

# SPIEGARE O PREDIRE?

Le Due Culture del Modelling  
Statistico nella Ricerca Psicometrica



# GLI OBIETTIVI DELLA PSICOLOGIA SCIENTIFICA

- Essere in grado di spiegare le cause che sottendono un comportamento;
- Essere in grado di predire comportamenti che ancora non si sono verificati.

...Ma questi obiettivi che relazione hanno tra loro?

# SPIEGARE E PREDIRE NELLA RICERCA PSICOLOGICA

Spiegazione e predizione vengono generalmente considerati obiettivi sovrapponibili.  
In altre parole, la pratica comune di ricerca sembra sottendere il seguente ragionamento:



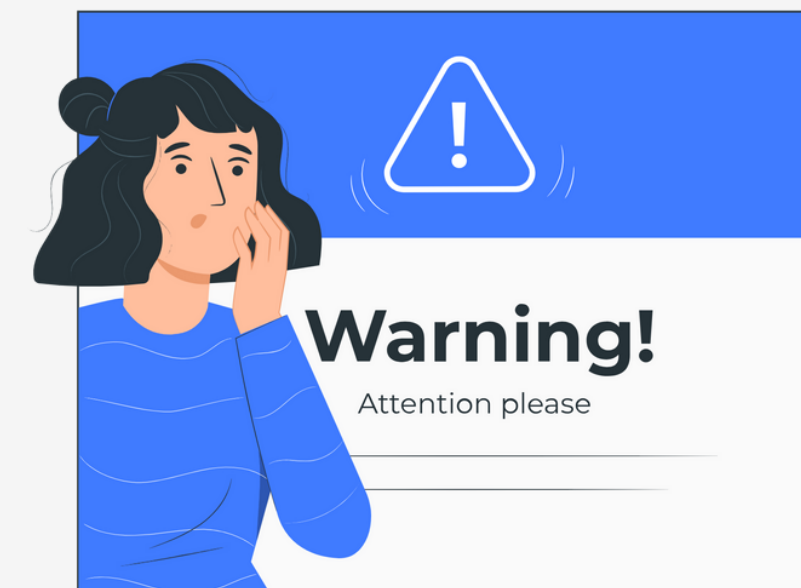
*"Se conosco tutte le cause di un comportamento e tutte le inter-relazioni tra queste cause, allora sarò in grado di predire con buona accuratezza quel comportamento."*

# SPIEGARE E PREDIRE NELLA RICERCA PSICOLOGICA

Se dal punto di vista filosofico spiegazione e predizione sembrano essere concetti compatibili...

Dal punto di vista statistico, non è corretto affermare che il modello che meglio approssima il processo di generazione dei dati è anche il modello che garantisce la migliore predizione di nuove osservazioni.

Inoltre, non è detto che sia possibile approssimare i fenomeni studiati abitualmente in ambito psicologico con modelli comprensibili alla ragione umana.



# LE DUE CULTURE DEL MODELLING STATISTICO

Breiman (2001) sostiene che esistono due culture nell'utilizzo del modeling statistico:

- una, assumendo un'ottica esplicativa, presume che i dati vengano generati da un modello stocastico e viene definita “Data Modeling Culture”;
- l'altra, assumendo un'ottica predittiva, utilizza modelli algoritmici, considera il meccanismo di generazione dei dati come sconosciuto e viene definita “Algorithmic Modeling Culture”.

Sulla stessa scia, Shmueli (2011) e Yarkoni & Westfall (2017) sostengono che spiegazione e predizione sono obiettivi distinti che danno vita a percorsi di ricerca differenti.

