

# Inférence causale sous modèles longitudinaux simplifiés

Lola Étiévant

Institut Camille Jordan, Université Claude Bernard Lyon 1

Journées du GDR Statistique et Santé  
Nantes, 27-28 septembre 2018

# Introduction

- Étudier l'effet causal de différents facteurs sur la survenue de pathologies, à partir de données observationnelles

# Introduction

- Étudier l'effet causal de différents facteurs sur la survenue de pathologies, à partir de données observationnelles
  - Inférence causale

# Introduction

- Étudier l'effet causal de différents facteurs sur la survenue de pathologies, à partir de données observationnelles

→ Inférence causale

$$X \longrightarrow Y$$

(CS)

# Introduction

- Étudier l'effet causal de différents facteurs sur la survenue de pathologies, à partir de données observationnelles

→ Inférence causale

$$X \longrightarrow Y$$

(CS)

$$ATE_{CS} := \mathbb{E}_{CS} (Y^{X=1} - Y^{X=0}),$$

# Introduction

- Étudier l'effet causal de différents facteurs sur la survenue de pathologies, à partir de données observationnelles

→ Inférence causale

$$X \longrightarrow Y$$

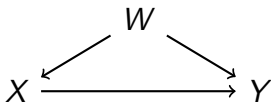
(CS)

$$\begin{aligned}ATE_{CS} &:= \mathbb{E}_{CS} (Y^{X=1} - Y^{X=0}), \\ &= \mathbb{E}(Y | X = 1) - \mathbb{E}(Y | X = 0).\end{aligned}$$

# Introduction

- Étudier l'effet causal de différents facteurs sur la survenue de pathologies, à partir de données observationnelles

→ Inférence causale



(CS)

$$ATE_{CS} = \sum_w [\mathbb{E}(Y \mid X = 1, W = w) - \mathbb{E}(Y \mid X = 0, W = w)] \times \mathbb{P}(W = w).$$

# Introduction

- Souvent, le vrai modèle causal fait intervenir des variables qui varient au cours du temps,  
→ Exemple : Obésité, consommation d'alcool...



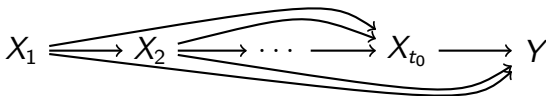
# Introduction

- Souvent, le vrai modèle causal fait intervenir des variables qui varient au cours du temps,

→ Exemple : Obésité, consommation d'alcool...

$$X \longrightarrow Y$$

(CS)

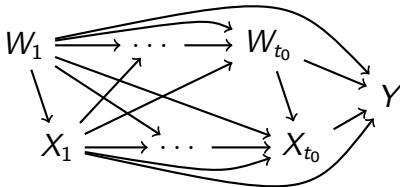


(L)

# Introduction

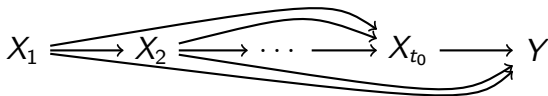
- Souvent, le vrai modèle causal fait intervenir des variables qui varient au cours du temps,

→ Exemple : Obésité, consommation d'alcool...



(L)

# Introduction

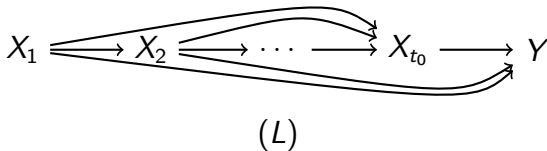


(L)

$$ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) := \mathbb{E}_L \left( Y^{\bar{X}_{t_0}=\bar{x}_{t_0}} - Y^{\bar{X}_{t_0}=\bar{x}_{t_0}^*} \right),$$

avec  $\bar{X}_{t_0} = (X_1, \dots, X_{t_0})$ ,  $\bar{x}_{t_0} = (x_1, \dots, x_{t_0})$  et  $\bar{x}_{t_0}^* = (x_1, \dots, x_{t_0})$ .

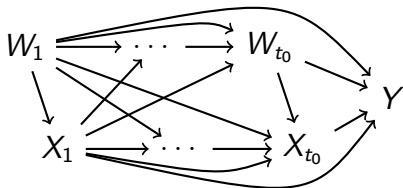
# Introduction



$$\begin{aligned}ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) &:= \mathbb{E}_L \left( Y^{\bar{X}_{t_0} = \bar{x}_{t_0}} - Y^{\bar{X}_{t_0} = \bar{x}_{t_0}^*} \right), \\ &= \mathbb{E} \left( Y \mid \bar{X}_{t_0} = \bar{x}_{t_0} \right) - \mathbb{E} \left( Y \mid \bar{X}_{t_0} = \bar{x}_{t_0}^* \right),\end{aligned}$$

avec  $\bar{X}_{t_0} = (X_1, \dots, X_{t_0})$ ,  $\bar{x}_{t_0} = (x_1, \dots, x_{t_0})$  et  $\bar{x}_{t_0}^* = (x_1, \dots, x_{t_0})$ .

