













Come calcolare i p-value dei coefficienti di regressione di LINESTX

PUBBLICATO LUGLIO 11, 2023 DI FRANCESCO BERGAMASCHI E DAVID BIANCONI

La funzione DAX <u>LINESTX</u> è stata introdotta recentemente in Power BI Desktop allo scopo di permettere la stima dei coefficienti di una regressione lineare. Da allora sono emersi diversi video e articoli su come usarla. Questo articolo non verterà su questo aspetto, di conseguenza. Per l'uso di <u>LINESTX</u> rimandiamo all'ottimo articolo (corredato di video) di sqlbi (<u>link</u>).

In questo articolo ci preoccuperemo, invece, come preannunciato in <u>questo articolo</u>, di calcolare i *p-value* dei coefficienti stimati. È, infatti, molto importante ricordare che le regressioni lineari non possono essere usate per fare stime se non si fanno delle verifiche di ipotesi. Non ci addentreremo nel complesso sistema teorico che sta dietro questa affermazione ma non possiamo non dire che ogni stima potrebbe essere un abbaglio.

La statistica è così: nulla è certo. Ogni affermazione ha un certo livello di confidenza. Senza la stima dei *p-value* (e altre verifiche sui residui che non tratteremo qui ma che sono più semplici in software dedicati come R e che possono essere fatti in Power BI usando la R-*visual*), usare una regressione lineare per prendere decisioni è molto pericoloso.

Useremo un modello molto semplice, come sempre. In figura 1 sono visibili tre tabelle, tra loro scollegate. I dati sono nella tabella importata *SampleData*. Le altre due tabelle (*Criteria* e *Linear Regression Output Enriched*) sono state introdotte per, rispettivamente, creare le soglie di non rifiuto dell'ipotesi nulla tra cui scegliere e mostrare l'output di *LINESTX*, che è una tabella, arricchita dai calcoli dei *p-value*. Il report che verrà creato è mostrato in figura 2, dove tutto è risolto attraverso misure e l'output arricchito dai *p-value* non è osservabile ma dove è visibile l'uso della tabella *Criteria* (riquadro rosso).

























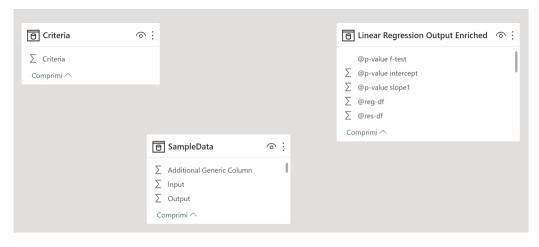


Figura 1

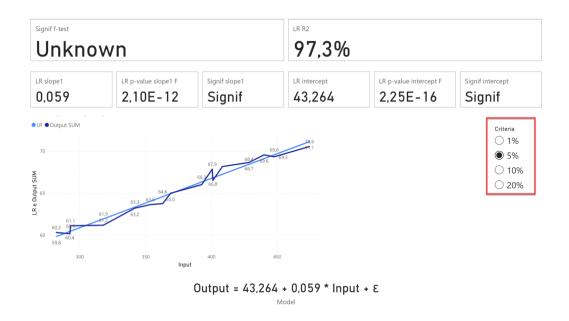


Figura 2

In figura 3 è mostrata la tabella SampleData. La tabella contiene un insieme di dati macroeconomici relativi al numero di persone occupate (Output) e al prodotto nazionale lordo o PNL – deflazionato ai valori del 1954 – (Input) rilevati negli Stati Uniti negli anni 1947-1962 (Additional Generic Column). Questo data set fa parte di uno studio più ampio realizzato da J. W. Longley (1967). Siamo interessati a spiegare l'andamento del numero di occupati (Output) in funzione del PNL (Input). Più avanti nell'articolo verrà mostrato che la tabella può avere un altro contenuto grazie ad un parametro di Power Query.