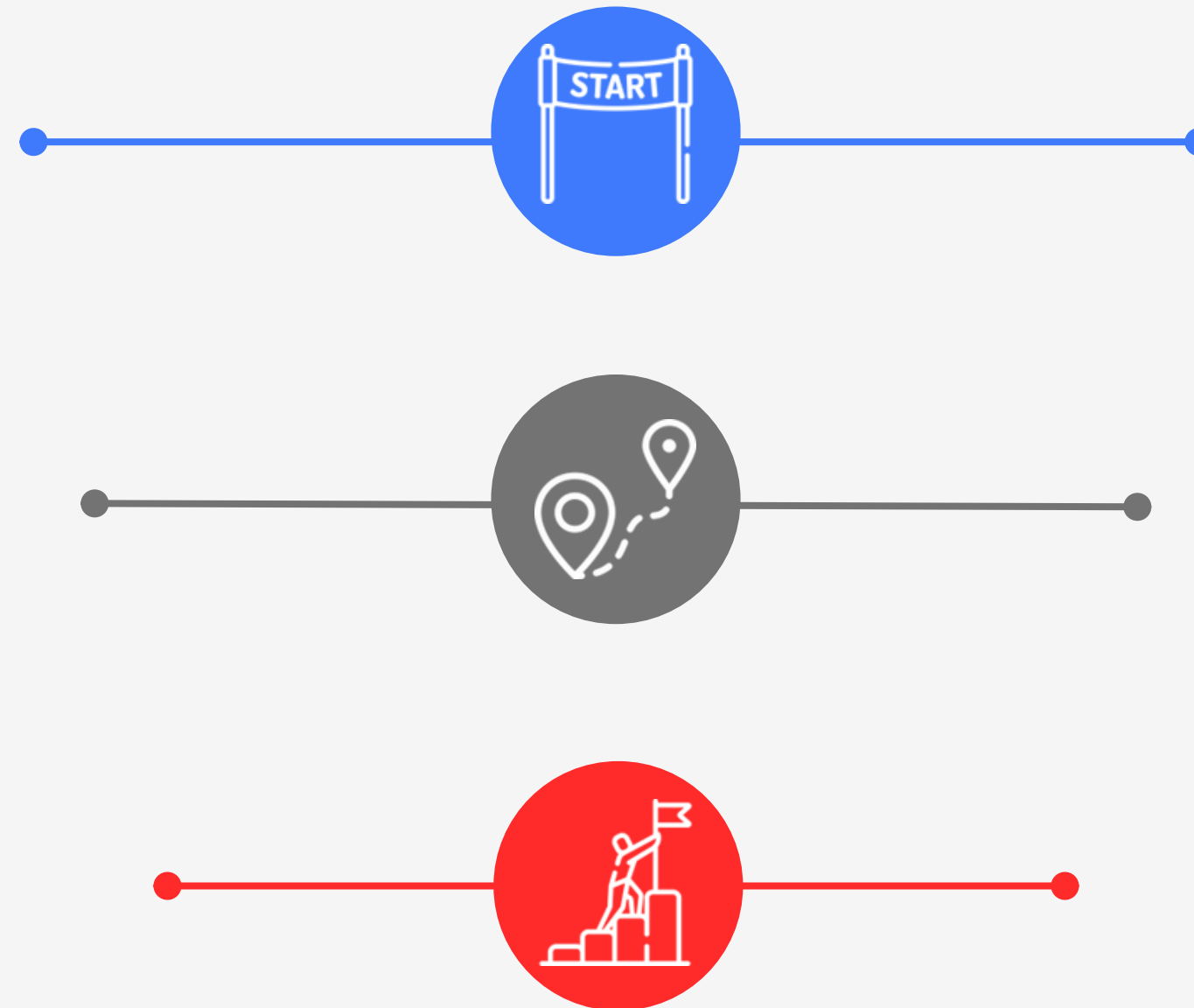
# SPIEGARE O PREDIRE

## Modelling Esplicativo

Il punto di partenza è il modello teorico.

I ricercatori sono interessati a testare le ipotesi di relazione tra le variabili e le loro intercorrelazioni.

L'obiettivo principale è stimare i parametri del modello usando l'inferenza statistica al fine di ottenere parametri che siano, in media, simili ai parametri della popolazione di riferimento.



## Modelling Predittivo

Il punto di partenza sono i dati.

I ricercatori sono interessati ad ottenere un algoritmo capace di generare predizioni accurate.

L'obiettivo principale è predire in maniera accurata valori di output per nuovi valori di input e migliorare la performance predittiva del modello.

# LA COSTRUZIONE DI UN MODELLO ESPLICATIVO

Le ipotesi di ricerca sono definite sulla base di costrutti teorici. Solitamente, è presente uno schema che illustra le relazioni causali tra i costrutti ipotizzate dai ricercatori.

Viene creato un collegamento tra il costrutto e le misurazioni direttamente osservabili utilizzando giustificazioni teoriche e basandosi sulla letteratura di riferimento: il costrutto viene, cioè, operazionalizzato.