# Introduction

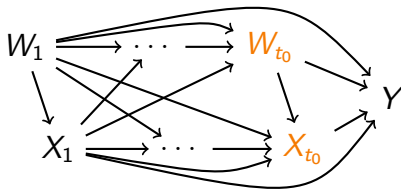


(L)

$$ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) = \sum_{\bar{w}_{t_0}} \left[ \mathbb{E}(Y \mid \bar{X}_{t_0} = \bar{x}_{t_0}, \bar{W}_{t_0} = \bar{w}_{t_0}) \right. \\ \times \prod_{t=1}^{t_0} \mathbb{P}(W_t = w_t \mid \bar{W}_{t-1} = \bar{w}_{t-1}, \bar{X}_{t-1} = \bar{x}_{t-1}) \\ \left. - \mathbb{E}(Y \mid \bar{X}_{t_0} = \bar{x}_{t_0}^*, \bar{W}_{t_0} = \bar{w}_{t_0}) \right. \\ \left. \times \prod_{t=1}^{t_0} \mathbb{P}(W_t = w_t \mid \bar{W}_{t-1} = \bar{w}_{t-1}, \bar{X}_{t-1} = \bar{x}_{t-1}^*) \right].$$

avec  $\bar{W}_{t_0} = (W_1, \dots, W_{t_0})$  et  $\bar{w}_{t_0} = (w_1, \dots, w_{t_0})$ .

# Introduction

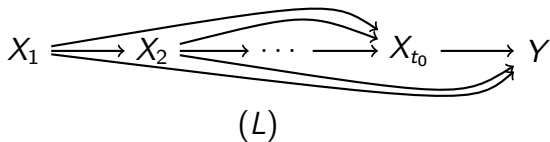


(L)

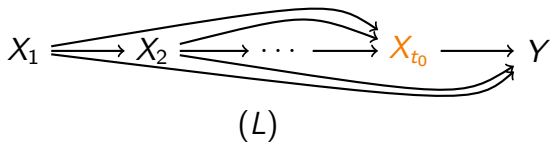
$$ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) = \sum_{\bar{w}_{t_0}} \left[ \mathbb{E}(Y \mid \bar{X}_{t_0} = \bar{x}_{t_0}, \bar{W}_{t_0} = \bar{w}_{t_0}) \right. \\ \times \prod_{t=1}^{t_0} \mathbb{P}(W_t = w_t \mid \bar{W}_{t-1} = \bar{w}_{t-1}, \bar{X}_{t-1} = \bar{x}_{t-1}) \\ \left. - \mathbb{E}(Y \mid \bar{X}_{t_0} = \bar{x}_{t_0}^*, \bar{W}_{t_0} = \bar{w}_{t_0}) \right. \\ \left. \times \prod_{t=1}^{t_0} \mathbb{P}(W_t = w_t \mid \bar{W}_{t-1} = \bar{w}_{t-1}, \bar{X}_{t-1} = \bar{x}_{t-1}^*) \right].$$

avec  $\bar{W}_{t_0} = (W_1, \dots, W_{t_0})$  et  $\bar{w}_{t_0} = (w_1, \dots, w_{t_0})$ .

# Introduction

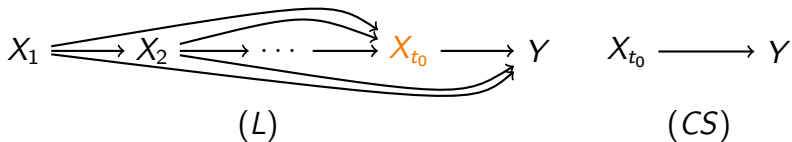


# Introduction

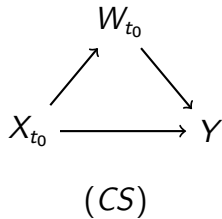
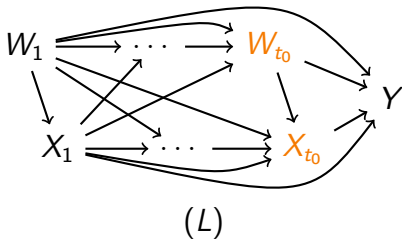




