

UNIVERSITÀ DEGLI STUDI DI MILANO-BICOCCA

Scuola di Economia e Statistica

Corso di laurea in

STATISTICA E GESTIONE DELLE INFORMAZIONI



**IL RUOLO DEI GRAFICI CAUSALI (DAG) NELLA
PIANIFICAZIONE, ANALISI E INTERPRETAZIONE
DEGLI STUDI OSSERVAZIONALI**

Relatore: Prof. Rino Bellocchio

Tesi di laurea di:

Sara Conti

Matr. N. 837969

Anno accademico 2020/2021

Indice

Lista delle figure	iii
Introduzione	1
1 Cenni statistici per comprendere meglio l'utilizzo dei DAG	3
1.1 Studi osservazionali	3
1.2 Effetti causali	5
1.3 Differenza tra correlazione, associazione e causalità	7
1.4 Classificazione del bias e confondimento	8
1.5 L'interazione	10
2 DAG	11
2.1 I grafi	11
2.2 I DAG	15
2.3 D-separation	17
2.4 ESC-DAGs	19
2.5 DAG e confondimento	21
2.6 DAG e selection bias	27
2.7 Interaction DAGs (IDAGs)	31
2.8 Come rappresentare i DAG?	33
Conclusione	37
Bibliografia	39

Elenco delle figure

2.1	Grafo rappresentativo di una rete stradale	11
2.2	Archi multipli	12
2.3	Loops	12
2.4	Grafo aciclico (sx) e grafo ciclico (dx)	12
2.5	Differenza grafo orientato e non orientato	13
2.6	Esempio di grafo con freccia bidirezionale	14
2.7	Esempio di grafo con arco non direzionale	14
2.8	Esempio di collisore (Z)	15
2.9	Esempio di DAG	16
2.10	Processo di traduzione di un ESC-DAG	20
2.11	Esempio di DAG in cui L è causa comune del trattamento e dell'outcome	21
2.12	Esempio di DAG in cui non vi sono cause comuni tra trattamento e outcome	24
2.13	Primo esempio di DAG con selection bias	28
2.14	Secondo esempio di DAG con selection bias	28
2.15	Terzo esempio di DAG con selection bias	29
2.16	Quarto esempio di DAG con selection bias	29
2.17	Esempio di DAG in cui non funziona la standardizzazione	30
2.18	Esempio di IDAG	32
2.19	Schermata principale di DAGitty	34
2.20	Esempio di utilizzo di DAGitty	35

Introduzione

Gli studi osservazionali, descrittivi o analitici, si basano sulla raccolta di informazioni necessarie per studiare le ipotesi di ricerca di interesse, spesso senza alcun intervento attivo dei ricercatori. In questo contesto hanno un ruolo molto importante i grafici causali (DAG), che permettono di rappresentare graficamente le interrelazioni tra le variabili di interesse, con lo scopo di suggerire le analisi statistiche appropriate. L'acronimo DAG abbrevia l'espressione "*directed acyclic graph*", cioè grafo orientato aciclico. I DAG non sostituiscono le analisi standard, le loro interpretazioni e i relativi commenti, ma ci aiutano a vedere il contesto studiato da un'altra prospettiva. Permettono infatti un'analisi qualitativa, rispetto alla consueta analisi quantitativa.

Questi grafi rappresentano essenzialmente relazioni causa-effetto tra le variabili considerate nello studio esaminato. Essi, infatti, raffigurano dei modelli causali strutturali attraverso nodi e archi orientati ("freccie" che collegano i nodi fra loro). Ogni nodo corrisponde ad una variabile e ogni arco riproduce la relazione causa-effetto studiata. L'impianto gerarchico della struttura suggerisce le relazioni tra le variabili studiate, permettendo di capire quali variabili dipendono dalle altre e quali, invece, sono indipendenti fra loro. I DAG aiutano anche a comprendere relazioni complicate dal confondimento, o dal *bias* da selezione, o dalle interazioni (in questo caso si parlerà di IDAGs), rivelandosi perciò un ottimo strumento di aiuto per i ricercatori, soprattutto dell'ambito biomedico.

Capitolo 1

Cenni statistici per comprendere meglio l'utilizzo dei DAG

1.1 Studi osservazionali

A differenza degli studi sperimentali, in cui i ricercatori intervengono sul trial o sulla sperimentazione clinica, per esempio assegnando in modo casuale diversi tipi di trattamenti a differenti pazienti, negli studi osservazionali si può semplicemente osservare un fenomeno, che non può essere in alcun modo modificato dall'intervento dei ricercatori. Non si può dunque intromettersi nell'esito di ciò che si sta studiando. Esistono due tipi di studi osservazionali: gli studi descrittivi e gli studi analitici.

Negli studi osservazionali descrittivi, i ricercatori non effettuano analisi statistiche, ma studiano la distribuzione dei fenomeni clinici. Essi hanno un importante ruolo nel riconoscimento di nuove malattie, di nuovi fattori di rischio, ma anche degli effetti collaterali dei farmaci o di altri interventi sanitari, come la valutazione di dispositivi medici, di tecniche chirurgiche e pure di programmi di screening. Si dividono in studi descrittivi e studi ecologici.

Gli studi descrittivi, a volte, sono rappresentati con mappe geografiche, che forniscono informazioni sulla continuità e la prossimità delle aree geografiche interessate da uno stesso fenomeno (oggetto dello studio), che può dunque interessare zone limitrofe oppure comportarsi in modo molto diverso anche a distanza di pochi chilometri. Gli obiettivi di questi studi sono per lo più esplorativi e per questo considerati dei "generatori di ipotesi", che sfruttano gli strumenti della demografia come i tassi di mortalità, d'incidenza e di prevalenza. Questi studi sono vantaggiosi perché l'evento è certo, l'informazione studiata è sempre disponibile (proviene per esempio dall'I-STAT), la sua rilevazione viene effettuata in maniera continuativa e sistematica e

sono inoltre disponibili degli ulteriori studi che ci permettono di comprendere la qualità delle informazioni considerate. Tuttavia, hanno anche degli svantaggi. Essi infatti, sottostimano il fenomeno, ad esclusione del caso in cui la letalità della malattia studiata sia pari al 100%. Ulteriori limiti sono l'impossibilità di misurare malattie non fatali e di informare sul rischio individuale di malattia.

Gli studi ecologici, invece, mettono in relazione l'intensità dell'esposizione ad un fattore (di rischio o di protezione) e la frequenza della malattia, considerando come unità statistica la popolazione e non l'individuo. Diventano problematici quando si confrontano diverse generazioni con differenti esposizioni e in più periodi temporali. Inoltre, portano fallacia ecologica, fornendo una misura media della popolazione che non descrive bene quella del singolo individuo, o, all'opposto, fallacia da individuo, valutando un legame tra esposizioni e malattia a livello individuale, senza tenere conto che alcuni fattori di rischio operano a livello di popolazione. Sono tuttavia molto utili quando si hanno scarse e contraddittorie conoscenze del fenomeno, se l'esposizione d'interesse ha una bassa variabilità all'interno della popolazione e un'alta variabilità tra le popolazioni studiate, quando la variabile di esposizione è impossibile da misurare a livello individuale e quando, dopo aver individuato una forte relazione a livello individuale, si vuole valutarne il legame a livello di popolazione per confermare il suo impatto sulla sanità pubblica.

Gli studi analitici, invece, si suddividono in studi di prevalenza o trasversali, detti anche cross-sectional, e studi longitudinali che possono essere caso-controllo o di coorte. Esistono anche studi ambi-direzionali, in cui l'associazione viene studiata prima con un approccio di coorte e poi con un approccio caso-controllo.

Negli studi di prevalenza un campione di individui è selezionato da una popolazione bersaglio e contattato in un preciso punto nel tempo. Gli individui del campione si classificano in base alla presenza / assenza della malattia e alla presenza / assenza dell'esposizione. Possono causare distorsioni temporali, quando gli individui modificano l'esposizione a causa della malattia, o distorsioni da incidenza-prevalenza, se i casi esposti sopravvivono meno rispetto a quelli generati dai non esposti. Essendo basati su casi prevalenti, piuttosto che incidenti, hanno una limitata possibilità di investigare relazioni etiologiche, non consentono di indagare su malattie di breve durata e non sono adatti per malattie ed esposizioni rare. In certi casi sono comunque vantaggiosi perché possono fornire una fotografia dettagliata dei bisogni di salute di una popolazione, al fine di pianificare interventi assistenziali e preventivi. Inoltre, sono più facili da condurre rispetto agli altri studi analitici, perché non richiedono un *follow-up*.

Negli studi di coorte i soggetti sono selezionati in base all'esposizione e sono seguiti nel tempo (*follow-up*) per valutare l'insorgere della malattia d'interesse. Tra le persone libere dalla malattia e a rischio di contrarla si prendono un campione di esposti e uno di non esposti, li si segue nel tempo e si osserva quanti si ammalano e quanti non si ammalano. Questi studi possono essere affetti dal *selection bias*, o distorsione da selezione, che si verifica quando si seleziona una popolazione di confronto non corretta. Spesso richiedono investimenti di molte risorse, soprattutto se condotti prospetticamente ed inoltre, i cambiamenti sull'esposizione durante il periodo di osservazione sono difficili da controllare. Anche in questo caso ci sono però vantaggi nell'usare questo tipo di studio: l'esposizione è misurata prima dell'insorgenza della malattia, possono essere studiate esposizioni rare tramite appropriata selezione da coorte in studio e può anche essere studiato un intero spettro di effetti conseguenti ad un'esposizione.

Infine, negli studi caso-controllo, i soggetti sono selezionati in base alla presenza o assenza della malattia ed è misurato il livello di esposizione pregressa al fattore d'interesse. Si suddividono in studi caso-controllo con base di popolazione, in cui si prendono tutti i malati e un campione di controlli dalla popolazione bersaglio, e con base ospedaliera, in cui si considera una popolazione bersaglio ricoverata in uno o più ospedali, si prendono tutti i malati ricoverati per la malattia d'interesse e un campione di non malati ricoverati per malattie diverse da quella d'interesse. Richiedono investimenti di minori risorse rispetto agli studi prospettici, possono essere studiate malattie rare e con lungo periodo di induzione-latenza e può essere studiato l'intero spettro delle esposizioni associate ad una malattia. Tuttavia, non si adattano allo studio di esposizioni rare, la selezione di un appropriato gruppo di controllo può essere problematica, la misura dell'esposizione pregressa è difficile da ottenere e la comparabilità tra casi e controlli, in termini di informazione sull'esposizione pregressa, è problematica.

1.2 Effetti causali

Dr. Eleanor Murray (attualmente assistente professore, di epidemiologia, presso la *Boston University School of Public Health*), in una delle sue presentazioni del 2019, parlando di inferenza causale e cercando di riassumerla in maniera molto semplice, afferma che *"causal inference is about what would have happened if the world had been just a little different. If we had a time machine, we could know. Instead,*

*we experiment. If we can't even experiment, we use statistics and assumptions to estimate what would have happened in a world where we could experiment".*¹ Secondo Dr. Murray esistono quindi moltissime variabili che in ogni momento potrebbero cambiare ciò che noi stiamo osservando. In particolare, l'inferenza casuale aiuta nell'analisi di ciò che vogliamo studiare, cercando di individuare e misurare i nessi di causalità.

Nel momento in cui noi umani pensiamo agli effetti causali, confrontiamo il risultato di una certa azione A e il risultato di quando questa azione non avviene. Se i due risultati sono diversi si dice che A ha un effetto causale o preventivo sull'esito dell'esperimento compiuto. Altrimenti, si dice che A non ha nessun effetto causale sul risultato. Gli epidemiologici, gli statistici o gli economisti chiamano l'azione di A intervento, esposizione o trattamento.