Dal momento, quindi, che nel modeling esplicativo le variabili sono viste come operazionalizzazioni dei costrutti, la scelta di quest'ultime è basata sul ruolo che ha il costrutto nella struttura teorica.

Dopo questi passaggi si procede a introdurre il modello statistico utilizzato e le sue assunzioni: il modeling esplicativo richiede modelli statistici interpretabili legati al modello teorico sottostante.

La valutazione del modello prevede diversi criteri per giudicarne la bontà di adattamento ai dati e la sua capacità esplicativa; tra questi, i più importanti sono la significatività dei parametri stimati, la coerenza interpretativa e la porzione di variabilità delle osservazioni spiegata dal modello.





# LA COSTRUZIONE DI UN MODELLO PREDITTIVO

Si definisce “modeling predittivo” il processo che, attraverso l'applicazione di modelli statistici o di algoritmi di “data mining” ai dati, consente di predire nuove o future osservazioni.

In questo tipo di approccio, non è necessario conoscere l'esatto ruolo di ogni variabile in termini di struttura causale sottostante: i criteri per la scelta delle variabili risiedono nella qualità dell'associazione tra predittore e risposta, nella qualità dei dati e nella disponibilità dei predittori al momento della predizione.

In questo caso, dunque, sono i dati a guidare il processo e non la formulazione teorica sottostante.

Nell'ambito del modeling predittivo, la validazione si concentra sulla generalizzazione, cioè sull'abilità del modello di predire accuratamente nuove osservazioni. Nel modeling predittivo, esistono casi in cui eliminare variabili con coefficienti bassi, anche se significative, aumenta l'accuratezza della predizione.





# E IN PSICOLOGIA?

Un gran numero di articoli di psicologia presentano in modo prominente la parola "predizione" nei loro titoli e nell'interpretazione dei loro risultati. Affermazioni come:

- *"L'impulsività predice il gioco d'azzardo problematico nei giovani uomini di basso status socioeconomico" (Vitaro, Arseneault, & Tremblay, 1999)*
- *"L'attività cerebrale predice quanto bene le esperienze visive verranno ricordate" (Brewer, Zhao, Desmond, Glover, & Gabrieli, 1998)*
- *"La selettività dei gesti precoci predice l'apprendimento del linguaggio in seguito" (Rowe & Goldin-Meadow, 2009)*

riflettono l'idea intuitiva che una vasta gamma di modelli statistici siano, in un certo senso, modelli predittivi.

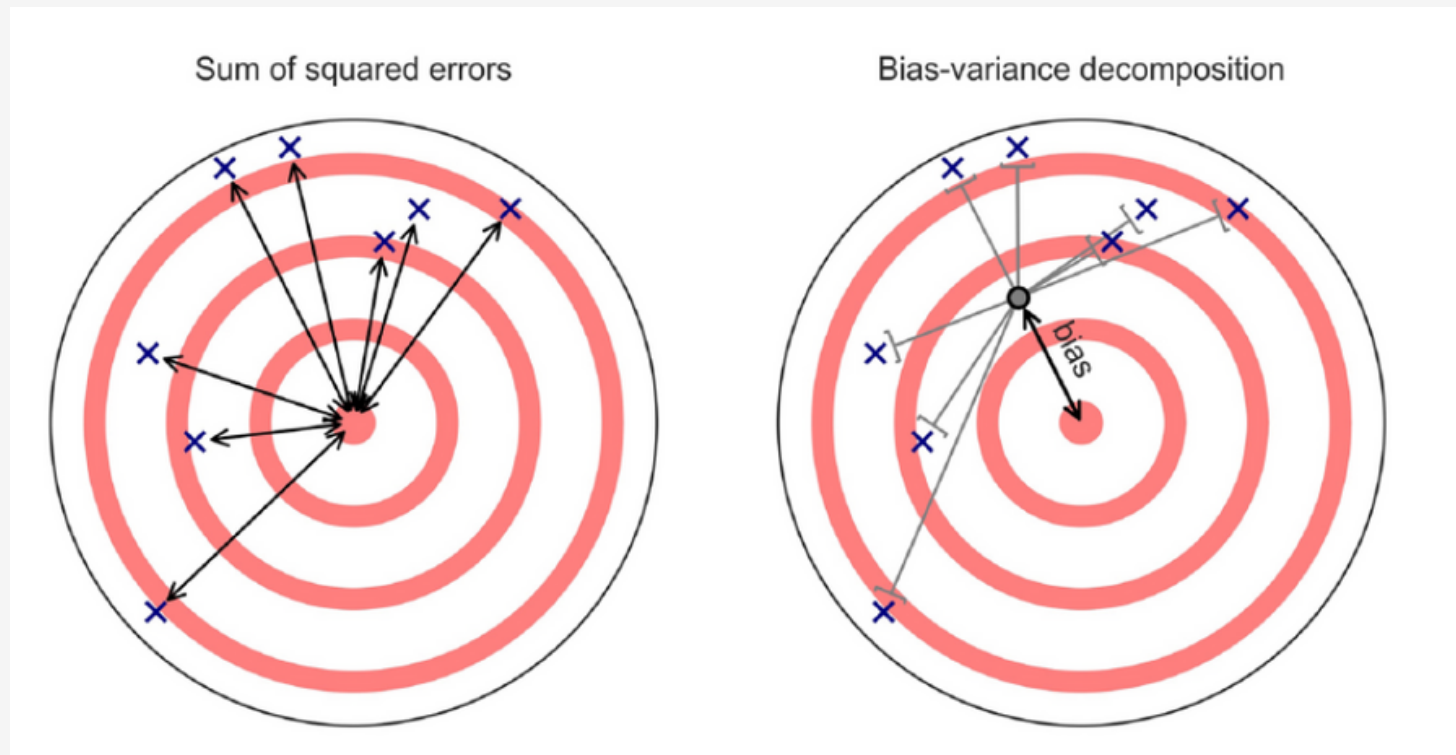
# E IN PSICOLOGIA?