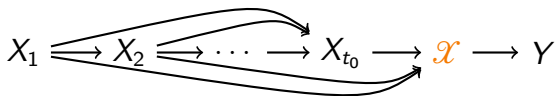
# Introduction



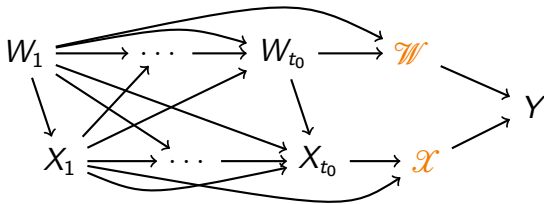
# Introduction



# Introduction

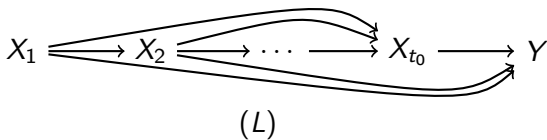


(L)

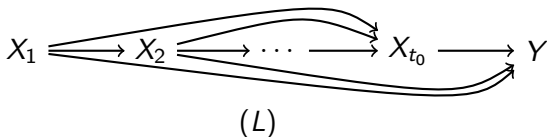


(L)

## Modèle causal longitudinal et données transversales



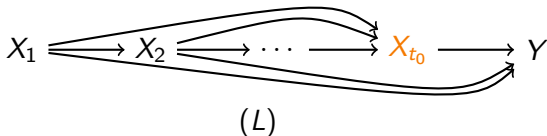
## Modèle causal longitudinal et données transversales



$$\begin{aligned}ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) &:= \mathbb{E}_L \left( Y^{\bar{X}_{t_0} = \bar{x}_{t_0}} - Y^{\bar{X}_{t_0} = \bar{x}_{t_0}^*} \right), \\ &= \mathbb{E} \left( Y \mid \bar{X}_{t_0} = \bar{x}_{t_0} \right) - \mathbb{E} \left( Y \mid \bar{X}_{t_0} = \bar{x}_{t_0}^* \right),\end{aligned}$$

avec  $\bar{X}_{t_0} = (X_1, \dots, X_{t_0})$ ,  $\bar{x}_{t_0} = (x_1, \dots, x_{t_0})$  et  $\bar{x}_{t_0}^* = (x_1, \dots, x_{t_0})$ .

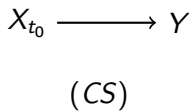
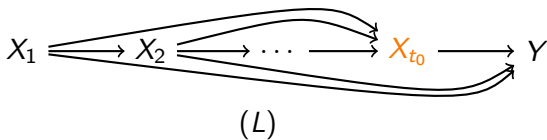
## Modèle causal longitudinal et données transversales



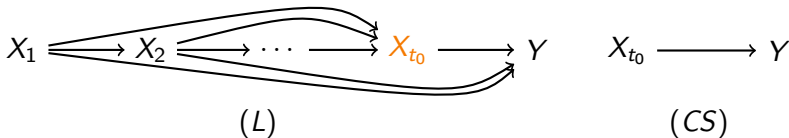
$$\begin{aligned}ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) &:= \mathbb{E}_L \left( Y^{\bar{X}_{t_0} = \bar{x}_{t_0}} - Y^{\bar{X}_{t_0} = \bar{x}_{t_0}^*} \right), \\ &= \mathbb{E} \left( Y \mid \bar{X}_{t_0} = \bar{x}_{t_0} \right) - \mathbb{E} \left( Y \mid \bar{X}_{t_0} = \bar{x}_{t_0}^* \right),\end{aligned}$$

avec  $\bar{X}_{t_0} = (X_1, \dots, X_{t_0})$ ,  $\bar{x}_{t_0} = (x_1, \dots, x_{t_0})$  et  $\bar{x}_{t_0}^* = (x_1, \dots, x_{t_0})$ .

## Modèle causal longitudinal et données transversales



## Modèle causal longitudinal et données transversales

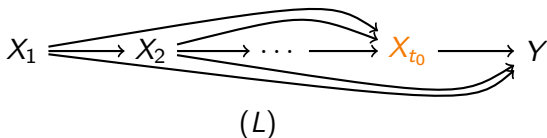


Sous le modèle (CS) nous avons

$$ATE_{CS} := \mathbb{E}_{CS} \left( Y^{X_{t_0}=1} - Y^{X_{t_0}=0} \right),$$



## Modèle causal longitudinal et données transversales



Sous le modèle (CS) nous avons

$$\begin{aligned}ATE_{CS} &:= \mathbb{E}_{CS} \left( Y^{X_{t_0}=1} - Y^{X_{t_0}=0} \right), \\ &= \mathbb{E}(Y \mid X_{t_0} = 1) - \mathbb{E}(Y \mid X_{t_0} = 0).\end{aligned}$$