Come ci spiega Miguel Hernan² (Professore di Biostatistica ed Epidemiologia presso *Harvard T.H. Chan School of Public Health* e Membro della Facoltà presso *Harvard-MIT Program in Health Sciences and Technology*), se consideriamo l'evento A dicotomo, indicando per esempio se è avvenuto o meno un certo trattamento (0: non trattato, 1: trattato) e Y l'*outcome* dicotomo che indica la morte (0) o la sopravvivenza (1), possiamo avere 4 situazioni diverse: $a=1$ e $Y=1$, $a=1$ e $Y=0$, $a=0$ e $Y=1$ ed infine $a=0$ e $Y=0$; A e Y sono variabili aleatorie, mentre a rappresenta il valore assunto. In certi casi alcune di queste 4 situazioni possono anche coincidere fra loro. In particolare, il trattamento A ha un effetto causale sull'*outcome* Y se l'*outcome* con $a=1$ è diverso dall'*outcome* con $a=0$. Si dice che $Y^{a=1}$ (la risposta nel momento in cui $a=1$) e $Y^{a=0}$ sono *outcome* potenziali o anche *outcome* controfattuali. Alcuni autori preferiscono il termine "potenziali" per enfatizzare che, in base al trattamento ricevuto, ciascuno di questi due *outcome* può essere potenzialmente osservato. Altri autori, invece, preferiscono il termine "controfattuali" per enfatizzare che questi *outcome* rappresentano situazioni che potrebbero non verificarsi. Per ogni individuo, uno degli *outcome* controfattuali, quello che corrisponde al trattamento che l'individuo ha ricevuto, è effettivamente fattuale. Gli effetti causali individuali sono definiti come contrasti dei valori degli *outcome* controfattuali, ma

¹L'inferenza causale tratta di qualcosa che accadrebbe se il mondo fosse giusto un po' diverso. Se avessimo una macchina del tempo potremmo conoscere ciò che sarebbe potuto succedere o che potrebbe accadere, ma non avendola si fanno molti esperimenti. Se non si possono nemmeno fare esperimenti, si usano la statistica e le assunzioni sulla stima di ciò che potrebbe accadere in un mondo dove si potrebbe sperimentare.

²Capitolo 1.1. Individual causal effect, pag 12, di Hernan M., Robins J., Causal Inference: What If, 2020, manoscritto inedito.

solo uno di questi *outcome* è osservato per ogni individuo, quello che corrisponde al trattamento che è stato effettivamente ricevuto dall'individuo. Tutti gli altri *outcome* controfattuali rimangono inosservati. A causa dei *missing data* strutturali, gli effetti individuali non possono essere identificati, cioè non possono essere espressi in funzione dei dati osservati.

Per definire un effetto causale individuale abbiamo bisogno di definire un *outcome* di interesse, le azioni $a=1$ e $a=0$ che devono essere comparate, e gli *outcome* $Y^{a=1}$ e $Y^{a=0}$ che devono essere confrontati. Tuttavia, identificare gli effetti causali individuali è genericamente impossibile, dunque ci si può soffermare sugli effetti causali medi in una popolazione di individui. In questo caso si ha bisogno di un *outcome* d'interesse, le azioni $a=0$ e $a=1$ che devono essere comparate, e una popolazione abbastanza definita di individui in cui si possono comparare $Y^{a=1}$ e $Y^{a=0}$.

Gli effetti causali possono essere rappresentati dalla differenza di rischio causale, dal rapporto tra rischi e dall'*odds ratio*.

$$\text{Differenza tra rischi : } Pr[Y^{a=1} = 1] - Pr[Y^{a=0} = 1]$$

$$\text{Rapporto tra rischi : } \frac{Pr[Y^{a=1}=1]}{Pr[Y^{a=0}=1]}$$

$$\text{Odds ratio : } \frac{Pr[Y^{a=1}=1]/Pr[Y^{a=1}=0]}{Pr[Y^{a=0}=1]/Pr[Y^{a=0}=0]}$$

Questi parametri quantificano la forza dello stesso effetto causale su scale diverse e dato che misurano tutte l'effetto causale, ci si può riferire a loro come misure di effetto. Ogni misura di effetto può essere utilizzata per scopi diversi.

1.3 Differenza tra correlazione, associazione e causalità

Prima di comprendere meglio cosa siano i DAG e come si utilizzino, è importante stabilire la differenza tra causalità, correlazione e associazione.

Quando la probabilità che Y sia pari a 1 è uguale sia nel caso in cui $A=1$ che nel caso in cui $A=0$, allora si dice che il trattamento A e l'*outcome* Y sono indipendenti, che A non è associato con Y o che A non predice Y . Le misure sopra citate (differenza di rischi, rapporto tra rischi e l'*odds ratio*) misurano la forza dell'associazione, quando esiste, e per questo si chiamano anche misure di associazione. Per A binario, Y e A non sono associati se e solo se non sono statisticamente correlati.

La correlazione ci mostra i pattern esistenti tra variabili che tendono a muoversi in parallelo, ma non è necessariamente indice di causalità. È infatti possibile

ritrovare una correlazione statistica significativa tra due variabili che tuttavia non hanno alcun legame di causalità. Esse potrebbero semplicemente essere associate ad una terza variabile causale che tende a verificarsi simultaneamente alle altre due variabili. Esiste addirittura un sito web ³ chiamato *Spurious Correlations*, che relaziona fra loro dati diversi facendo notare che la correlazione non indica sempre causalità. Alcuni esempi riportati nell'articolo "Ecco il generatore di correlazioni assurde" di Sandro Iannaccone (2014) sono davvero divertenti: *Spurious Correlations* illustra come ci sia una forte correlazione tra il consumo di margarina e il tasso di divorzi, oppure tra la crescita delle colonie di api e il tasso dei matrimoni.

La causalità, invece, spiega il rapporto causa-effetto tra due variabili. Ad esempio, non svolgere attività fisica e mangiare molto e male, potrebbe causare problemi cardiaci. L'inferenza relativa alla causalità è relativa al "what if", ovvero a domande come "Quale potrebbe essere il rischio SE tutti fossero trattati?" oppure "Quale potrebbe essere il rischio SE tutti non fossero trattati?", mentre l'inferenza relativa all'associazione si riferisce a domande come "Qual è il rischio nei trattati?" o "Qual è il rischio nei non trattati?".

Questi rapporti causali sono reali ed importanti da analizzare ed è quindi fondamentale stabilire con certezza se le analisi svolte nel proprio studio portano a semplici correlazioni oppure sono vere e proprie causalità.

1.4 Classificazione del bias e confondimento

La parola "*bias*" è frequentemente usata da coloro che lavorano sull'inferenza causale. Ci sono molti usi del termine *bias*. Si dice che c'è un *bias* sistematico quando i dati sono insufficienti per identificare l'effetto causale anche con una dimensione infinita del campione. Informalmente parliamo di *bias* sistematico anche riferendoci a qualsiasi associazione strutturale tra il trattamento e l'*outcome* che non nasce da un effetto causale del trattamento sull'*outcome* nella popolazione d'interesse.

La fonte principale di *bias* è la mancanza di scambiabilità (dal punto di vista interpretativo la condizione di scambiabilità coincide con l'idea di analogia o equivalenza tra i valori osservabili, molto più efficacemente della condizione di indipendenza stocastica e uguale distribuzione degli stessi) tra i trattati e i non trattati. La mancanza di *bias* implica che la misura di associazione nella popolazione sia una stima

³<https://www.tylervigen.com/spurious-correlations>

consistente della corrispondente misura di effetto nella popolazione. In particolare, la mancanza di scambiabilità può derivare da due diverse cause strutturali:

- Cause comuni: quando il trattamento e l'esito condividono una causa comune, la misura dell'associazione generalmente differisce dalla misura dell'effetto. Molti epidemiologi usano il termine confondente per riferirsi a questo concetto;
- Condizionamento sugli effetti comuni: questa struttura è la fonte di *bias* che molti epidemiologi chiamano *bias* di selezione.

Esistono anche altre fonti di *bias*, per esempio l'errore di misurazione. Infatti, spesso nella pratica è previsto un certo grado di errore di misurazione. Questo tipo di *bias* è detto *bias* di misurazione o di informazione.

In epidemiologia si parla di distorsione da confondimento quando si osserva un'associazione tra una data esposizione e una malattia per effetto di una terza variabile (o gruppi di variabili) usualmente chiamata variabile di confondimento o confondente. Affinché una variabile sia un potenziale confondente deve essere legata all'esposizione, deve essere un fattore di rischio per l'evento e non deve trovarsi sul *path-way* causale tra l'esposizione e l'evento. Un esempio di confondimento spesso riconosciuto è l'età. Confrontare gruppi o singoli soggetti che hanno età molto diverse potrebbe portare a conclusioni errate, in quanto l'età ha molti effetti sulla persona e quindi, per esempio, un ragazzo di 18 anni non sarà mai confrontabile con un adulto di 65 anni. Quando si è in presenza di confondimento le stime di associazione strato-specifiche sono tra di loro molto simili e la loro sintesi si differenzia dalla stima grezza. Questo significa che la variabile confondente non è distribuita in modo omogeneo nei gruppi e ciò modifica in modo rilevante la stima dell'associazione tra il tipo di intervento medico e l'outcome.

Negli esperimenti randomizzati il trattamento viene assegnato dal lancio di una moneta, ma negli studi osservazionali il trattamento può essere determinato da molti fattori. Se questi fattori influenzano il rischio di sviluppare il risultato, gli effetti di quei fattori si intrecciano con l'effetto del trattamento. Diciamo quindi che c'è confusione, che è solo una forma di mancanza di scambiabilità tra il trattato e il non trattato. La confusione è spesso vista come il principale difetto degli studi osservazionali. In presenza di confusione, vale che "l'associazione non è causalità", anche se la popolazione in studio è arbitrariamente grande.

Questi tipi di *bias* (confondimento, *bias* da selezione, *bias* di misurazione) possono sorgere sia negli studi osservazionali, di cui ci occupiamo, ma anche negli esperimenti randomizzati.

1.5 L'interazione

Associata al confondimento vi è spesso l'interazione, che però non rientra in alcun tipo di *bias*. Bisogna saper distinguere quando si parla del primo e quando si parla della seconda. In statistica si ha un'interazione quando, considerando ad esempio il modello $Y_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$, il parametro γ_{ij} è diverso da zero, ossia l'effetto sulla media di Y, dovuto ai livelli di A, cambia al variare dei livelli di B, e viceversa. Il termine interazione misura gli effetti prodotti congiuntamente dal livello i-esimo di A e il livello j-esimo di B, senza che tali effetti possano imputarsi distintamente ai due fattori, in quanto dovuti alla loro combinazione.

Come detto, l'interazione richiede l'intervento congiunto di due o più trattamenti e quando ciò avviene, l'identificazione dell'interazione ci permette di adottare gli interventi più efficaci.

Più tecnicamente, una variabile V si dice modificatore di un effetto A su Y quando l'effetto causale medio di A su Y varia in base ai livelli di V. Quindi quando diciamo che V modifica l'effetto di A non stiamo considerando V e A come variabili sullo stesso piano, perché solo A è considerata una variabile su cui potremmo ipoteticamente intervenire. Al contrario, la definizione di interazione tra A e E mette sullo stesso piano entrambi i trattamenti (A e E), come si intende da quanto detto nelle righe precedenti sull'interazione. Se la scambiabilità, la positività e la consistenza sono le tre condizioni chiave per identificare un effetto causale medio di un trattamento A su un *outcome* Y, dato che l'interazione è l'effetto congiunto di due (o più) trattamenti, per identificarla sono necessarie scambiabilità, positività e consistenza per entrambi i trattamenti (o più).

