisher_c
2. >>> from pgmpy.estimators.CITests import chi_square
3.
4. >>> fisher_c(dag_bic_fixed, df, ci_test=chi_square)
5. 6.136802177536538e-11
6.
7. >>> cancer_model = get_example_model('cancer')
8. >>> sim_df = cancer_model.simulate(int(1e4))
9. >>> fisher_c(cancer_model, sim_df, ci_test=chi_square)
10. 0.9464350296111621
```

# Model Evaluation: Correlation Score

- Compares whether variables that are correlated in data are also correlated in the model.
- Computes the F1-score based on this comparison.

```
1. >>> from pgmpy.metrics import correlation_score
2. >>> correlation_score(dag_bic_start, df)
3. 0.5128205128205128
4. >>> correlation_score(dag_bic_fixed, df)
5. 0.5789473684210527
```

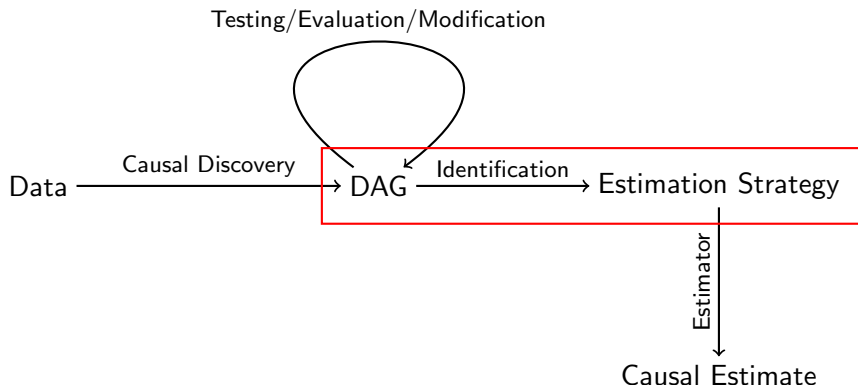
# Model Evaluation: Structure Score Metrics

- Scores the network structure based on how well they fit to data.
- Useful for comparing multiple models.

```
1. >>> from pgmpy.metrics import structure_score
2. >>> structure_score(model=dag_bic_start,
3.                      data=df,
4.                      scoring_method='bic')
5. -11472.356725048614
6. >>> structure_score(model=dag_bic_fixed,
7.                      data=df,
8.                      scoring_method='bic')
9. -50993.46795221821
```

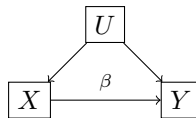
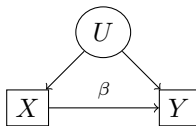


# Identification



- After decided the DAG, we can start estimating causal effects.
- **Identification:** Is the causal effect of interest estimable?
- Everything is identified if all variables are observed.

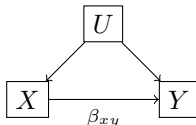
# Identification: do-calculus



- Theoretically, do-calculus provides a complete solution to identification.
- do-calculus can give an estimand for every identified causal effect.
- However, no efficient algorithms to get these estimands using do-calculus.
- In practice, we rely on a set of simpler identification strategies that work in special cases.

# Identification: Backdoor Criterion

- For cases when biasing paths can be blocked by conditioning.
- Finds adjustment variables to block confounding paths.



$$\beta_{xy} : Y \sim X + U$$

## `get_all_backdoor_adjustment_sets( $X$ , $Y$ )`

[source]

Returns a list of all adjustment sets per the back-door criterion.

A set of variables  $Z$  satisfies the back-door criterion relative to an ordered pair of variables  $(X_i, X_j)$  in a DAG  $G$  if:

- i. no node in  $Z$  is a descendant of  $X_i$ ; and
- ii.  $Z$  blocks every path between  $X_i$  and  $X_j$  that contains an arrow into  $X_i$ .

**Parameters:**

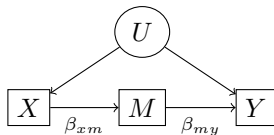
- $X$  ([\*str\* \(variable name\)](#)) – The cause/exposure variables.
- $Y$  ([\*str\* \(variable name\)](#)) – The outcome variable.

**Returns:**

- **frozenset** (A frozenset of frozensets)

# Identification: Front-door Criterion

- Used for scenarios when there are mediator variables.
- Double application of backdoor criterion.



$$\beta_{xm} : M \sim X$$

$$\beta_{my} : Y \sim M + X$$

$$\beta_{xy} = \beta_{xm}\beta_{my}$$

**get\_all\_frontdoor\_adjustment\_sets**( $X, Y$ )

[\[source\]](#)

Identify possible sets of variables,  $Z$ , which satisfy the front-door criterion relative to given  $X$  and  $Y$ .

$Z$  satisfies the front-door criterion if:

- $Z$  intercepts all directed paths from  $X$  to  $Y$
- there is no backdoor path from  $X$  to  $Z$
- all back-door paths from  $Z$  to  $Y$  are blocked by  $X$

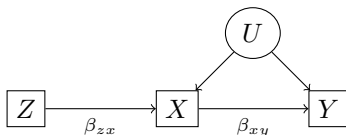
**Parameters:**

- **X** (*str* (variable name)) – The cause/exposure variables.
- **Y** (*str* (variable name)) – The outcome variable

**Returns:** **frozenset**

# Identification: Instrumental Variables

- When there are variables in the model that are correlated with the outcome variable only through the exposure variable.



$$\beta_{zx}\beta_{xy} : Y \sim Z$$

$$\beta_{zx} : X \sim Z$$

$$\beta_{xy} = \frac{\beta_{zx}\beta_{xy}}{\beta_{zx}}$$

**get\_ivs**(*X*, *Y*, *scaling\_indicators*={})

[source]

Returns the Instrumental variables(IVs) for the relation  $X \rightarrow Y$

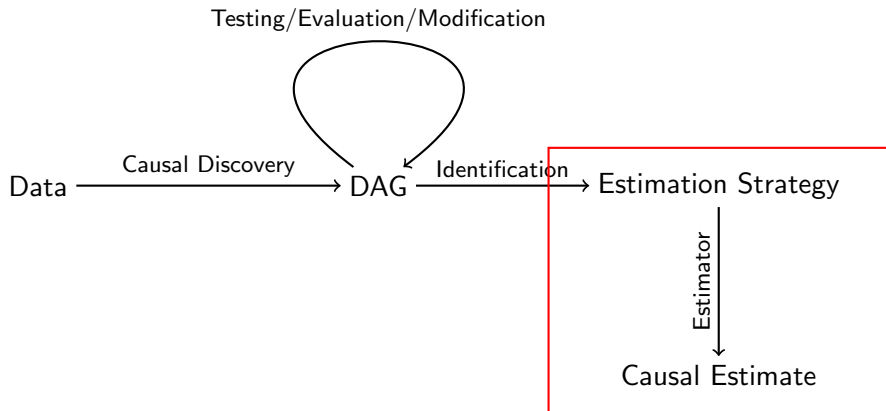
**Parameters:**

- X (node)** – The variable name (observed or latent)
- Y (node)** – The variable name (observed or latent)
- scaling\_indicators** (*dict* (optional)) – A dict representing which observed variable to use as scaling indicator for the latent variables. If not given the method automatically selects one of the measurement variables at random as the scaling indicator.

**Returns:** **set** – The set of Instrumental Variables for  $X \rightarrow Y$ .

**Return type:** {str}

# Estimation



- Identification methods give the estimand.
- Any estimator can be used to make the estimates.
- Usually linear models are used for their interpretability.

# Some other features.

# Simulations

- Parameterized DAGs are generative models.
- Simulated data can be used to evaluate methods.
- Can be used for approximate inference.
- Helpful in explaining concepts like confounding, bias, etc.