## Modèle causal longitudinal et données transversales

Sous le modèle (L)

$$\begin{aligned}ATE_{CS} &= \sum_{\bar{x}_{t_0-1}} \sum_{\bar{x}_{t_0-1}^*} ATE_L((\bar{x}_{t_0-1}, 1); (\bar{x}_{t_0-1}^*, 0)) \\ &\quad \times \mathbb{P}(\bar{X}_{t_0-1} = \bar{x}_{t_0-1} \mid X_{t_0} = 1) \\ &\quad \times \mathbb{P}(\bar{X}_{t_0-1} = \bar{x}_{t_0-1}^* \mid X_{t_0} = 0).\end{aligned}$$

## Modèle causal longitudinal et données transversales

Sous le modèle (L)

$$\begin{aligned}ATE_{CS} = & \sum_{\bar{x}_{t_0-1}} \sum_{\bar{x}_{t_0-1}^*} ATE_L((\bar{x}_{t_0-1}, 1); (\bar{x}_{t_0-1}^*, 0)) \\ & \times \mathbb{P}(\bar{X}_{t_0-1} = \bar{x}_{t_0-1} \mid X_{t_0} = 1) \\ & \times \mathbb{P}(\bar{X}_{t_0-1} = \bar{x}_{t_0-1}^* \mid X_{t_0} = 0).\end{aligned}$$

En général, ce n'est pas une mesure pertinente,

## Modèle causal longitudinal et données transversales

Sous le modèle ( $L$ )

$$\begin{aligned}ATE_{CS} &= \sum_{\bar{x}_{t_0-1}} \sum_{\bar{x}_{t_0-1}^*} ATE_L((\bar{x}_{t_0-1}, 1); (\bar{x}_{t_0-1}^*, 0)) \\ &\quad \times \mathbb{P}(\bar{X}_{t_0-1} = \bar{x}_{t_0-1} \mid X_{t_0} = 1) \\ &\quad \times \mathbb{P}(\bar{X}_{t_0-1} = \bar{x}_{t_0-1}^* \mid X_{t_0} = 0).\end{aligned}$$

En général, ce n'est pas une mesure pertinente,

→ Les termes tels que  $\mathbb{E}_L \left( Y^{\bar{X}_{t_0}=(0,\dots,0,1)} - Y^{\bar{X}_{t_0}=(1,\dots,1,0)} \right)$   
ont des poids non négatifs.

## Modèle causal longitudinal et données transversales Le cas d'exposition "stable"

On suppose que  $X_t = 1 \Rightarrow X_{t'} = 1$  pour tout  $t' \geq t$ .

Dans ce cas

$$ATE_{CS} = \sum_{i=0}^{t_0-1} ATE_L((\mathbf{0}_i, \mathbf{1}_{t_0-i}); \mathbf{0}_{t_0}) \\ \times \mathbb{P}(\bar{X}_{t_0-1} = (\mathbf{0}_i, \mathbf{1}_{t_0-i-1}) \mid X_{t_0} = 1).$$

## Modèle causal longitudinal et données transversales

### Le cas d'exposition "stable"

On suppose que  $X_t = 1 \Rightarrow X_{t'} = 1$  pour tout  $t' \geq t$ .

Dans ce cas

$$ATE_{CS} = \sum_{i=0}^{t_0-1} ATE_L((\mathbf{0}_i, \mathbf{1}_{t_0-i}); \mathbf{0}_{t_0}) \\ \times \mathbb{P}(\bar{X}_{t_0-1} = (\mathbf{0}_i, \mathbf{1}_{t_0-i-1}) \mid X_{t_0} = 1).$$

→ Moyenne pondérée des effets causaux longitudinaux comparant l'historique "jamais exposé" à tous les historiques "exposé".

## Modèle causal longitudinal et données transversales

### Le cas d'exposition "stable"

On suppose que  $X_t = 1 \Rightarrow X_{t'} = 1$  pour tout  $t' \geq t$ .