Capitolo 2

DAG

2.1 I grafi

Il grafo è una struttura relazionale composta da un insieme finito di oggetti detti nodi, o vertici, e da un insieme di relazioni tra coppie di oggetti detti archi, o spigoli. Attraverso i grafi si possono schematizzare moltissime strutture reali, da una rete stradale (nodi: incroci, archi: strade; esempio in Figura 2.1)

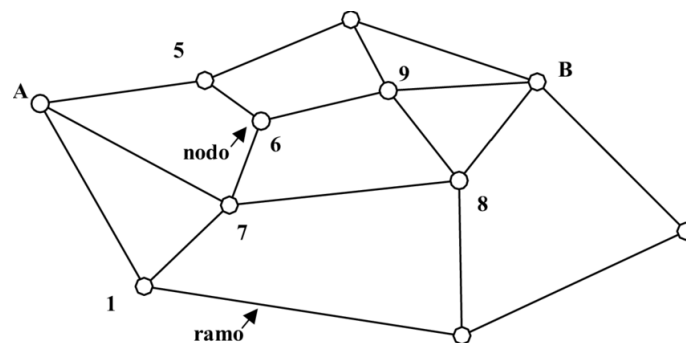


FIGURA 2.1: Grafo rappresentativo di una rete stradale

Fonte: Canale S., Leonardi S., Nicosia F., Le intersezioni stradali in ambito extraurbano, 2021, Figura 1, pag. 1

https://www.researchgate.net/publication/265047078_LE_INTERSEZIONI_STRADALI_IN_AMBITO_EXTRAURBANO

ad un programma di calcolo (nodi: istruzioni, arco: esiste se le istruzioni possono essere eseguite in successione) oppure una struttura di dati (nodi: dati semplici, archi: legami tra i diversi dati).

I nodi si possono considerare come dei punti che possono essere o meno connessi da alcuni segmenti o frecce. Un *path* tra due nodi X e Y è una sequenza di nodi che inizia con X e finisce con Y, nella quale ogni nodo è connesso all'altro attraverso un

arco. Due nodi sono adiacenti se sono connessi da un arco e due archi sono adiacenti se hanno un nodo in comune. Un nodo non collegato con altri è detto nodo isolato o nodo singolo, invece archi isolati non esistono. Si può parlare anche di archi multipli (esempio in Figura 2.2) e *loops* (esempio in Figura 2.3), i primi collegano gli stessi due nodi, mentre gli altri partono ed arrivano sempre nello stesso nodo.

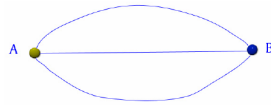


FIGURA 2.2: Archi multipli

Fonte: Desmatron, *Teoria dei Grafi*, 2004. Pag. 10, Figura 2.4

https://www.matematicamente.it/staticfiles/teoria/geometria/teoria_dei_grafi.pdf

https://www.matematicamente.it/staticfiles/teoria/geometria/teoria_dei_grafi.pdf

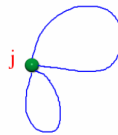


FIGURA 2.3: Loops

Fonte: Desmatron, *Teoria dei Grafi*, 2004. Pag. 10, Figura 2.5

https://www.matematicamente.it/staticfiles/teoria/geometria/teoria_dei_grafi.pdf

https://www.matematicamente.it/staticfiles/teoria/geometria/teoria_dei_grafi.pdf

Un grafo semplice non ha né archi multipli, né *loops*. Un grafo è completo se c'è un arco tra ogni coppia di nodi nel grafo. Il numero di nodi di un grafo costituisce il suo ordine, mentre la sua dimensione è data dal numero di archi che possiede.

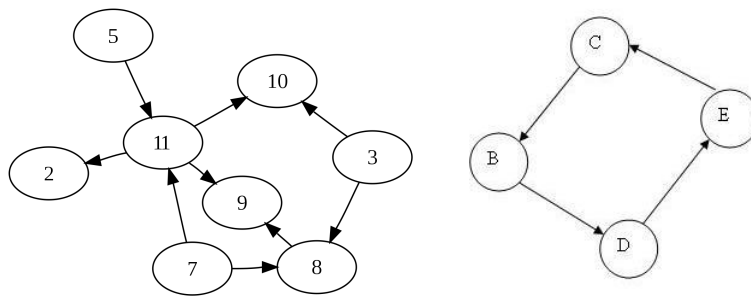


FIGURA 2.4: Grafo aciclico (sx) e grafo ciclico (dx)

Fonte 1 (grafo aciclico): https://it.wikipedia.org/wiki/Digrafo_aciclico,

Fonte 2 (grafo ciclico): http://ciropersico.altervista.org/studenti/strutture_dati/grafico/caratteristiche.htm

Un grafo è definito aciclico se non contiene cicli, cioè quando il nodo iniziale non è collegato a quello finale, il punto di partenza non coincide con quello di arrivo.

Non esiste dunque alcun *path* da un nodo a sè stesso. In caso contrario, si parla di grafo ciclico. Vediamo un esempio di grafo aciclico e uno di grafo ciclico nella Figura 2.4.

Gli archi possono essere orientati o meno. Un arco orientato esce da un nodo ed entra in un altro, con una direzione indicata da una freccia. Essi possono essere quindi percorsi in un solo verso. Passare dal nodo A al nodo B non è dunque la stessa cosa che passare dal nodo B al nodo A. Un grafo in cui tutti gli archi sono orientati è detto grafo orientato. In particolare, si ha una testa (o nodo di arrivo) e una coda (o nodo di partenza). Considerata una coppia (i, j) di nodi in un grafo orientato, i è detto predecessore di j e j è detto successore di i . Nella Figura 2.5 notiamo la differenza tra un grafo orientato e un grafo non orientato.

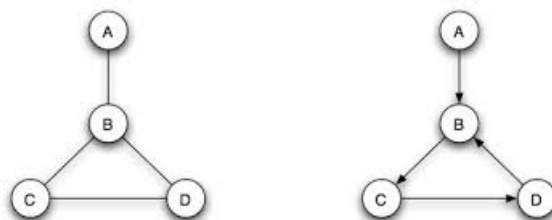


FIGURA 2.5: Differenza grafo orientato e non orientato

A sinistra un grafo non orientato, a destra un grafo orientato

Fonte: <http://www.di-srv.unisa.it/professori/lg/RS/SN-Grafi.pdf>, Pag. 1

Se gli archi sono orientati, in ogni grafo possiamo notare una sorta di gerarchia fra gli elementi che lo compongono. Il nodo da cui parte un arco orientato è chiamato genitore del nodo in cui l'arco arriva, mentre il nodo in cui l'arco entra è chiamato figlio del nodo da cui l'arco inizia. Se due nodi sono connessi da un *path* orientato, allora il primo nodo è un antenato di tutti i nodi del *path* e ogni nodo del *path* è un discendente del primo nodo. Per esempio, nel grafo di Figura 2.6, **1** è antenato e, in particolare, genitore di **2**, mentre **2** è figlio di **1** e quindi suo discendente.

Una freccia bidirezionale (con due teste) che connette due variabili in un grafo è spesso usata per indicare che le due variabili condividono uno o più antenati, ma gli antenati e le loro relazioni non sono state disegnate nel grafo. Sempre considerando la Figura 2.6, possiamo notare come tra **1** e **4** ci sia una freccia bidirezionale, ad indicare che questi due nodi hanno in comune degli antenati che non sono stati raffigurati.

Un arco non direzionale, cioè senza una freccia, è usato a volte per indicare che due variabili sono associate per motivi diversi dalla condivisione di un antenato o

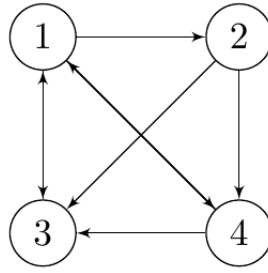


FIGURA 2.6: Esempio di grafo con freccia bidirezionale

Fonte: <https://tex.stackexchange.com/questions/304974/how-to-create-2-arrows-rather-than-double-headed-arrow-with-tikz>

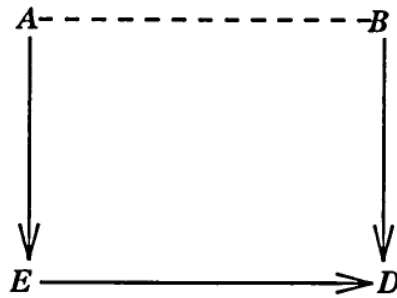


FIGURA 2.7: Esempio di grafo con arco non direzionale

Fonte: Greenland S, Pearl J, Robins JM., Causal diagrams for epidemiologic, in "Epidemiology", a. X, n. 1, Gennaio 1999, Figura 3, Pag. 3

dall'influenza reciproca. Ne vediamo un esempio in figura 2.7, dove A e B sono collegati da un arco tratteggiato, che non è in alcun modo orientato.

In questi grafi possiamo distinguere le catene oppure i forks.

Un esempio di catena è il seguente:

$$A \rightarrow B \rightarrow C.$$

Le frecce hanno lo stesso verso, una entra in B e l'altra esce da B. B media l'effetto di A su C.

Un esempio di fork invece, è il seguente:

$$A \leftarrow B \rightarrow C.$$

Le frecce hanno due versi differenti, escono entrambe da B, una verso A e una verso C. In questo caso, B è una causa comune ad A e C.

Esiste infine il fork invertito o collisore:

$$A \rightarrow B \leftarrow C.$$

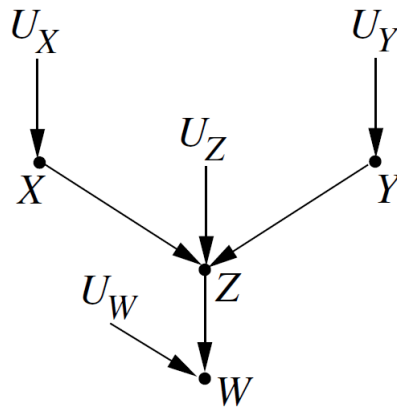


FIGURA 2.8: Esempio di collisore (Z)

Fonte: Pearl J., Glymour M., Jewell N., Causal inference in statistics a Primer, New York, John Wiley, 2016, Figura 2.4, Pag. 72.

Stavolta entrambe le frecce puntano su B. A e C sono indipendenti se B non è condizionato; l'esclusione di B blocca il percorso tra A e C.

In generale possiamo definire *backdoor path* da X ad Y un percorso che collega X a Y, se la punta della freccia va verso X. Un percorso collide in una variabile X se il *path* entra in X (dove arriva la punta della freccia) ed in questo caso X è detto collisore del *path*. Nella Figura 2.8, un percorso *backdoor* potrebbe essere quello da U_X a Z. Inoltre, sempre in questa figura, notiamo come Z sia un collisore, in quanto tutte le frecce dei suoi antenati puntano su di essa e quindi, collidono in Z.

Un collisore è infatti un nodo in cui due o più punte di freccia si incontrano e un percorso collide nel momento in cui si chiude con il collisore. Un percorso è *blocked* se ha uno o più collisori, altrimenti è *unblocked*. I collisori richiedono molta attenzione, perché condizionando inavvertitamente un collisore, nel tentativo di rimuovere il confondimento, si introduce effettivamente il confondimento. Inoltre, a cascata, quando si condiziona un collisore o i suoi discendenti, i genitori potrebbero diventare dipendenti, modificando così le relazioni di partenza tra le variabili in studio.