La ricerca psicologica, così come la ricerca nell'ambito delle scienze sociali, si è concentrata prevalentemente sullo sviluppo di modelli esplicativi, interpretandoli spesso come modelli predittivi.

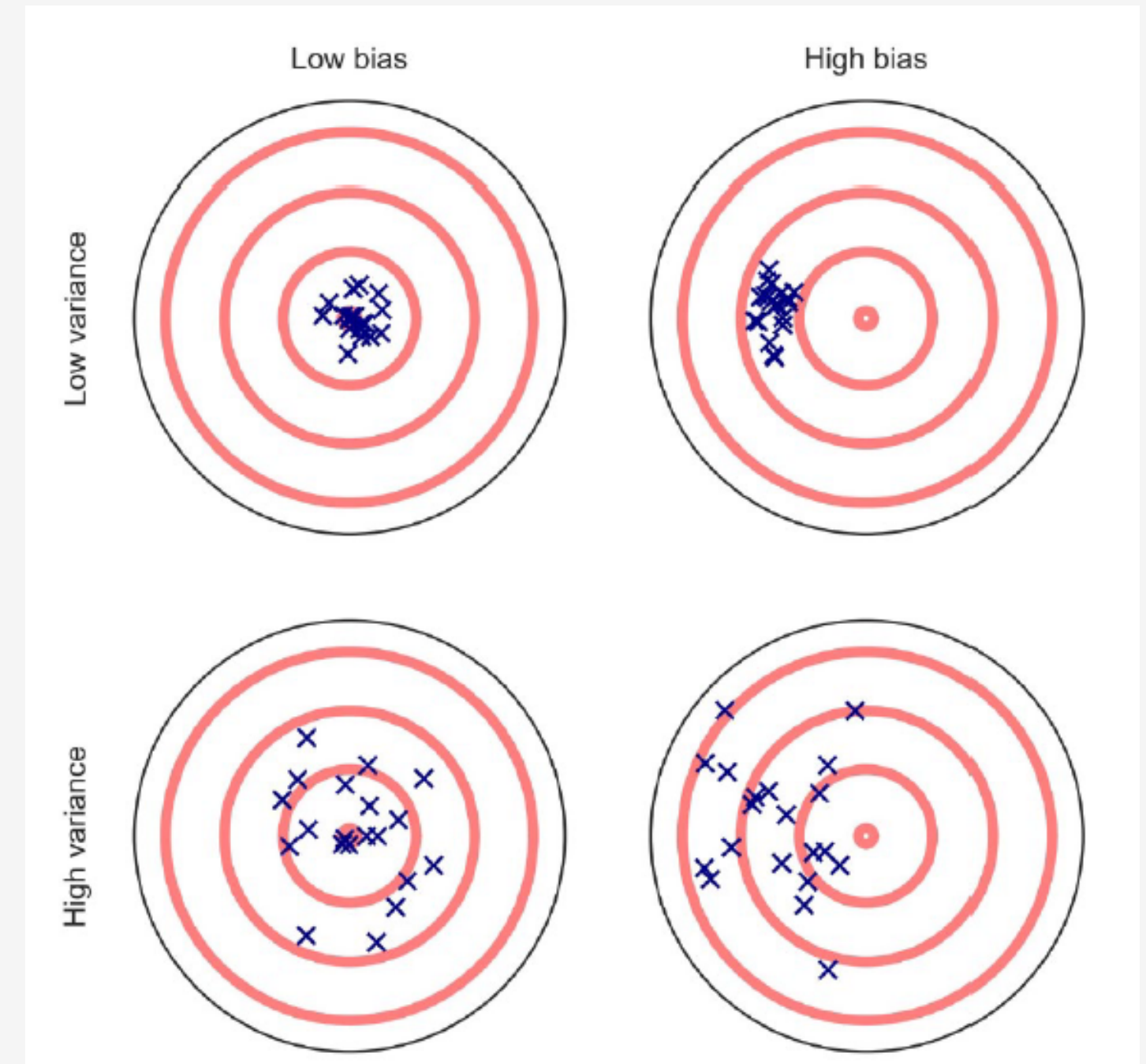
Questo può dar vita a risultati fuorvianti.