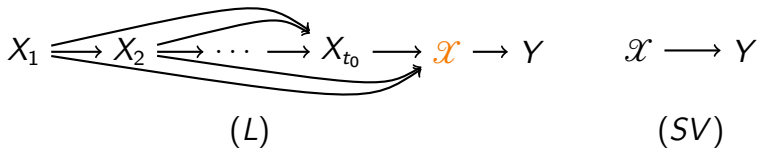
Dans ce cas

$$ATE_{CS} = \sum_{i=0}^{t_0-1} ATE_L((\mathbf{0}_i, \mathbf{1}_{t_0-i}); \mathbf{0}_{t_0}) \\ \times \mathbb{P}(\bar{X}_{t_0-1} = (\mathbf{0}_i, \mathbf{1}_{t_0-i-1}) \mid X_{t_0} = 1).$$

→ Moyenne pondérée des effets causaux longitudinaux comparant l'historique "jamais exposé" à tous les historiques "exposé".

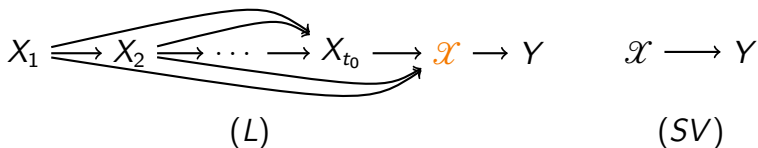
→ Même résultat si un facteur de confusion qui ne varie pas au cours du temps est présent.

## Modèle faisant intervenir des variables résumés de l'exposition passée



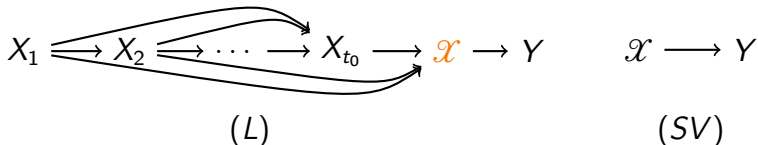


## Modèle faisant intervenir des variables résumés de l'exposition passée



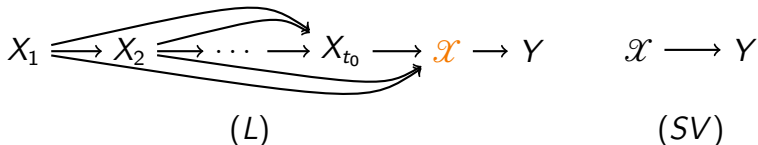
- $\mathcal{X}$  est une mesure "résumé" de l'exposition passée.

## Modèle faisant intervenir des variables résumés de l'exposition passée



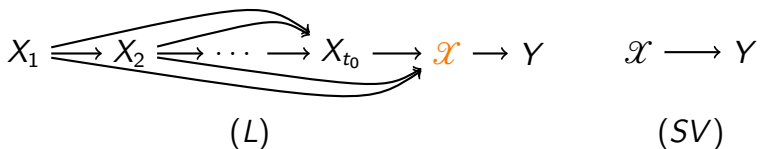
- $\mathcal{X}$  est une mesure "résumé" de l'exposition passée.
- $\bar{X}_{t_0}$  affecte  $Y$  au travers de  $\mathcal{X}$  seulement.

## Modèle faisant intervenir des variables résumés de l'exposition passée



- $\mathcal{X}$  est une mesure "résumé" de l'exposition passée.
- $\bar{X}_{t_0}$  affecte  $Y$  au travers de  $\mathcal{X}$  seulement.
- $\bar{X}_{t_0}$  peut être négligé et le modèle (SV) est correct.

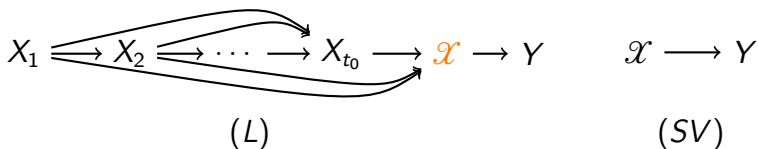
## Modèle faisant intervenir des variables résumés de l'exposition passée



- $\mathcal{X}$  est une mesure "résumé" de l'exposition passée.
- $\bar{X}_{t_0}$  affecte  $Y$  au travers de  $\mathcal{X}$  seulement.
- $\bar{X}_{t_0}$  peut être négligé et le modèle (SV) est correct.

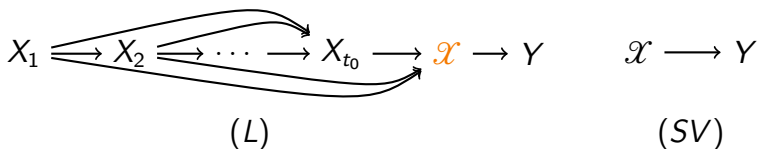
→ L'effet total estimé sous le modèle (SV) est pertinent, c'est celui que l'on aurait obtenu en ayant observé  $\bar{X}_{t_0}$ .