2.2 I DAG

L'acronimo DAG, come già accennato, deriva dalle parole inglesi "*directed acyclic graph*", tradotte in italiano in "grafo orientato aciclico". Si tratta quindi di un grafo con tutti gli archi direzionali e in cui il nodo iniziale non è connesso con quello finale.

I DAG hanno molte applicazioni oltre all'inferenza causale, ma noi ci focalizzeremo su quelli causali.

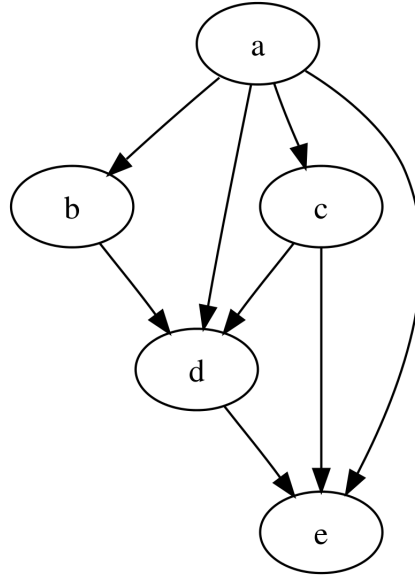


FIGURA 2.9: Esempio di DAG

Fonte: https://en.wikipedia.org/wiki/Directed_acyclic_graph

Si possono quindi chiamare anche diagrammi causali o modelli causali grafici. Nei diagrammi causali si adotta la convenzione che il tempo "scorre" da sinistra verso destra (o dall'alto verso il basso) e dunque una variabile posta più a sinistra (o più in alto) sarà temporaneamente precedente ad una posta più a destra (o più in basso). Per esempio, nella Figura 2.9, **a** è temporaneamente precedente al resto delle variabili.

Nei DAG causali, come accennato, ogni nodo è una variabile e ogni arco collega le variabili in relazioni causa-effetto. In particolare, una freccia da una variabile all'altra indica che la prima variabile causa la seconda e che il valore della prima variabile è parte della funzione che determina il valore della seconda. Unendo dunque l'aspetto grafico e quello teorico, possiamo dire che se una variabile X è figlia di una variabile Y , allora Y è una causa diretta di X . Inoltre, se X è un discendente di Y , allora Y è una causa potenziale di X .

Ogni diagramma causale è associato ad un *Structural Causal Model* (SCM), o modello causale strutturale, che aiuta a descrivere particolari caratteristiche del mondo e come queste possono interagire fra loro. Nello specifico, questi modelli descrivono in che modo vengono assegnati i valori alle variabili di nostro interesse.

I DAG sono utili perchè, anche se contengono meno informazioni degli SCM, consentono un'immediata conoscenza relativa alle relazioni causali, che, essendo di tipo qualitativo, emergono con maggiore evidenza osservando un grafico, piuttosto che guardando un modello SCM. In aggiunta, permettono di esprimere in modo molto efficiente le distribuzioni congiunte. Infatti, rappresentare questo tipo di distribuzioni attraverso tabelle, oppure specificandone il modello che le descrive, non è sempre immediato, facile e soprattutto possibile.

Le nuove teorie sui diagrammi per l'inferenza causale sono cresciute con discipline come la computer science e l'artificial intelligence. I DAG servono per rappresentare le chiavi dei concetti causali. Nello specifico, la presenza di una freccia che va da una variabile ad un'altra, indica che si conosce un effetto causale diretto per almeno un individuo. Al contrario, la mancanza di una freccia significa che conosciamo che le due variabili non sono connesse da alcun tipo di effetto causale, per nessun individuo della popolazione. Inoltre, un diagramma causale standard non distingue se una freccia rappresenta un effetto protettivo oppure un effetto dannoso, nè come due cause interagiscono fra di loro quando una variabile ha più cause.

2.3 D-separation

Nella maggior parte dei modelli grafici alcune variabili hanno più percorsi che le connettono e ogni percorso attraverserà una varietà di catene, *fork* e collisori. Esiste un criterio o processo che può essere applicato a un modello causale grafico di qualsiasi complessità per poter prevedere le dipendenze condivise da tutti gli insiemi di dati generati da quel grafico. Questo criterio è il processo della *d-separation* ($d = \text{"directional"}$). Esso ci permette di determinare, per qualsiasi coppia di nodi, se i nodi sono d-connessi o d-separati. Nel primo caso esiste un percorso che collega i due nodi, mentre nel secondo non esiste alcun percorso tra i nodi. Due nodi sono d-separati se le variabili che essi rappresentano sono definitivamente indipendenti; due nodi sono d-connessi se, molto probabilmente, sono dipendenti.

Due nodi X e Y sono d-separati se ogni percorso tra loro è *blocked*; se anche solo un percorso tra loro è *unblocked*, X e Y sono d-connessi.

Ci sono alcuni tipi di nodi che possono bloccare un percorso, in base al fatto che si stia considerando una *d-separation* condizionata o meno. Se però, nessuna variabile è condizionata, gli unici nodi che possono bloccare un percorso sono i collisori. Così, se ogni percorso tra due nodi ha un collisore, allora i due nodi

considerati non possono essere incondizionatamente dipendenti, ma devono essere marginalmente indipendenti.

Se invece si condizionano un insieme di nodi Z , allora i successivi tipi di nodi che possono bloccare un percorso possono essere:

- Un collisore che non è condizionato (cioè non condizionato in Z) e che non ha discendenti in Z ;
- Una catena o un fork il cui nodo centrale è in Z .

Se Z blocca ogni percorso tra due nodi, allora i due nodi sono d -separati, condizionati a Z , e quindi indipendenti in base a Z .

Se si condiziona rispetto ad un collisore, esso "sblocca" il percorso che unisce due variabili e dunque il percorso tra *blocked* diventa *unblocked* e due variabili che prima erano d -separate, ora saranno d -connesse, condizionate al collisore che è stato "sbloccato".

Grazie alla definizione di *d-separation* si possono dunque osservare anche grafi complessi e, indipendentemente dal modello che lega le variabili, la *d-separation* ci permette di identificare sempre le indipendenze nei dati generati dal modello.

Sappiamo che se ogni condizione di *d-separation* nel modello corrisponde ad un'indipendenza condizionale nei dati, allora nessun ulteriore test può confutare il modello.

Ci sono altri metodi per testare l'idoneità di un modello. Il modo più utilizzato in statistica è quello dei test d'ipotesi rispetto all'intero modello. Tuttavia, se il modello non è totalmente specificato, abbiamo bisogno di stimare i parametri prima di valutarne l'idoneità. Questo può essere fatto solo quando si assume un modello lineare e normale. Anche con questo metodo, però, ci sono diversi problemi. Per prima cosa, se alcuni parametri non possono ugualmente essere stimati, allora la distribuzione congiunta non si può stimare e il modello non può essere testato. Questo avviene quando alcuni errori sono correlati fra loro oppure quando alcune variabili non sono osservate.

In secondo luogo, questa procedura testa il modello globalmente. Se si scopre che il modello non spiega bene i dati, non ci sono modi per determinarne il motivo e dunque non è possibile specificare quali nodi dovrebbero essere rimossi o aggiunti per migliorare il *fit*.

Per ultimo, quando si testa un modello globalmente, il numero di variabili considerate può essere elevato, e se ci sono rumori nei dati e/o variazioni dovute al campionamento associati a ciascuna variabile, allora il test non sarà affidabile.

La *d-separation* ha molti vantaggi rispetto a questo metodo di test globale. Per prima cosa, non è parametrico, cioè non si basa su funzioni specifiche che collegano le variabili, ma utilizza solo il grafico del modello in questione. Inoltre, testa i modelli a livello locale, invece che globale. Ciò permette di identificare aree specifiche in cui il modello ipotizzato è difettoso e di ripararle, piuttosto che iniziare da zero su un modello completamente nuovo. Per questo motivo, se per qualsiasi ragione non possiamo identificare il coefficiente in un'area del modello, possiamo lo stesso ottenere alcune informazioni incomplete sul resto del modello.

2.4 ESC-DAGs

I DAG sono quindi strumenti molto importanti per l'analisi, ma il loro utilizzo è spesso problematico. Talvolta sono semplicistici, vengono modificati per adattarsi ai dati disponibili e non vengono presentati quasi mai negli studi che li utilizzano. In generale, vi è una mancanza di “linee guida” per la costruzione e l'uso di DAG. Tuttavia, esistono tre punti condivisi da tutti coloro che costruiscono questi tipi di grafi:

- La teoria e la conoscenza di base dovrebbero avere un ruolo fondamentale nella costruzione dei DAG.
- È molto più importante il fatto che due nodi non siano connessi fra loro, piuttosto che includere un nodo nel DAG. Questo implica che i ricercatori dovrebbero “lavorare a ritroso” da un DAG “satturo” (modello teorico e ipotetico), in cui tutte le variabili sono interconnesse, ed eliminare solo le connessioni ritenute impossibili.
- I metodi basati sui dati, come la selezione *stepwise*, sono molto utilizzati e possono offrire valore aggiunto alla costruzione di DAG, tuttavia possono anche indurre *bias* aggiustando erroneamente per mediatori o collisori.

Ferguson e i suoi colleghi propongono una revisione metodologica per approcciarsi alla costruzione dei DAG, chiamata “*Evidence Synthesis for Constructing Directed Acyclic Graphs*” e riassunta in ESC-DAGs. Gli ESC-DAGs sistemizzano il modo in cui la conoscenza di base viene utilizzata per determinare quali variabili e connessioni tra le variabili sono incluse (lavorando a ritroso da un DAG saturo). Questo tipo di approccio fa leva sulla letteratura empirica, prima traducendo i risultati empirici in DAG e poi sintetizzandoli in uno o più DAG “integrati”. Tale

metodo viene applicato agli studi identificati da una ricerca in letteratura. Idealmente si tratta di una revisione sistematica o una revisione di revisioni sistematiche. Tuttavia, poiché gli ESC-DAGs potrebbero produrre DAG molto complessi con dozzine di variabili, le ricerche sistematiche dovrebbero essere limitate alle relazioni focali della domanda di ricerca.

I tre principali processi degli ESC-DAGs sono “mappatura”, “traduzione” e “integrazione”. Nello specifico, ogni studio identificato da una ricerca in letteratura (svolta esaminando numerosi articoli su motori di ricerca di letteratura scientifica biomedica) passa attraverso questi tre processi fondamentali.

Il processo di mappatura produce un DAG che rappresenta la conclusione dello studio. Questo DAG è “saturato”, con un arco per ogni coppia di nodi. Chiamiamo questo output “grafico implicito” (IG) dello studio, che funge da modello strutturale trasparente per la traduzione in un DAG.

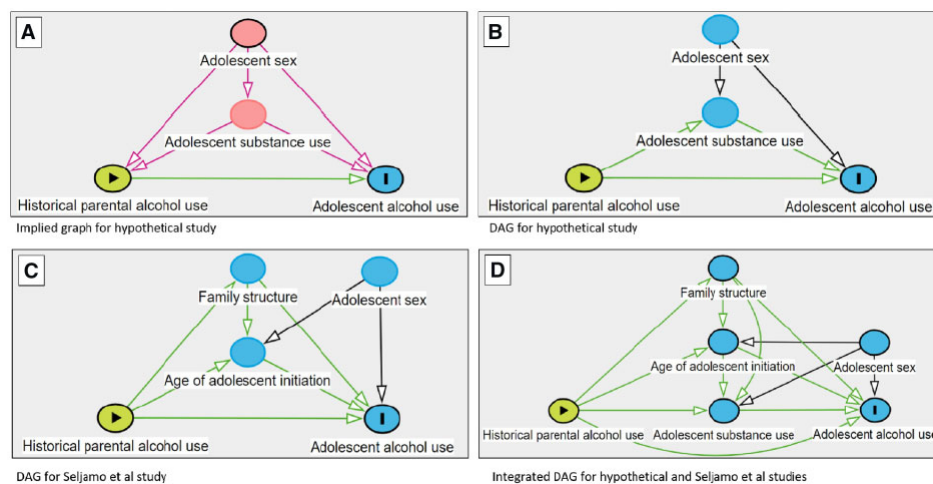


FIGURA 2.10: Processo di traduzione di un ESC-DAG

Fonte: Ferguson K., McCann M., Katikireddi S., Thomson H., Green M., Smith D, Lewsey J, Evidence synthesis for constructing directed acyclic graphs (ESC-DAGs): a novel and systematic method for building directed acyclic graphs, in “International Journal of Epidemiology” a. XLIX, n. 1, 2020, Figura 1, Pag. 4.