# DECOMPORRE L'ERRORE IN BIAS E VARIANZA



La somma totale degli errori al quadrato può essere considerata come composta di due termini: il *bias* che cattura la tendenza sistematica di un modello a deviare dal punteggio vero in modo prevedibile e la *varianza* che rappresenta le deviazioni delle osservazioni dalla previsione attesa del modello.



# IL TRADE-OFF BIAS/VARIANZA

Con il termine “bias” si definisce un particolare tipo di errore: la tendenza per un modello a produrre costantemente risposte sbagliate in una particolare direzione (ad esempio, stime che sono costantemente troppo alte).

Il bias può essere arginato con la varianza, che si riferisce alla misura in cui i parametri di un modello tenderanno a deviare dalla loro media nei diversi set di dati. A parità di altre condizioni, infatti, quando aumenta la varianza di uno stimatore, diminuisce il suo bias.

L'approccio predittivo cerca di minimizzare la combinazione di bias e varianza, sacrificando occasionalmente l'accuratezza teorica per una migliore precisione empirica.

Nell'approccio esplicativo, i parametri vengono stimati in modo tale che le stime siano, in media, uguali ai valori dei parametri della popolazione di riferimento.

Per i ricercatori appartenenti alla cultura psicometrica classica, le stime o le previsioni che non soddisfano questa proprietà sono da evitare e sono definite distorte.





# IL TRADE-OFF BIAS/VARIANZA

L'approccio esplicativo dà la priorità alla minimizzazione del bias. Tuttavia, poiché l'errore di previsione totale è uguale alla somma di bias e varianza, questo approccio corre il rischio di produrre modelli che sono essenzialmente inutili per la previsione, poiché la varianza è troppo grande.

È importante sottolineare che anche l'utilità delle teorie in esame diminuisce notevolmente, perché operare in un regime ad alta varianza implica che i modelli che si ricavano dai dati del proprio campione sono altamente instabili e possono cambiare drasticamente a causa di cambiamenti relativamente piccoli nei dati.

PERCHÈ?



# L'OVERFITTING

Si consideri un modello di regressione con due predittori, specificato nel modo seguente:

$$y = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i \quad (1)$$

Dopo aver stimato i parametri, otteniamo questo risultato:

$$y = 1.60 + 0.65x_{1i} + 0.60x_{2i} + \epsilon_i \quad (2)$$

Un ricercatore potrebbe affermare che, ad esempio, utilizzando un set di predittori demografici e di personalità è in grado di “predire” il 45% della variabilità del livello d’istruzione.

Implicitamente, questo significa essere in grado di predire, ragionevolmente in maniera accurata, il livello di istruzione di una persona non appartenente al campione estratta casualmente dalla popolazione di riferimento.