## Modèle faisant intervenir des variables résumés de l'exposition passée



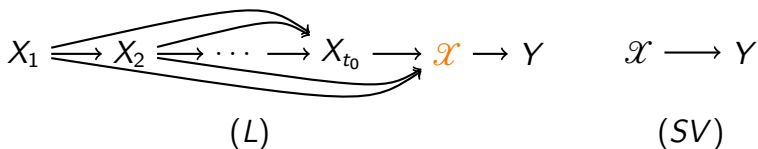
$$\begin{aligned}ATE_{SV}(x; x^*) &:= \mathbb{E}_{SV}(Y^{\mathcal{X}=x} - Y^{\mathcal{X}=x^*}), \\ &= \mathbb{E}(Y \mid \mathcal{X} = x) - \mathbb{E}(Y \mid \mathcal{X} = x^*),\end{aligned}$$

## Modèle faisant intervenir des variables résumées de l'exposition passée



$$\begin{aligned}ATE_{SV}(x; x^*) &:= \mathbb{E}_{SV}(Y^{\mathcal{X}=x} - Y^{\mathcal{X}=x^*}), \\ &= \mathbb{E}(Y \mid \mathcal{X} = x) - \mathbb{E}(Y \mid \mathcal{X} = x^*), \\ &= \mathbb{E}_L(Y^{\mathcal{X}=x} - Y^{\mathcal{X}=x^*}),\end{aligned}$$

## Modèle faisant intervenir des variables résumés de l'exposition passée

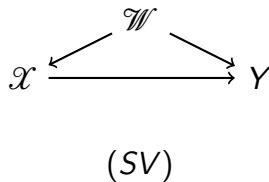
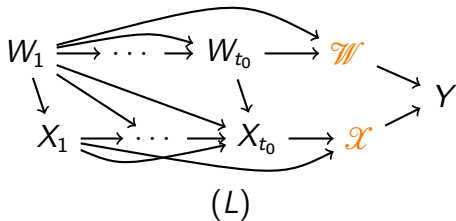


$$\begin{aligned}
 ATE_{SV}(x; x^*) &:= \mathbb{E}_{SV}(Y^{\mathcal{X}=x} - Y^{\mathcal{X}=x^*}), \\
 &= \mathbb{E}(Y \mid \mathcal{X} = x) - \mathbb{E}(Y \mid \mathcal{X} = x^*), \\
 &= \mathbb{E}_L(Y^{\bar{X}_{t_0}=x} - Y^{\bar{X}_{t_0}=x^*}), \\
 &= \mathbb{E}_L(Y^{\bar{X}_{t_0}=\bar{x}_{t_0}} - Y^{\bar{X}_{t_0}=\bar{x}_{t_0}^*}),
 \end{aligned}$$

pour tous  $\bar{x}_{t_0}$  et  $\bar{x}_{t_0}^*$  tels que  $\mathcal{X} = x$  et  $\mathcal{X} = x^*$ , respectivement.

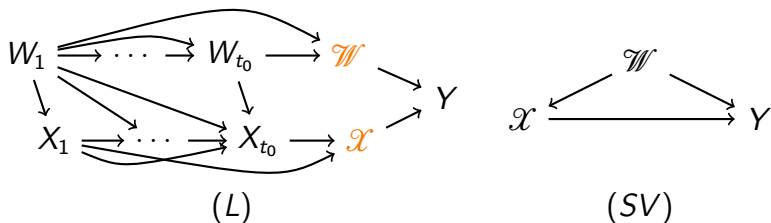
→ Car tous les historiques  $\bar{x}_{t_0}$  donnant lieu à  $\mathcal{X} = x$  ont le même effet sur  $Y$ .

Même modèle, mais avec un facteur de confusion variant au cours du temps



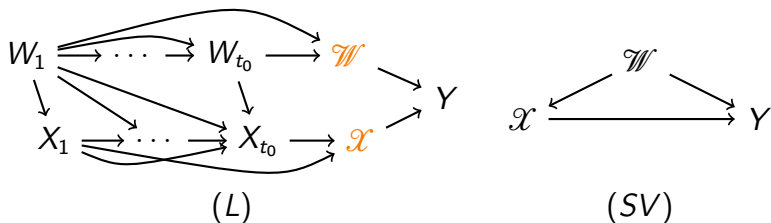


Même modèle, mais avec un facteur de confusion variant au cours du temps



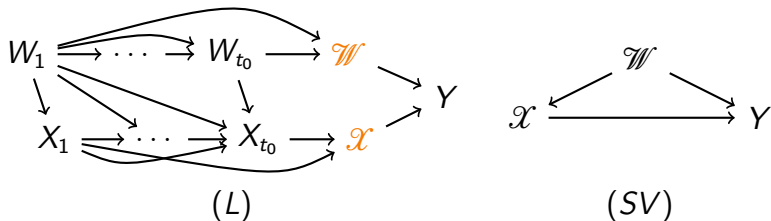
$$ATE_{SV}(x; x^*) = \sum_w [\mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x) - \mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x^*)] \times \mathbb{P}(\mathcal{W} = w),$$