La traduzione, che possiamo osservare in Figura 2.10, implica invece la valutazione delle caratteristiche causali di ciascuna connessione (archi diretti) di questo DAG. Gli archi diretti sono compilati in un indice. Questo step crea il DAG per lo studio.

La fase di integrazione combina gli archi diretti dall’indice in un diagramma. Questo step si suddivide in due: la fase di sintesi e quella di ricombinazione. La prima serve per combinare i DAG tradotti in uno solo, sintetizzando tutti i bordi

diretti indicizzati. La seconda, invece, combina i nodi per ragioni pratiche, ovvero per ridurre la complessità, oppure sostanziali, cioè per stabilire la coerenza. L'output finale consiste in uno o più DAG integrati. Il processo decisionale è registrato in un "registro delle decisioni", da fornire come appendice agli ESC-DAG.

Come gli autori si augurano, essendo un metodo molto semplice, dovrebbe essere un ulteriore stimolo affinché i ricercatori utilizzino i DAG per migliorare la ricerca sulla salute della popolazione.

2.5 DAG e confondimento

Dato che i diagrammi causali sono di aiuto nel rappresentare diverse fonti di associazione, possiamo usarli per classificare il *bias* sistematico in base alla sua fonte.

La struttura del confondimento, il *bias* dovuto alle cause comuni tra trattamento e *outcome*, può essere rappresentato usando i diagrammi causali. Vediamone un esempio in figura 2.11.

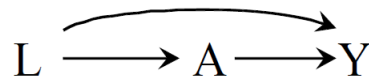


FIGURA 2.11: Esempio di DAG in cui L è causa comune del trattamento e dell'outcome

Fonte: Hernan M., Robins J., Causal Inference: What If, 2020, Figura 7.1, Pag. 92.
<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

In questo DAG A è il trattamento, Y è l'*outcome* e L è la causa comune tra *outcome* e trattamento. Questo DAG mostra due fonti di associazione tra il trattamento e l'*outcome*:

- Il percorso $A \rightarrow Y$ che rappresenta l'effetto causale di A su Y;
- Il percorso $A \leftarrow L \rightarrow Y$ tra A e Y che include la causa comune L. Questo percorso che collega A e Y attraverso la causa comune L è un esempio di *backdoor path*.

Se la causa comune non esistesse, allora il solo percorso tra trattamento e *outcome* sarebbe $A \rightarrow Y$ e dunque l'intera associazione tra A e Y sarebbe dovuta all'effetto causale di A su Y. Il rapporto associativo tra rischi coincide con il rapporto causale

tra rischi e dunque l'associazione coincide con la causa. Riprendendo quanto detto nel paragrafo d'introduzione al confondimento, la presenza di una causa comune crea una fonte aggiuntiva di associazione tra il trattamento e l'*outcome*, che noi chiamiamo confondimento per l'effetto di A su Y. A causa del confondimento il rapporto associativo tra rischi non è uguale a quello causale e così l'associazione non coincide con la causa.

Il confondimento può essere messo in relazione con la scambiabilità. In presenza di scambiabilità, come in un esperimento marginalmente randomizzato in cui tutti gli individui hanno la stessa probabilità di ricevere il trattamento, l'effetto causale medio può essere identificato senza aggiustamento per nessuna variabile. Quando non vi è scambiabilità, ma è presente la scambiabilità condizionata dall'effetto comune tra il trattamento e l'*outcome*, come in un esperimento condizionalmente randomizzato nel quale la probabilità di ricevere il trattamento varia in base al valore della causa comune L, l'effetto causale medio può essere identificato. Tuttavia, l'identificazione dell'effetto causale nella popolazione richiede l'aggiustamento per la variabile L, standardizzando o pesando. La scambiabilità condizionale, però, permette di identificare gli effetti causali condizionati per ogni valore di L, stratificando.

Nella pratica, se crediamo che ci sia confondimento, ci chiediamo se si possa determinare, se esiste, un insieme di covariate misurate L, per le quali vale la scambiabilità condizionale. Rispondere a questa domanda è possibile se si conosce il DAG causale che ha generato i dati. Ci sono due diversi approcci per sfruttare i DAG e rispondere quindi alla domanda. Il primo approccio consiste nell'applicare il criterio *backdoor* al DAG causale, il secondo, invece, richiede di trasformare il DAG in uno SWIG (altro tipo di diagramma causale, dove i nodi rappresentano le variabili controfattuali). L'approccio degli SWIG è più diretto, ma richiede anche molti più macchinari e non lo approfondiremo. Ci soffermiamo invece sul criterio *backdoor* applicato ai DAG.

Un set di covariate L soddisfa il criterio *backdoor* se tutti i percorsi *backdoor* tra A e Y sono bloccati condizionando su L e L non contiene variabili che sono discendenti del trattamento A. Questo tipo di criterio si utilizza quando ci sono dei DAG complessi e ci permette di capire se esiste il confondimento e se c'è un insieme di variabili misurate sufficiente per identificare l'effetto causale di X su Y.

Per determinare quando il confondimento esiste, utilizziamo due semplici passaggi. Per prima cosa, eliminiamo tutte le frecce che partono dall'esposizione, poi, nel grafico rimanente, determiniamo dove c'è qualsiasi percorso *backdoor* sbloccato

che parte dall'esposizione e arriva all'*outcome*. Se esiste questo tipo di percorso, la relazione causale è confusa dall'effetto delle altre variabili, se invece non esiste, non c'è confondimento.

È inoltre possibile utilizzare un algoritmo, chiamato *backdoor test for sufficiency*, per controllare dove un set di variabili è sufficiente per l'aggiustamento. I primi due passaggi da seguire, dato un sottoinsieme di variabili che non contiene discendenti dell'esposizione o dell'*outcome*, sono:

- Eliminare tutte le frecce che partono dall'esposizione;
- Disegnare archi non orientati per collegare qualsiasi coppia di variabili che condivide un figlio che è all'interno del sottoinsieme o ha un discendente nel sottoinsieme.

Infine, nel nuovo grafico, determinare se c'è qualsiasi percorso *backdoor* non bloccato che parte dall'esposizione e arriva all'*outcome*, che evita di passare attraverso qualsiasi nodo nell'insieme dei fattori di stratificazione. Se non è stato trovato nessun percorso, il confondimento è controllato dai fattori proposti, se invece c'è un percorso, la stratificazione su questi fattori non è sufficiente per rimuovere tutto il confondimento.

Mettiamo ora in relazione il criterio *backdoor* con il confondimento. Il criterio *backdoor* è soddisfatto in queste due situazioni:

- Nessuna causa comune tra il trattamento e l'*outcome*;
- Cause comuni tra il trattamento e l'*outcome*, ma il sottoinsieme L di non discendenti di A misurati è sufficiente per bloccare tutti i percorsi *backdoor*

Nel primo caso, non essendoci alcuna causa comune, non esiste nessun percorso *backdoor* che deve essere bloccato, quindi l'insieme di variabili che soddisfano il criterio *backdoor* è vuoto e si dice che non c'è confondimento. Nel secondo caso, si dice che c'è confondimento, ma non c'è confondimento residuo la cui eliminazione richiederebbe un aggiustamento per variabili non misurate. Si dice quindi che non c'è confondimento non misurato.

Il criterio *backdoor* non risponde a domande riguardanti la forza e la direzione del confondimento. È logicamente possibile che alcuni percorsi *backdoor* non bloccati siano deboli e così inducano un piccolo *bias*, o che percorsi *backdoor* molto forti inducano un *bias* in una direzione opposta, costituendo un debole *bias* netto. A

causa del confondimento non misurato è importante considerare la direzione attesa e la grandezza del *bias*.

Un ulteriore caso è esemplificato in Figura 2.12.

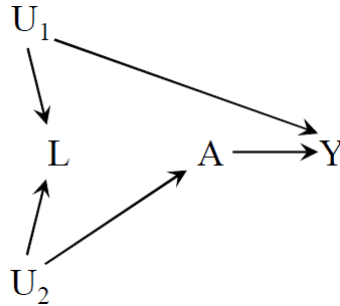


FIGURA 2.12: Esempio di DAG in cui non vi sono cause comuni tra trattamento e outcome

Fonte: Hernan M., Robins J., Causal Inference: What If, 2020, Figura 7.4, Pag. 96.

<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

In questo DAG non ci sono cause comuni tra il trattamento A e l'*outcome* Y e quindi non c'è confondimento. Il percorso *backdoor* tra A e Y attraverso L è bloccato perché L è un collisore del percorso. Così, l'associazione tra A e Y è dovuta all'effetto tra A e Y e perciò l'associazione è di tipo causale. Aggiustare per L però, potrebbe indurre *bias* poichè condizionando per un collisore si aprirebbe il percorso *backdoor* tra A e Y , che era precedentemente bloccato dal collisore stesso. In questo caso l'associazione non sarebbe più di causalità. Questo è un esempio in cui vale la scambiabilità incondizionata, ma non quella condizionale: viene identificato l'effetto causale medio, ma non gli effetti causali condizionati entro i livelli di L . Chiamiamo distorsione da selezione la distorsione risultante dall'effetto condizionale, perché deriva dalla selezione sull'effetto comune L di due variabili marginalmente indipendenti (U_1 e U_2), una delle quali è associata ad A e l'altra a Y .

Se al DAG precedente aggiungessimo un percorso $L \rightarrow A$, il *bias* sarebbe intrattabile in quanto la presenza della freccia tra L e A creerebbe un percorso *backdoor* aperto, generando un confondimento. Condizionando su L , però, sarebbe bloccato il percorso *backdoor*, ma nello stesso tempo si aprirebbe un nuovo percorso *backdoor* nel quale L è un collisore.

Il confondimento può dunque essere visto come il *bias* dovuto al percorso *backdoor* aperto tra A e Y , perciò è qualsiasi *bias* sistematico che potrebbe essere eliminato dall'assegnazione casuale di A . Questa particolare definizione implica che l'esistenza di confondimento dipende dal metodo di analisi.

Una volta che il DAG è conosciuto, si può semplicemente applicare il criterio *backdoor* per determinare quali variabili devono essere aggiustate. Il tradizionale approccio per gestire il confondimento, basato principalmente su associazioni osservate piuttosto che su una conoscenza causale precedente, etichetta prima le variabili che soddisfano determinate condizioni associative come fattori di confondimento e poi impone che questi cosiddetti fattori di confondimento siano adeguati all'analisi. Si dice che è presente confondimento quando la stima aggiustata differisce dalla stima non aggiustata. In particolare, come già illustrato, un confondente è una variabile che soddisfa le tre seguenti condizioni:

- È associata al trattamento;
- È associata al risultato condizionato al trattamento;
- Non si trova su un percorso causale tra il trattamento e l'esito.