# L'OVERFITTING

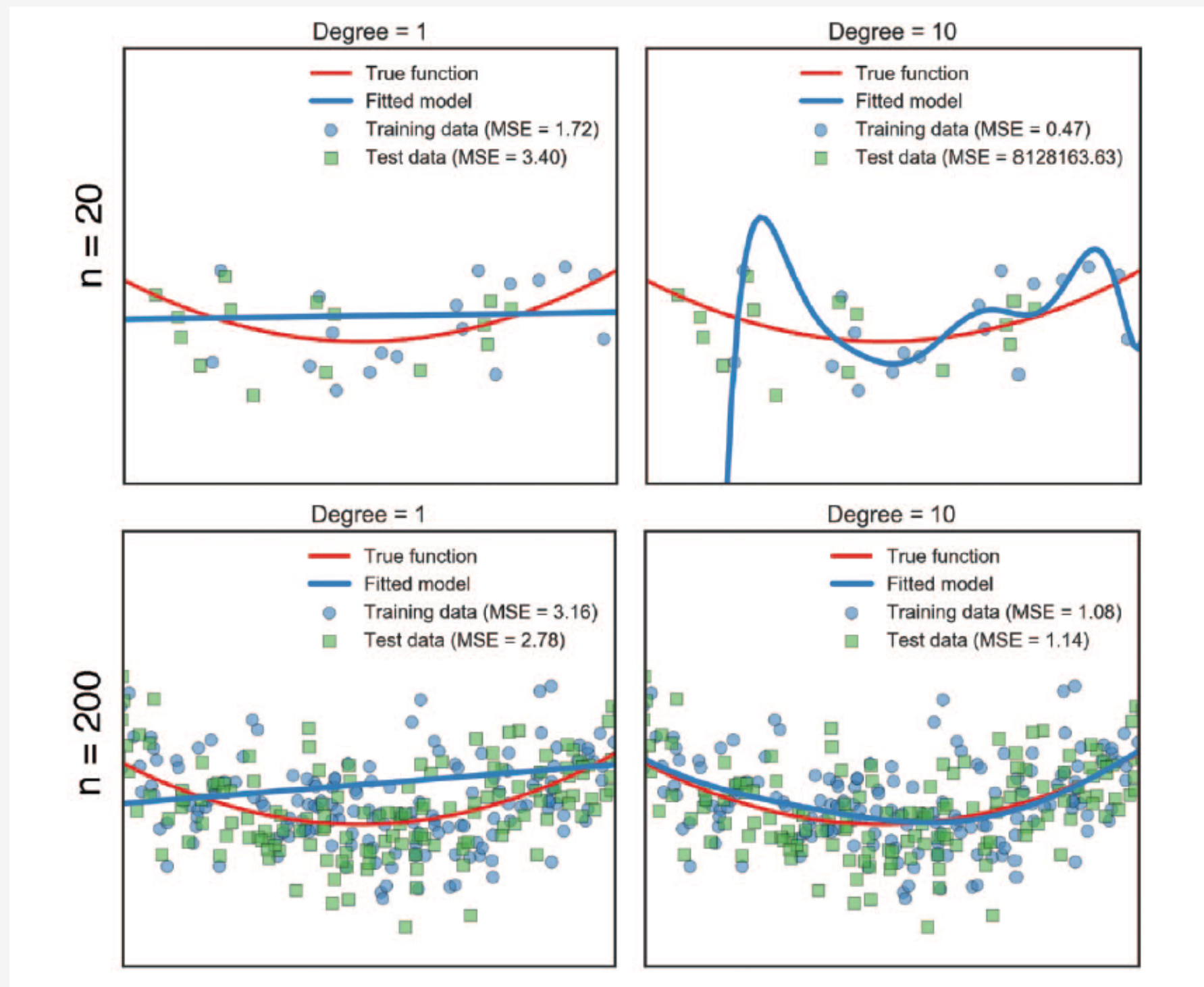
È opportuno notare che le stime dei parametri ottenute nel campione in cui viene costruito il modello non avranno prestazioni simili se applicate ad altri campioni estratti dalla stessa popolazione, dal momento che i parametri vengono stimati in modo da minimizzare la somma degli errori al quadrato in quel particolare campione (metodo dei minimi quadrati).

La statistica  $R^2$ , comunemente utilizzata per valutare la bontà di adattamento del modello ai dati, fornisce indicazioni sulla riduzione media della somma degli errori al quadrato in campioni casuali ripetuti simili a quello utilizzato per costruire il modello stimando ogni volta nuovi valori dei parametri. In altre parole, ci dà informazioni circa la bontà di adattamento di un modello lineare ai dati (equazione 1).