Même modèle, mais avec un facteur de confusion variant au cours du temps



$$\begin{aligned}
 ATE_{SV}(x; x^*) &= \sum_w [\mathbb{E}(Y \mid \mathcal{W} = w, \mathcal{X} = x) \\
 &\quad - \mathbb{E}(Y \mid \mathcal{W} = w, \mathcal{X} = x^*)] \times \mathbb{P}(\mathcal{W} = w), \\
 &= ATE_L(x; x^*),
 \end{aligned}$$

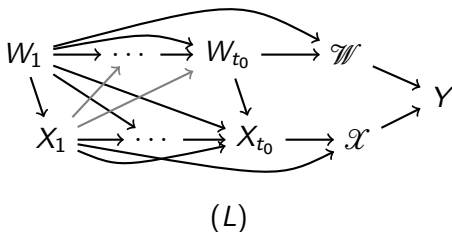
Même modèle, mais avec un facteur de confusion variant au cours du temps



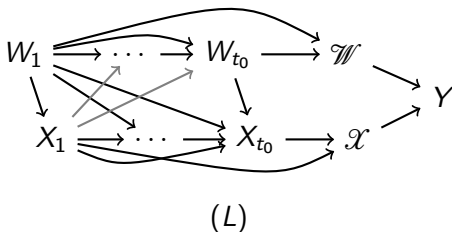
$$\begin{aligned}
 ATE_{SV}(x; x^*) &= \sum_w [\mathbb{E}(Y \mid \mathcal{W} = w, \mathcal{X} = x) \\
 &\quad - \mathbb{E}(Y \mid \mathcal{W} = w, \mathcal{X} = x^*)] \times \mathbb{P}(\mathcal{W} = w), \\
 &= ATE_L(x; x^*), \\
 &= ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*),
 \end{aligned}$$

pour tous  $\bar{x}_{t_0}$  et  $\bar{x}_{t_0}^*$  tels que  $\mathcal{X} = x$  et  $\mathcal{X} = x^*$ , respectivement.

Avec un facteur de confusion variant au cours du temps, affecté par l'exposition

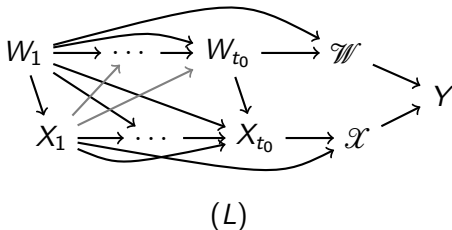


Avec un facteur de confusion variant au cours du temps, affecté par l'exposition



- Typiquement,  $ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) \neq ATE_L(\bar{x}'_{t_0}; \bar{x}'_{t_0}^*)$  même pour  $\bar{x}_{t_0}$  et  $\bar{x}'_{t_0}$  tels que  $\mathcal{X} = x$  et pour  $\bar{x}_{t_0}^*$  et  $\bar{x}'_{t_0}^*$  tel que  $\mathcal{X} = x^*$ .
- $ATE_L(x; x^*) = \mathbb{E}_L(Y^{\mathcal{X}=x} - Y^{\mathcal{X}=x^*})$  n'a pas vraiment de sens.

Avec un facteur de confusion variant au cours du temps, affecté par l'exposition

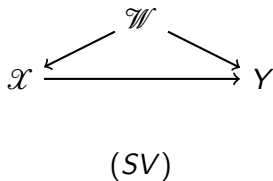
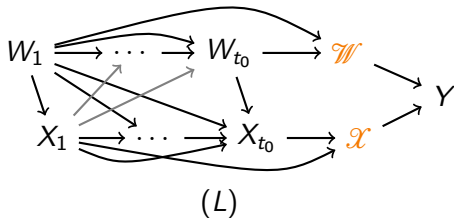


- Il semble plus naturel de considérer

$$\sum_{\bar{x}_{t_0}} \sum_{\bar{x}_{t_0}^*} ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*) \times \mathbb{P}(\bar{X}_{t_0} = \bar{x}_{t_0} \mid \mathcal{X} = x) \times \mathbb{P}(\bar{X}_{t_0} = \bar{x}_{t_0}^* \mid \mathcal{X} = x^*),$$