Questo approccio tradizionale, però, può portare ad un adattamento inadeguato. Si dice quindi che queste tre condizioni siano condizioni necessarie, ma non sufficienti. Se consideriamo sempre la figura 2.12, con l'approccio tradizionale diremmo che c'è associazione con il trattamento, che c'è associazione con l'*outcome* condizionato al trattamento e che L non giace sul percorso tra trattamento e *outcome*. Dunque, si potrebbe concludere che L è un confondente da aggiustare anche in assenza di confondimento. Ma come già visto, aggiustare per L causerebbe la presenza di *bias*. Questo esempio ci illustra come il criterio di associazione (o statistico) non sia sufficiente per capire se è presente confondimento. Per eliminare il problema, si potrebbe sostituire la condizione associativa "è associato all'*outcome* condizionato al trattamento" con la condizione strutturale "è una causa dell'*outcome*".

L'approccio tradizionale porta gli investigatori ad adeguarsi alle variabili, quando l'aggiustamento è dannoso. Il problema legato a questo approccio deriva dalla definizione dei confondenti in assenza di una conoscenza causale sufficiente sulle fonti del confondimento e quindi impone l'adeguamento per quei cosiddetti confondenti. In particolare, se le stime aggiustate e quelle non aggiustate differiscono, l'approccio tradizionale dichiara l'esistenza di elementi confondenti.

Al contrario, un approccio strutturale inizia identificando esplicitamente le fonti di confusione, le cause comuni del trattamento e il risultato che, se fossero tutti misurati, sarebbero sufficienti per aggiustare il fattore di confusione, e quindi identifica un insieme sufficiente di variabili di aggiustamento. Questo approccio chiarifica che

includere una particolare variabile in un certo insieme dipende dalle variabili già incluse nell'insieme considerato. Dato un DAG causale, il confondimento è un concetto assoluto dove il confondente è un concetto relativo.

Un approccio strutturale al confondimento enfatizza che l'inferenza causale derivata dai dati osservazionali richiede una conoscenza causale a priori. Questa conoscenza causale è riassunta in un DAG causale che codifica le credenze dei ricercatori o le assunzioni circa la rete causale. Certamente, non c'è garanzia che questo DAG sia corretto e quindi è possibile che l'insieme di variabili scelto dai ricercatori non riesca ad eliminare i fattori di confondimento o introduca *bias* di selezione. Tuttavia, l'approccio strutturale ha due importanti vantaggi: previene le incongruenze tra credenze e azioni ed esplicita le assunzioni dei ricercatori relative al confondimento, così da poter essere esplicitamente criticate da altri investigatori.

In assenza di randomizzazione, l'inferenza causale si basa sull'assunzione non verificabile che abbiamo misurato un insieme di variabili L , che è un insieme sufficiente per aggiustare il confondente, cioè un insieme di non discendenti del trattamento A che include abbastanza variabili per bloccare tutti i percorsi *backdoor* da A a Y . Sotto questa ipotesi di scambiabilità condizionale, data L , la standardizzazione e la ponderazione (IP) possono essere utilizzate per calcolare l'effetto causale medio nella popolazione. Tuttavia, la standardizzazione e la ponderazione non sono gli unici metodi possibili per aggiustare il confondimento negli studi osservazionali. I metodi che aggiustano per alcuni confondenti L possono essere classificati in due categorie:

- Metodi G: standardizzazione, ponderazione e la stima-g. Questi metodi (dove g sta per "generalizzati") spiegano la scambiabilità condizionale dato L per stimare l'effetto causale di A su Y nell'intera popolazione o in un sottogruppo della popolazione.
- Metodi basati sulla stratificazione: stratificazione (includendo restrizioni) e *matching*. Questi metodi spiegano la scambiabilità condizionale dato L per stimare l'associazione tra A e Y in un sottogruppo definito da L .

I metodi G simulano l'associazione A - Y nella popolazione nel caso in cui non esistano percorsi *backdoor* che coinvolgono le variabili misurate L . In questi metodi si sceglie di aggiustare per il confondimento, perché i metodi basati sulla stratificazione potrebbero portare al *selection bias*.

Tutti i metodi citati richiedono la scambiabilità condizionale data da L . Il confondimento, però, può essere talvolta gestito con metodi che non richiedono scambiabilità

condizionale. Alcuni esempi sono la differenza nelle differenze, la stima delle variabili strumentali, il criterio della porta d'ingresso e altri. Sfortunatamente, questi metodi richiedono ipotesi alternative che, come la scambiabilità condizionale, non sono verificabili. Pertanto, in pratica, la validità delle stime degli effetti risultanti non è garantita. Inoltre, questi metodi non possono essere generalmente utilizzati per domande causali che coinvolgono trattamenti che variano nel tempo. Di conseguenza, questi metodi sono esclusi dalla considerazione per molti problemi di ricerca.

Il raggiungimento della scambiabilità condizionale può essere un obiettivo non realistico in molti studi osservazionali, ma la conoscenza di esperti sulla struttura causale può essere utilizzata per avvicinarsi il più possibile a tale obiettivo. Pertanto, negli studi osservazionali, i ricercatori misurano molte variabili (non discendenti dal trattamento) nel tentativo di garantire che il trattato e il non trattato siano condizionatamente scambiabili. La speranza è che, anche se possono esistere cause comuni, le variabili misurate siano sufficienti per bloccare tutti i percorsi *backdoor*. Tuttavia, non vi è alcuna garanzia che questo tentativo abbia successo, il che rende l'inferenza causale dai dati osservati un'impresa rischiosa.

La conoscenza degli esperti può essere inoltre utilizzata per evitare aggiustamenti per variabili che possono introdurre bias. Per lo meno, gli investigatori dovrebbero generalmente evitare aggiustamenti per le variabili influenzate dal trattamento o dal risultato.

In generale, dunque, non è necessario includere tutte le cause delle variabili nel DAG, ma è importante includere tutte le variabili chiave, altrimenti si giungerebbe a interpretazioni errate.

2.6 DAG e selection bias

Dopo essere stato citato numerose volte, è giunto il momento di spiegare cosa sia il *selection bias*. Questo tipo di *bias* si verifica in seguito al processo di selezione degli individui nell'analisi. A differenza del confondimento, non è dovuto alla presenza di cause comuni di trattamento e di esito e può sorgere sia in esperimenti randomizzati che in studi osservazionali. Come il confondimento, tuttavia, consiste in una mancanza di scambiabilità tra il trattato e il non trattato.

Ci focalizziamo in particolare sul bias che sorgerebbe anche se il trattamento avesse un effetto nullo sull'*outcome*, si parla cioè di “*selection bias* sotto il nullo”.

La struttura del *selection bias* può essere rappresentata da un DAG come quello in Figura 2.13.

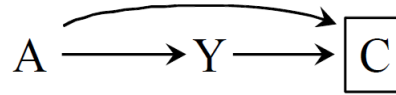


FIGURA 2.13: Primo esempio di DAG con selection bias

Fonte: Hernan M., Robins J., Causal Inference: What If, 2020, Figura 8.1, Pag. 108.

<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

In questo DAG A è il trattamento, Y l'*outcome* e C l'effetto comune. Questo DAG ci mostra due fonti di associazione tra il trattamento e l'*outcome*. La prima è il percorso aperto $A \rightarrow Y$, che rappresenta l'effetto causale di A su Y, la seconda è il percorso aperto $A \rightarrow C \leftarrow Y$, che collega A e Y attraverso il loro effetto comune C. Un'analisi condizionata su C porterà generalmente ad un'associazione tra A e Y. Ci si riferisce a questa associazione indotta tra il trattamento e l'*outcome* come *selection bias* dovuto al condizionamento su C. A causa del *bias*, il rapporto associativo tra rischi non sarà uguale al rapporto causale tra rischi e dunque l'associazione non coincide con la causa. Se invece non si condizionasse per C, allora i rapporti tra rischi coinciderebbero e l'associazione corrisponderebbe con la causa.

Si riportano ora altri esempi di DAG in cui vi è *selection bias*.

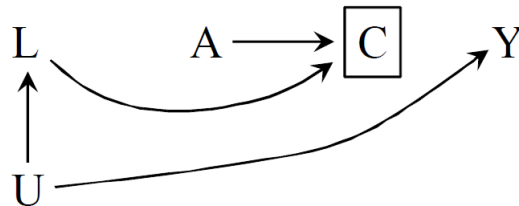


FIGURA 2.14: Secondo esempio di DAG con selection bias

Fonte: Hernan M., Robins J., Causal Inference: What If, 2020, Figura 8.3, Pag. 109.

<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

Nel DAG in Figura 2.14, vi è *bias* perché è stato condizionato per un effetto comune tra trattamento e *outcome*. Questo *bias* si verifica indipendentemente dal fatto che sia presente una freccia da A a Y, è cioè un *bias* di selezione sotto lo zero.

Nel DAG in Figura 2.15, il *selection bias* deriva dal condizionamento su C, che è un effetto comune del trattamento e una causa dell'*outcome*, piuttosto che dell'*outcome* stesso.

Nel DAG in Figura 2.16 vediamo la situazione più nello specifico. Si ha un trattamento A che ha effetti diretti sui sintomi L. Restringere lo studio ponendo $C=0$ (in

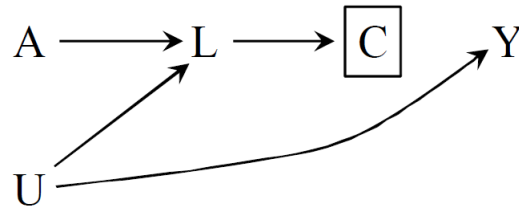


FIGURA 2.15: Terzo esempio di DAG con selection bias

Fonte: Hernan M., Robins J., Causal Inference: What If, 2020, Figura 8.4, Pag. 109.

<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

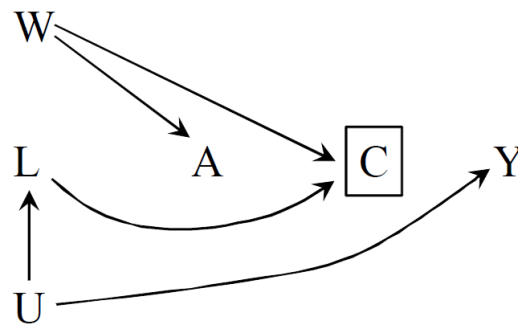


FIGURA 2.16: Quarto esempio di DAG con selection bias

Fonte: Hernan M., Robins J., Causal Inference: What If, 2020, Figura 8.5, Pag. 109.

<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

questo particolare caso si tratta degli individui incensurati, cioè coloro che rimangono nello studio e non sono quindi persi al *follow up* o per altri motivi) implica nuovamente il condizionamento dell'effetto comune C su A e U; si introduce quindi un'associazione tra il trattamento e l'*outcome*.

In tutti i casi, il *bias* è il risultato della selezione su un effetto comune di altre due variabili nel DAG, cioè di un collisore. Useremo il termine *selection bias* per riferirci a tutti i *bias* che derivano dal condizionamento su un effetto comune di due variabili, una delle quali è il trattamento o una causa del trattamento e l'altra è l'*outcome* o una causa dell'*outcome*.

Questo tipo di *bias* si può verificare negli studi retrospettivi, in cui i dati sul trattamento A sono raccolti dopo il verificarsi dell'*outcome* Y e negli studi prospettici, in cui i dati sul trattamento A sono raccolti prima del verificarsi dell'*outcome* Y. Inoltre, il *selection bias* si può verificare sia in studi osservazionali che in studi sperimentali. Gli individui, sia negli esperimenti randomizzati che negli studi osservazionali, possono essere persi al *follow-up* o abbandonare lo studio prima che il loro esito sia accertato.