Tuttavia, le prestazioni dell'equazione 1 sono virtualmente sempre una stima eccessivamente ottimistica delle prestazioni dell'equazione 2, dal momento che i parametri stimati in un dato campione sono specificamente selezionati in modo da minimizzare la somma degli errori al quadrato in quel particolare campione.



# LA PREDIZIONE OUT-OF-SAMPLE



È possibile utilizzare la predizione out-of-sample per identificare errori di specificazione del modello senza correre il rischio di adattare eccessivamente il modello ai dati specifici.

Questo problema di sovra-adattamento del modello ai dati si definisce overfitting e si riferisce, in particolare, alla tendenza dei modelli statistici a considerare erroneamente il rumore campione-specifico (cioè le fluttuazioni casuali dei dati) come se fosse informazione.

Anche le analisi condotte con metodi predittivi incontrano problemi di generalizzazione.

In alcuni casi, si sviluppano modelli predittivi complessi che su piccoli set di dati raggiungono un'accuratezza di classificazione quasi perfetta, tuttavia, questo risultato non si replica su nuovi dati non utilizzati per sviluppare il modello, indicando così una scarsa capacità di generalizzazione. Anche in questo caso, si parla di overfitting.

# MODELLI PREDITTIVI E MODELLI ESPLICATIVI

Non esiste una linea netta di demarcazione tra metodi esplicativi e metodi predittivi, in quanto ad essere diverso è l'approccio con cui il modello viene costruito e valutato. L'analisi di regressione può essere un metodo esplicativo o predittivo, a seconda di come viene costruito e validato il modello.



# L'UTILITÀ DEI METODI PREDITTIVI NELLA RICERCA PSICOLOGICA

Per le ragioni discusse fino ad ora, non sembra possibile costruire dei modelli che contemporaneamente massimizzino la spiegazione e la predizione di un fenomeno. Tuttavia, riportare le performance predittive di un modello esplicativo può fornire diversi vantaggi, tra cui:

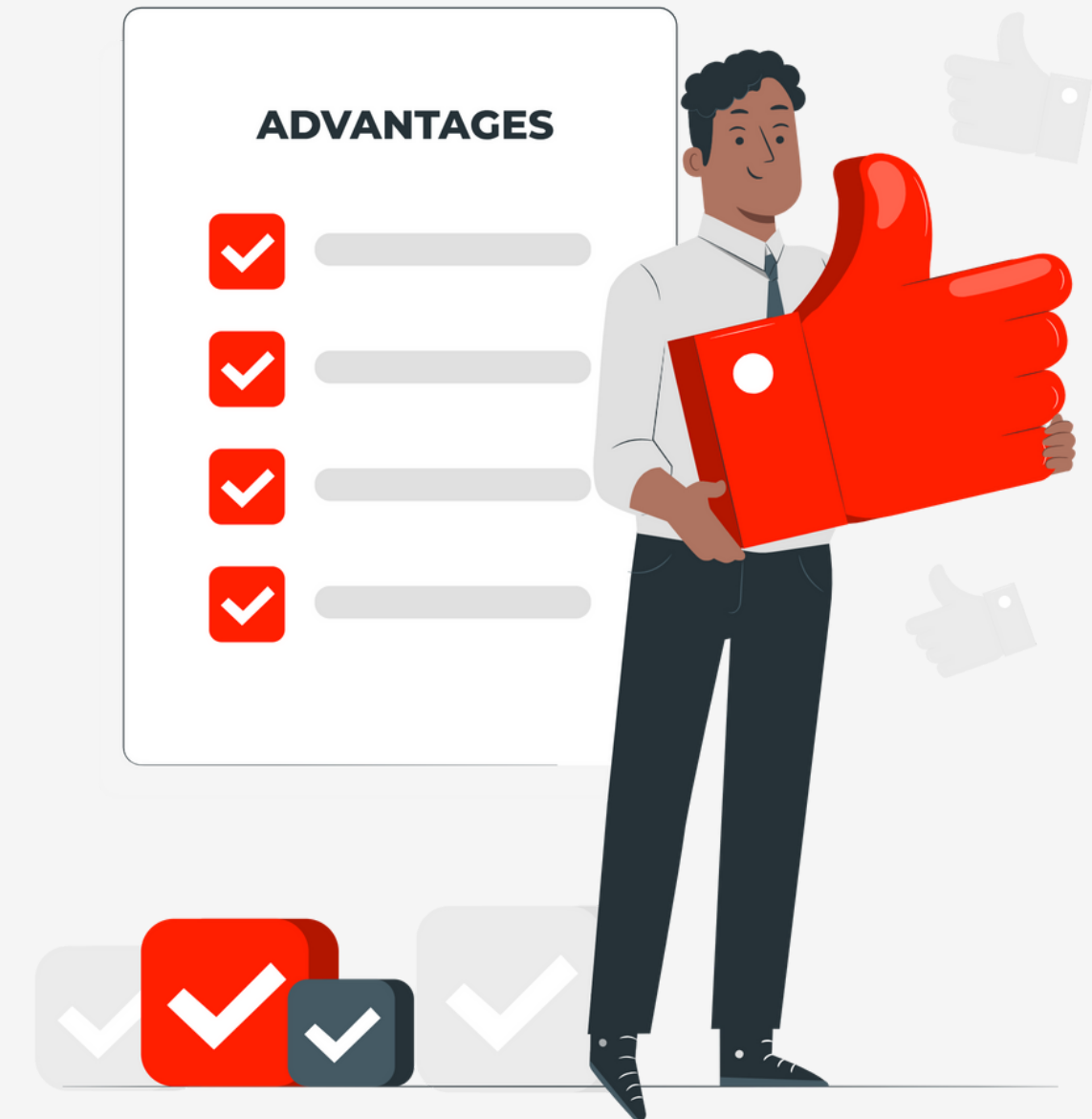
- la possibilità di scegliere tra diversi modelli esplicativi quello più predittivo
- verificare l'effettiva utilità del modello esplicativo proposto