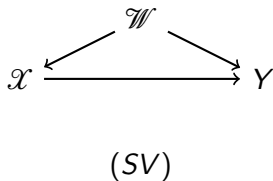
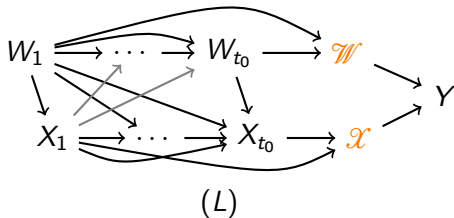
comme l'effet causal d'intérêt.

Avec un facteur de confusion variant au cours du temps, affecté par l'exposition



$$ATE_{SV}(x; x^*) = \sum_w [\mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x) - \mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x^*)] \times \mathbb{P}(\mathcal{W} = w),$$

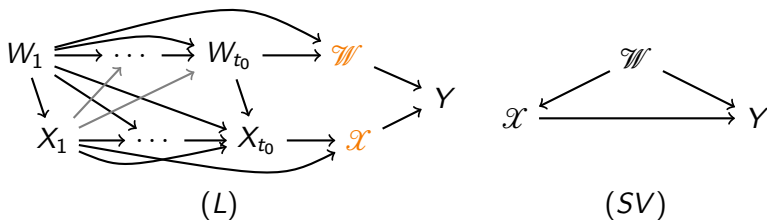
Avec un facteur de confusion variant au cours du temps, affecté par l'exposition



$$\begin{aligned}
 ATE_{SV}(x; x^*) &= \sum_w [\mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x) \\
 &\quad - \mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x^*)] \times \mathbb{P}(\mathcal{W} = w), \\
 &= ATE_L(x; x^*),
 \end{aligned}$$



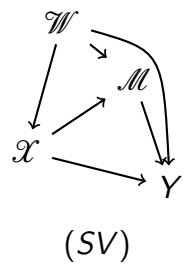
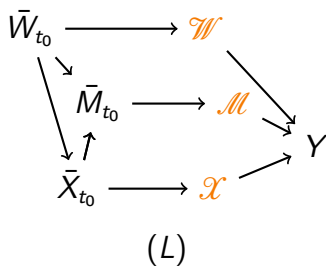
Avec un facteur de confusion variant au cours du temps, affecté par l'exposition



$$\begin{aligned}
 ATE_{SV}(x; x^*) &= \sum_w [\mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x) \\
 &\quad - \mathbb{E}(Y | \mathcal{W} = w, \mathcal{X} = x^*)] \times \mathbb{P}(\mathcal{W} = w), \\
 &= ATE_L(x; x^*), \\
 &\neq ATE_L(\bar{x}_{t_0}; \bar{x}_{t_0}^*),
 \end{aligned}$$

pour généralement tous  $\bar{x}_{t_0}$  et  $\bar{x}_{t_0}^*$  tels que  $\mathcal{X} = x$  et  $\mathcal{X} = x^*$ .


# Avec facteur de confusion et médiateurs variant au cours du temps



## Conclusion

- Nature longitudinale des variables généralement négligée en pratique.

## Conclusion

- Nature longitudinale des variables généralement négligée en pratique.
  - L'inférence basée sur une mesure des expositions à un seul temps conduit à des quantités non liées aux vrais effets causaux sous les modèles longitudinaux.
-  Interprétation des résultats de telles analyses.

## Conclusion


- Nature longitudinale des variables généralement négligée en pratique.
- L'inférence basée sur une mesure des expositions à un seul temps conduit à des quantités non liées aux vrais effets causaux sous les modèles longitudinaux.
- ⚠ Interprétation des résultats de telles analyses.
- Exception : estimation de l'effet total d'une exposition stable, en l'absence de facteur de confusion variant au cours du temps.

## Conclusion

- En l'absence de facteur de confusion variant au cours du temps et affecté par l'exposition d'intérêt, l'inférence basée sur les mesures résumé retourne des quantités pertinentes pour l'estimation des effets totaux.

## Conclusion

- En l'absence de facteur de confusion variant au cours du temps et affecté par l'exposition d'intérêt, l'inférence basée sur les mesures résumé retourne des quantités pertinentes pour l'estimation des effets totaux.

 Les médiateurs sont généralement négligés lorsque l'on s'intéresse aux effets totaux, mais il en existe certainement. Dès qu'un facteur de confusion qui varie au cours du temps existe aussi, les mesures "résumé" ne sont plus suffisantes.