In generale, se la randomizzazione protegge contro il confondimento, ma non contro

il *selection bias*, dall'altra parte la selezione avviene dopo la randomizzazione. Dall'altra parte, negli esperimenti randomizzati nessun *bias* nasce dalla selezione nello studio prima che sia assegnato il trattamento.

Gli statistici e gli economisti spesso usano il termine “*selection bias*” per riferirsi sia al confondimento che al *selection bias* vero e proprio, perché in entrambi i casi il *bias* è dovuto alla selezione. Nel confondimento si tratta della selezione degli individui in un trattamento, mentre nel *selection bias* della selezione degli individui all'interno dell'analisi.

A volte questo tipo di *bias* può essere evitato utilizzando un adeguato disegno. Tuttavia in presenza di perdite al *follow up*, oppure quando viene svolta un'autoselezione e in generale ci sono dati mancanti, non si può intervenire e si verifica comunque il *selection bias*, che deve essere esplicitamente corretto nell'analisi. Talvolta questa correzione può essere ottenuta tramite ponderazione, che si basa sull'assegnazione di un peso a ciascun individuo selezionato, in modo che l'individuo conti nell'analisi non solo per se stesso, ma anche per quelli come lui, cioè con gli stessi valori di L e A , che non sono stati selezionati. Il peso è l'inverso della sua probabilità di selezione.

Qualcuno potrebbe dire che la ponderazione non è necessaria per aggiustare il *selection bias* e che si potrebbe procedere con la standardizzazione. Tuttavia, questo procedimento a volte funziona, come nel caso del DAG in figura 2.16; altre volte, invece, non funziona, come nella figura 2.17.

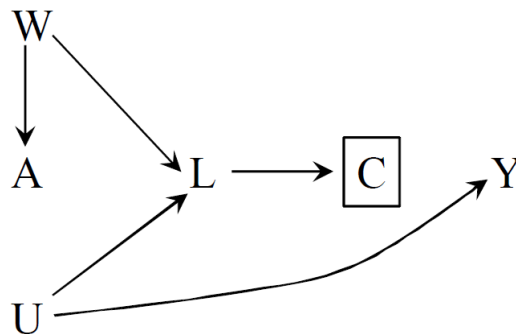


FIGURA 2.17: Esempio di DAG in cui non funziona la standardizzazione
 Fonte: Hernan M., Robins J., Causal Inference: What If, 2020, Figura 8.6, Pag. 109.
<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>

2.7 Interaction DAGs (IDAGs)

Come già detto, i DAG sono un ottimo aiuto per i ricercatori che provano a capire la natura delle relazioni causali e le conseguenze del condizionamento rispetto a diverse variabili, tuttavia, essi non considerano la possibile interazione presente tra le variabili. Per questo motivo, sono stati creati gli IDAG ($I = \text{"interaction"}$), che includono un nodo per un effetto causale e non per un *outcome*. Attraverso questi grafi si possono spiegare le interazioni confuse, quelle totali, dirette e indirette ed essi ci permettono di distinguere tra i meccanismi causali o non, dietro alle variazioni degli effetti. Inoltre, possono essere usati anche per chiarire i meccanismi che compromettono la generalizzabilità e per determinare quali variabili tenere in considerazione per rendere i risultati validi per la popolazione target.

I DAG, come visto nei paragrafi precedenti, possono illustrare concetti come il confondimento e il *selection bias*, distinguere effetti totali, diretti ed indiretti e essere usati per determinare quali variabili condizionano l'analisi empirica. Sono non parametrici e dunque non sono rilevanti per la costruzione del grafico se le variabili interagiscono fra di loro. Gli IDAG invece, sono intuitivi e fondati sulla teoria dell'inferenza causale, nello stesso modo dei DAG, ma ci permettono anche di comprendere come le diverse variabili influenzano l'effetto dell'*outcome*.

Il concetto di interazione, già accennato precedentemente, si riferisce ad alcuni effetti congiunti. L'idea generale di interazione è che l'effetto di una variabile dipende dal livello di un'altra variabile presente nel modello. La definizione di interazione è spesso espressa attraverso un *outcome* potenziale. Come sappiamo, nei modelli causali strutturali un *outcome* potenziale, o controfattuale, è un risultato che, per un insieme di fattori di fondo predeterminati, prevale quando si costringe una variabile (o più) nel modello ad assumere valori particolari.

In questo nuovo tipo di grafi denoteremo un effetto causale di A su Y con ΔY_A . Nel momento in cui ΔY_A è una variabile che dipende dalle altre variabili in maniera causale, può essere inclusa nel grafo causale. Se esiste un'interazione tra qualche variabile e A, si esprimerà tale rapporto con una freccia diretta da questa variabile a ΔY_A . Nella Figura 2.18 sono presenti le variabili X (genotipo) e A (chirurgia bariatrica), che influenzano Y (perdita di peso), vi è inoltre presenza di interazione. Dalla freccia che collega X a ΔY_A si può effettivamente notare la presenza della relazione tra X e A. L'effetto di A è modificato da Q (colore dei capelli), ma non c'è interazione tra A e Q. Non vi è infatti alcuna freccia che collega le due variabili (A e Q).

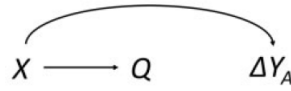


FIGURA 2.18: Esempio di IDAG

Fonte: Nilsson A., Bonander C., Stromberg U., Bjork J., A directed acyclic graph for interactions, in “International Journal of Epidemiology”, a. 00, n. 00, 2020, Figura 2.B, Pag. 4.

L’IDAG è dunque abbastanza simile allo standard DAG, ad eccezione del nodo dell’*outcome*, che viene sostituito da un nodo che rappresenta l’effetto causale, mentre il nodo che rappresenta il trattamento A non è incluso. Questi grafi sono solitamente raffigurati basandosi su letterature precedenti (si tratta di articoli pubblicati precedentemente, riguardanti il fenomeno in studio), che nel caso degli IDAG devono contenere evidenze sulla presenza di interazioni tra i trattamenti studiati.

Quando è presente un’interazione, questa potrebbe dipendere dalla scala con cui gli effetti casuali sono stati misurati, che può essere ad esempio di tipo additivo e dunque basata su differenze, oppure di tipo moltiplicativo e perciò basata su rapporti. Due variabili che influenzano un risultato, infatti, interagiranno sempre su alcune scale. L’aspetto del grafo (IDAG) dipende quindi dalla scala scelta e, per questo, alcune variabili possono puntare verso ΔY_A in alcune versioni dell’IDAG, ma non in altre.

Per semplicità si può assumere che non ci siano interazioni che coinvolgono A e che per questo considereremo solamente ΔY_A e non, per esempio, ΔY_Q . Tuttavia, si potrebbero anche considerare diverse situazioni ed analizzarne ognuna in modo diverso, per comprendere ancora meglio le relazioni tra le variabili e dunque il modello finale.

Per stimare le interazioni tra due variabili (per esempio Q e A) si può usare la stratificazione, oppure la stima di una regressione su tutti i dati, incluso il prodotto dei termini $Q \times A$.

Quando gli IDAG sono usati insieme ai DAG, ci danno indicazioni su come eseguire le stime. Per quanto riguarda il confondimento, un criterio sufficiente affinché questo non sia presente è che entrambe le variabili che interagiscono non siano confuse. In generale, le variabili dovranno essere condizionate per assicurarsi che A e Q non abbiano alcuna relazione di confusione. Per capire quali variabili condizionare si possono osservare gli standard DAG. Tuttavia, questi ultimi non forniscono indicazioni sulla misura in cui è necessaria la stratificazione o l’inclusione dei termini del prodotto $A \times Q$.

Consideriamo invece una nuova situazione, in cui si è interessati ad esaminare non più l'interazione di per sé, ma l'effetto complessivo, come un effetto causale medio. Se le interazioni sono comunque presenti, la selezione del campione causerà spesso problemi di generalizzabilità, poiché l'effetto causale medio nel campione selezionato può differire da quello nella popolazione target. In generale questo problema sorgerà se la selezione dipende da variabili che influenzano l'effetto causale in esame.

I DAG standard sono informativi sui pregiudizi che potrebbero sorgere a causa del campionamento non causale, indipendentemente dalla misura dell'effetto scelta, tuttavia non sono informativi sul fatto che, per una misura di effetto scelta, ci siano effettivamente interazioni rispetto alle variabili da cui dipende la selezione e quindi se la generalizzabilità sia di fatto compromessa. Di conseguenza, la loro utilità è limitata. Gli IDAG possono invece essere usati per ottenere valide stime di interazione che non siano confuse ed anche valide stime esterne, nel senso di generalizzabilità delle stime degli effetti complessivi. Essi sono parametrici e dunque un po' meno generali dei DAG standard, tuttavia gli autori di "*A directed acyclic graph for interactions*"⁴ ritengono che questo sia un vantaggio che riduce il divario tra teoria e stima.

2.8 Come rappresentare i DAG?

I DAG sono strumenti didattici, ci fanno capire quali dati raccogliere e modellare ed anche come testare le relazioni e falsificare le affermazioni di causalità. Appaiono dunque come uno strumento utile da utilizzare di più nei propri studi.

Dagitty.net è un sito molto utile per il nostro scopo. "DAGitty è un ambiente basato su browser per la creazione, la modifica e l'analisi di diagrammi causali (noti anche come grafi aciclici diretti o reti bayesiane causali). L'attenzione si concentra sull'uso di diagrammi causali per ridurre al minimo i pregiudizi negli studi empirici di epidemiologia e altre discipline"⁵.

Si può lanciare direttamente online sul proprio browser, o si può scaricare sul proprio dispositivo per poterlo utilizzare offline, inoltre è anche disponibile un pacchetto R chiamato "*dagitty*".

In figura 2.19 vediamo ciò che si può osservare, se si decide di lanciarlo direttamente online.

⁴Nilsson A., Bonander C., Stromberg U., Bjork J., *A directed acyclic graph for interactions*, in "International Journal of Epidemiology", a. 00, n. 00, 2020.