D'altro canto, mantenere anche una valutazione esplicativa di un modello predittivo è utile in tutti quegli ambiti scientifici, inclusa la psicologia, dove l'interpretazione di un fenomeno dal punto di vista teorico rimane una questione importante.

# SPEIGARE E PREDIRE: UN PASSO VERSO L'INTEGRAZIONE

Ad ogni modo, integrare metodi tipici dell'approccio predittivo nell'analisi dei dati psicologici può aiutare a migliorare aspetti in cui il modeling esplicativo manifesta dei limiti. In particolare:

- può aiutare a superare delle assunzioni di base tipiche dell'approccio esplicativo (ad esempio, la linearità delle relazioni).
- offre la possibilità di utilizzare strumenti di analisi innovativi molto comuni in altri campi scientifici dove l'approccio predittivo è quello predominante
- può mitigare problemi relativi alla replicabilità della ricerca e alla generalizzazione dei risultati.



# CONCLUSIONI

- Spiegare e predire sono due obiettivi diversi, che danno vita a percorsi di ricerca differenti.
- La ricerca psicologica si è concentrata per lungo tempo sui modelli esplicativi inferendo da questi il potere predittivo dei propri modelli.
- Non è possibile costruire un modello che contemporaneamente spieghi e predica in maniera perfetta un fenomeno, ma è possibile ottenere modelli che contemporaneamente spiegano un fenomeno e lo predicono almeno in maniera adeguata.
- Integrare modelli predittivi in ambito psicologico può aiutare la ricerca a superare assunzioni di base tipiche dell'approccio esplicativo e ad avere una valutazione più completa dei modelli che descrivono un fenomeno.



# BIBLIOGRAFIA

Breiman, L. (2001). Statistical Modeling: The Two Cultures (with comments and a rejoinder by the author). *Statistical Science*, 16.

<https://doi.org/10.1214/ss/1009213726>

Hofman, J. M., Watts, D. J., Athey, S., Garip, F., Griffiths, T. L., Kleinberg, J., ... & Yarkoni, T. (2021). Integrating explanation and prediction in computational social science. *Nature*, 595(7866), 181-188.

Shmueli, G. (2011). To Explain or to Predict? *Statistical Science*, 25.

<https://doi.org/10.1214/10-STS330>

Yarkoni, T., & Westfall, J. (2017). Choosing Prediction Over Explanation in Psychology: Lessons From Machine Learning. *Perspectives on Psychological Science: A Journal of the Association for Psychological Science*, 12(6), 1100–1122. <https://doi.org/10.1177/1745691617693393>