```
1. >>> from pgmpy.utils import get_example_model
2. >>> sachs_model = get_example_model('sachs')
3. >>> sachs_model.simulate(int(1e4))
4.
```

	PIP2	P38	Jnk	PKA	Plcg	Mek	Akt	PKC	Erk	PIP3	Raf
6. 0	LOW	LOW	AVG	AVG	LOW	AVG	LOW	LOW	AVG	HIGH	LOW
7. 1	LOW	LOW	AVG	LOW	LOW	HIGH	HIGH	LOW	HIGH	AVG	AVG
8. 2	LOW	LOW	LOW	AVG	LOW	LOW	LOW	AVG	AVG	HIGH	LOW
9. 3	LOW	LOW	LOW	AVG	LOW	AVG	AVG	AVG	AVG	HIGH	AVG
10. 4	AVG	HIGH	HIGH	LOW	AVG	LOW	LOW	LOW	LOW	HIGH	HIGH



# Simulations

```
1. simulate(n_samples=10,  
2.         do=None,  
3.         evidence=None,  
4.         virtual_evidence=None,  
5.         virtual_intervention=None,  
6.         include_latents=False,  
7.         partial_samples=None,  
8.         seed=None,  
9.         show_progress=True)
```

# Extensibility

- Causal Inference is a very active field of research.
- pgmpy offers easy ways to extend/modify algorithms.
- Methods accept custom functions as argument.
- Base classes make it easier to implement new algorithms and can be plugged into other functionality.


```
1. from pgmpy.estimators import PC
2.
3. def random_ci(X, Y, Z, data, boolean=True):
4.     return np.random.choice([True, False])
5.
6. est = PC(df)
7. est.estimate(ci_test=random_ci)
```

# Future Plans

- Focus on implementing more practical methods such as bootstrapping, model comparison methods.
- Improve support for mixed data, i.e., combination of categorical, ordinal, and continuous variables.
- Add more commonly used causal discovery, and identification algorithms.

# Thank you

: [pgmpy/pgmpy](https://github.com/pgmpy/pgmpy)

: [ankurankan@gmail.com](mailto:ankurankan@gmail.com)