⁵<http://dagitty.net/>

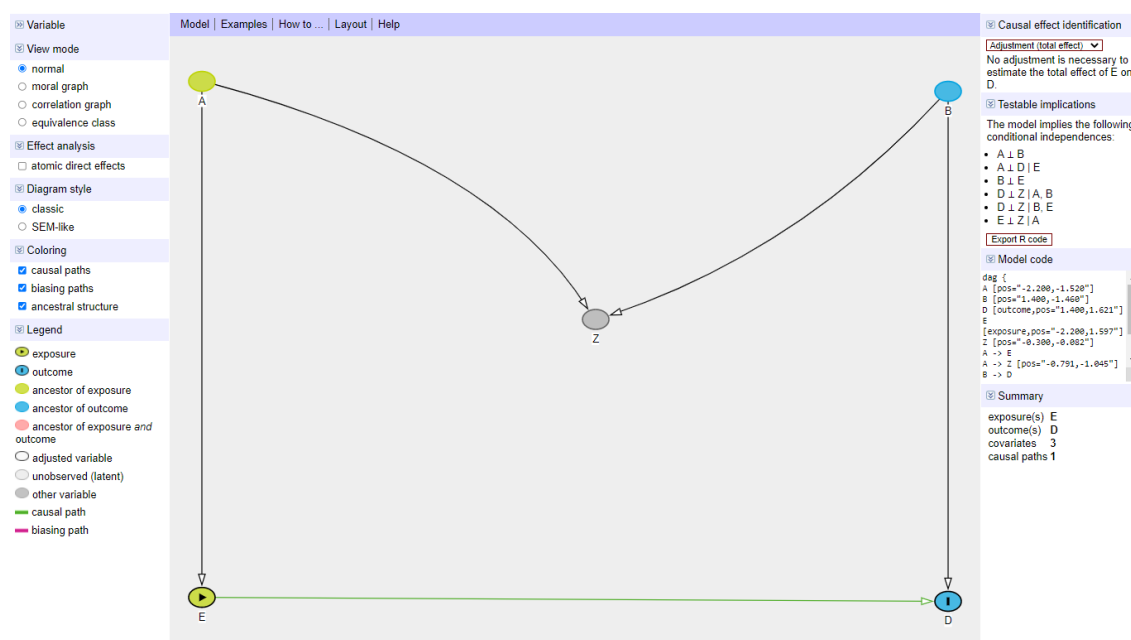


FIGURA 2.19: Schermata principale di DAGitty
Fonte: screenshot della pagina <http://dagitty.net/dags.html>

Ci sono moltissime possibilità di modificare il proprio DAG a proprio piacimento e soprattutto in base allo studio che si vuole raffigurare. Si può decidere se la variabile selezionata è un' esposizione oppure un *outcome* ed anche se è aggiustata oppure non osservata. Si può modificare la vista, che può essere normale, morale, di correlazione oppure con classi equivalenti. Si può scegliere se considerare degli effetti atomici diretti e lo stile del diagramma. Si possono anche cambiare i colori, vedere la legenda di quanto raffigurato, identificare la causa degli effetti identificati e quindi se l'aggiustamento è totale o diretto. Si osservano le condizioni che il modello implica, il codice del modello ed anche il riassunto che esplicita le variabili di esposizione, l'*outcome*, il numero di covariate e il numero di percorsi causali. Inoltre, il sito ha anche un *help* che aiuta e guida chi è alle prime armi, e un *learn* che riprende alcuni aspetti teorici utili per comprendere meglio la rappresentazione dei DAG.

Vediamo insieme un esempio di utilizzo di DAGitty.

È noto che il fumo aumenti l'incidenza di tumore ai polmoni e che l'età sia una variabile confondente. Decidiamo dunque di inserire la variabile fumo come esposizione, il tumore ai polmoni come *outcome* e la variabile età come variabile di aggiustamento. Sulla destra della nostra schermata (Figura 2.20) notiamo che effettivamente, per stimare l'effetto del fumo sul tumore al polmone, si deve aggiustare per l'età. Inoltre, nel *summary* si vede appunto come il fumo sia l'esposizione e il tumore ai polmoni

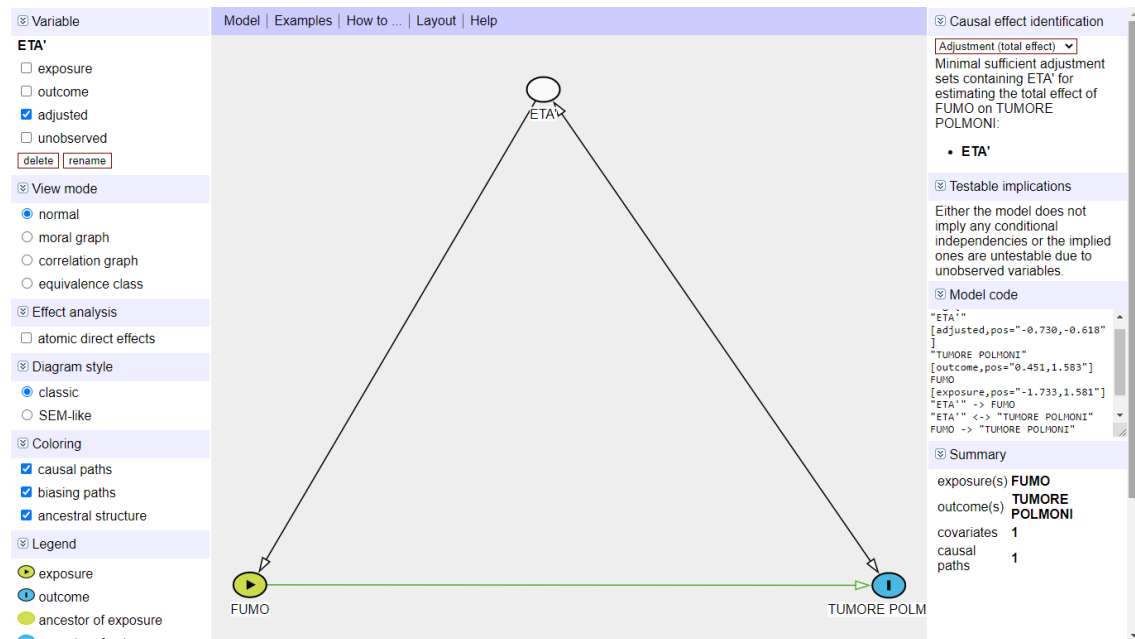


FIGURA 2.20: Esempio di utilizzo di DAGitty
Fonte: screenshot della pagina <http://dagitty.net/dags.html>

sia l'*outcome*, ed anche che ci siano una sola covariata e un solo percorso causale. Si tratta di un esempio molto semplice, ma permette subito di capire l'intuitività e la flessibilità dello strumento di cui si dispone.

Conclusione

I DAG aiutano a chiarire i problemi concettuali e a migliorare la comunicazione tra i ricercatori, dato che riassumono la conoscenza e le assunzioni in modo intuitivo. Come afferma Hernan nel suo libro "*What if*", "l'uso di grafici nei problemi legati all'inferenza causale rende più facile seguire un consiglio sensato: disegna le tue assunzioni prima di trarre le tue conclusioni".

Questi grafici permettono di codificare le assunzioni relative alla causalità e di decidere a priori quali variabili devono essere aggiustate nell'analisi e quali no. Tuttavia, molti continuano ad usare l'approccio tradizionale per rispondere ai problemi clinici, in quanto spesso richiedono DAG complicati e sono dunque scoraggiati dalla loro apparente complessità. Ci sono ulteriori criticità, tra cui la difficoltà nel capire la grandezza dell'effetto di una relazione causale tra due variabili o la relazione temporale tra i fattori considerati.

Ricorrere ai DAG è comunque utile per individuare le covariate da includere nello studio analizzato, in modo tale da minimizzare la grandezza del *bias* nelle stime prodotte.

Bibliografia e Sitografia

- [1] Cartabellotta N., *Pillole di metodologia della ricerca*, 2010,
https://www.evidence.it/articoli/pdf/2010_2_3.pdf
- [2] Dallolio L., Bellocchio R., Richiardi L., Fantini M. and the Causal Inference In Epidemiology (ICE) SISMEC Working Group, *Using directed acyclic graphs to understand confounding in observational studies*, in “biomedical statistics and clinical epidemiology”, a. III, n. 2, 2009.
- [3] Desmatron, *Teoria dei Grafi*, 2004
https://www.matematicamente.it/staticfiles/teoria/geometria/teoria_dei_grafi.pdf
- [4] Ferguson K., McCann M., Katikireddi S., Thomson H., Green M., Smith D, Lewsey J, *Evidence synthesis for constructing directed acyclic graphs (ESC-DAGs): a novel and systematic method for building directed acyclic graphs*, in “International Journal of Epidemiology” a. XLIX, n. 1, 2020.
- [5] Greenland S, Pearl J, Robins JM., *Causal diagrams for epidemiologic*, in “Epidemiology”, a. X, n. 1, Gennaio 1999.
- [6] Hernan M., Robins J., *Causal Inference: What If*, 2020, manoscritto inedito.
<https://www.hsph.harvard.edu/miguel-hernan/causal-inference-book/>
- [7] Iannaccone S., *Ecco il generatore di correlazioni assurde*, 2014
https://www.wired.it/scienza/lab/2014/05/12/generatore-correlazioni-assurde/?refresh_ce=
- [8] Murray E., *A cartoon Guide to Causal Inference*, 2019,
<https://www.slideshare.net/EleanorMurray8/>
- [9] Nilsson A., Bonander C., Stromberg U., Bjork J., *A directed acyclic graph for interactions*, in “International Journal of Epidemiology”, a. 00, n. 00, 2020.
- [10] Pearl J., Glymour M., Jewell N., *Causal inference in statistics a Primer*, New York, John Wiley, 2016.

- [11] Poletti G., *grafi e strutture. Appunti di Teoria dei Grafi*, Università di Ferrara.
http://www.unife.it/lettere/filosofia/comunicazione/insegnamenti/tecnologie_informatiche_multimediali/archivio/aa-2012-2013-1/materiale-didattico/dispense-e-link/note-di-teoria-dei-grafi
- [12] Roderick A. Rose, *Pearls of Wisdom in Causal Analysis: The Directed Acyclic Graph*, 2014
https://tatetalks.web.unc.edu/wp-content/uploads/sites/7733/2014/11/Rose_Pearl_DAG.pdf
- [13] Wedlin A., *Didattica per il corso di inferenza Bayesiana. Un'introduzione all'inferenza su processi stocastici*, Università di Trieste
https://moodle2.units.it/pluginfile.php/13219/mod_resource/content/1/DIDATTICA.pdf
- [14] *Correlazione o causalità?*, Statistics Knowledge Portal. Un'introduzione gratuita alla statistica online.
https://www.jmp.com/it_it/statistics-knowledge-portal/what-is-correlation/correlation-vs-causation.html#:~:text=La%20correlazione%20%C3%A8%20indice%20della,%C3%A8%20necessariamente%20indice%20di%20causalit%C3%A0%E2%80%9D [visitata 14.5.2021]
- [15] <http://dagitty.net/> [visitata il 16.5.2021]

Ringraziamenti

Il mio percorso universitario non è ancora finito, ma eccomi arrivata al raggiungimento del mio primo obiettivo: la laurea triennale.

Ci tengo a ringraziare il mio relatore, il prof. Rino Bellocco, che mi ha proposto di trattare i DAG nella mia tesi, permettendomi di conoscere e comprendere meglio l'argomento. Lo ringrazio inoltre per avermi fatto conoscere il linguaggio latex, con cui ho scritto questo elaborato.

Ringrazio i miei genitori, che mi hanno supportata in questa scelta universitaria, per me molto difficile, consentendomi di fermarmi un anno dopo il liceo per capire davvero cosa volessi studiare.

Ringrazio i miei fratelli, che mi hanno sopportata ad ogni "sclero" prima di un esame (e purtroppo per voi, dovrete continuare a farlo).

Ringrazio i miei nonni e i miei zii per esserci sempre stati, anche solo col pensiero.

Ringrazio Martina, che da una mia proposta di aiutarla a studiare Algebra in modo che io potessi ripassarla, è diventata la mia costante compagna di studi, di ansie, ma anche di molte gioie.

Ringrazio Leonardo, che più di tutti sa cosa significa avere a che fare con la me in sessione e riesce sempre a tranquillizzarmi.

Ringrazio Laura, che ogni volta che le dicevo che avevo un esame di lì a poco mi diceva "ma cosa ti lamenti a fare che poi prendi 30", poi magari il 30 non lo prendevo, ma l'esame lo passavo sempre. Grazie per il tuo prezioso sostegno.

Ringrazio Margherita, sempre presente nel momento del bisogno, per uno dei miei sfoghi, semplicemente per fare due chiacchiere e soprattutto, compagna di camminate necessarie per allievare un po' la tensione degli esami.

Ringrazio Diletta, che anche se non ci sentivamo spesso, quando succedeva, ci sfogavamo l'una con l'altra e ci spronavamo ad andare avanti, che mancava sempre meno all'obiettivo.

Sono stati tre anni importanti per me, tre anni di impegno, di studio, di sacrificio, ma anche di emozioni, di soddisfazioni, di nuovi incontri e nuove amicizie.

Grazie quindi a tutti voi che avete creduto e continuerete a credere in me.