www.kubisco.com



info@kubisco.com



Pag. 2 di 13













gni (
=	Ò
=.	Č
	7
=	1
Ζij	`.
ıtilizz	-
zat	-
₩.	
0	6
a d	٠,
Φ	
dei conteni	
Φ	71101101
Ξ.	2
Ω	
0	7
⊇.	
Ж.	-
~	=
7	
≒	
iti è teni	Œ
Φ,	=
-	7
Φ	2
\supset	
_	C
₹	
U	2
a	
_	τ
to a verit	יוויקטנט שטטני
×	Ċ
rificar	5
<u> </u>	
ρý	(
=	5
a verificare	
0)	-
au	- 5
⇌	Ş
tono	וומסכומ
⋾	- 5
0	2
⊐	2
ma	-
7	ĉ
ımen	Ć
ā	~
namente l'as	זה ממשנו ממנטו כטטו כטווכי.
₹.	7
jej.	2
_	'n
מ	-
S	-
Ω	(
Φ	- (
⊇	2
Za	-
മ	۶
assenza di	2
=	-
Φ	7
rrori	,
o)	
~	_
	7
Φ	•
assenza di errori e la coe	0 00
മ	
C	ò
0	-
Ō	ā
	4
<u>ē</u>	2
7	7
Za	. `
ш	90100
nza ris	ć
g	_
þe	5
Ō.	٠
ťξ	٥
0	-
0)	7
≝.	(
∇	-
ai propri	טלטטנו ליכו טכטליו
0	(
∇	_
⊐.	7
0	-
ca	2
S	Č
	ĉ
ppri casi di a	טטו מו ממנוכו.
	5
appli	2
0	
℧	
=	
S	
'n.	
applicazior	
0	
ione.	
Φ	

Input 🔻	Output 🔻	Additional Generic Column
282,276	60,323	1947
293,137	61,122	1948
292,578	60,171	1949
317,988	61,187	1950
341,970	63,221	1951
353,720	63,639	1952
369,076	64,989	1953
363,112	63,761	1954
392,756	66,019	1955
400,746	67,857	1956
408,458	68,169	1957
401,215	66,513	1958
428,689	68,655	1959
440,106	69,564	1960
447,859	69,331	1961
474,674	70,551	1962

Figura 3

Sviluppo

L'output della funzione *LINESTX* è un tabella, che riporta una singola riga con diverse colonne. Ognuno dei valori rappresenta un termine che descrive il modello statistico stimato tramite la regressione. Ecco i diversi output forniti nella stima di un modello come il sequente:

y = intercept + slope1 * x1 + slope2 *x2 + ... + slopeN * xN + ε

- SlopeN (un output per ogni espressione x)
- Intercept
- StandardErrorSlopeN (un output per ogni espressione x)
- StandardErrorIntercept
- CoefficientOfDetermination (R2)
- StandardError
- **FStatistic**
- DegreesOfFreedom
- RegressionSumOfSquares
- ResidualSumOfSquares



www.kubisco.com



info@kubisco.com



Pag. 3 di 13













Nella lista non sono presenti i *p-value* né del test F (intero modello) né delle stime dei singoli coefficienti stimati (*SlopeN*). Tuttavia, è possibile calcolare i p-value dei singoli coefficienti stimati come segue.

Per prima cosa, osserviamo l'output generato da *LINESTX*, nella tabella calcolata *Linear Regression Output Enriched*. Essa è generabile con il seguente codice:

Linear Regression Output = LINESTX (SampleData, [Output], [Input])

La tabella creata è visibile in figura 4



Tra i vari risultati, *Intercept* fornisce la stima dell'intercetta della retta di regressione. Avendo un solo regressore (*Input*) per cercare di spiegare *Output*, la tabella contiene soltanto una colonna di stima di coefficienti di regressione (*Slope1*). Tale valore è una stima di pendenza (quantifica, per ogni incremento unitario dell'input, di quanto aumenta l'output).

Ecco come arricchire la tabella con i p-value dell'intercetta e della pendenza:

```
Linear Regression Output Enriched =

VAR SampleTable = SampleData

VAR SampleSize =

COUNTROWS ( SampleData )

VAR NumOfRegressors = 1

RETURN

ADDCOLUMNS (

ADDCOLUMNS (

ADDCOLUMNS (

LINESTX ( SampleTable, [Output], [Input] ),

"@t-test slope1", DIVIDE ( [Slope1], [StandardErrorSlope1] ),

"@t-test intercept", DIVIDE ( [Intercept], [StandardErrorIntercept] ),

"@reg-df", NumOfRegressors
),
```



www.kubisco.com



info@kubisco.com



Pag. 4 di 13















Ecco le colonne addizionali che sono state calcolate (figura 5), con in evidenza i *p-value* di interesse:



Figura 5

Le colonne Linear Regression Output Enriched[@t-test slope1] e Linear Regression Output Enriched[@t-test intercept] rappresentano un passo intermedio per arrivare al calcolo dei p-value dell'intercetta e del regressore; le colonne Linear Regression Output Enriched[@res-df] servirebbero al calcolo del p-value del test F che, tuttavia, non è possibile al momento e da qui la scelta di attribuire il valore "Unknown" a questa colonna. Il p-value del test F è, in effetti, il più importante perché chiarisce la significatività dell'intero modello; la ragione dell'attuale impossibilità del calcolo è la mancanza in DAX della funzione F.DIST.RT che è invece disponibile in Excel. Non ci addentriamo ulteriormente in questi complessi concetti in questo articolo.

Adesso che possiamo disporre dei *p-value* dei coefficienti del modello di regressione, passiamo alla sintesi delle misure e della tabella *Criteria*, che permettono di creare il report mostrato in figura 2.

La tabella *Criteria* è stata creata usando il tasto "Immettere i dati" in Power BI Desktop ed ha i tipici valori soglia che si usano in Statistica: 1%, 5%, 10%, 20%.

Ecco le misure:

Signif f-test = "Unknown"



www.kubisco.com



info@kubisco.com



Pag. 5 di 13















```
LR R2 =
VAR Line =
  LINESTX (ALLSELECTED (SampleData), [Output], [Input])
VAR Result =
  SELECTCOLUMNS (Line, [CoefficientOfDetermination])
RETURN
  Result
LR slope1 =
VAR Line =
  LINESTX ( ALLSELECTED ( SampleData ), [Output], [Input] )
VAR slope =
  SELECTCOLUMNS (Line, [Slope1])
RETURN
  slope
LR p-value slope1 =
VAR SampleSize =
  COUNTROWS (SampleData)
VAR Line =
  ADDCOLUMNS (
    ADDCOLUMNS (
      LINESTX (ALLSELECTED (SampleData), [Output], [Input]),
      "at-test slope1", DIVIDE ([Slope1], [StandardErrorSlope1])
    ),
    "@p-value slope1", T.DIST.2T ( ABS ( [@t-test slope1] ),                    SampleSize – 2 )
VAR Result =
  SELECTCOLUMNS (Line, [@p-value slope1])
RETURN
  Result
LR p-value slope1 F =
FORMAT ([LR p-value slope1], "Scientific")
```



www.kubisco.com



info@kubisco.com



Pag. 6 di 13







```
LR p-value slope1 F =
FORMAT ([LR p-value slope1], "Scientific")
Signif slope1 =
IF (
  [LR p-value slope1] < SELECTEDVALUE ( Criteria[Criteria] ),
  "Signif",
  "NOT Signif"
)
LR intercept =
VAR Line =
  LINESTX (ALLSELECTED (SampleData), [Output], [Input])
VAR Result =
  SELECTCOLUMNS (Line, [Intercept])
RETURN
  Result
LR p-value intercept =
VAR SampleSize =
  COUNTROWS (SampleData)
VAR Line =
  ADDCOLUMNS (
    ADDCOLUMNS (
      LINESTX (ALLSELECTED (SampleData), [Output], [Input]),
      "@t-test intercept", DIVIDE ([Intercept], [StandardErrorIntercept])
    ),
    '@p-value intercept", T.DIST.2T ( ABS ( [@t-test intercept] ), SampleSize – 2 )
VAR Result =
  SELECTCOLUMNS (Line, [@p-value intercept])
RETURN
```



Result

www.kubisco.com

#kubisco



info@kubisco.com



Pag. 7 di 13



www.kubisco.com











Pag. 8 di 13



```
LR p-value intercept F =
FORMAT ([LR p-value intercept], "Scientific")
Signif intercept =
IF (
  [LR p-value intercept] < SELECTEDVALUE ( Criteria[Criteria] ),
  "Signif",
  "NOT Signif"
Output SUM =
SUM (SampleData[Output])
LR =
VAR Line =
  LINESTX ( ALLSELECTED ( SampleData ), [Output], [Input] )
VAR slope =
  SELECTCOLUMNS (Line, [Slope1])
VAR intercept =
  SELECTCOLUMNS (Line, [Intercept])
VAR x =
  SELECTEDVALUE (SampleData[Input])
VAR Result = x * slope + intercept
RETURN
  Result
Model =
IF (
  [Signif slope1] = "Signif"
    && [Signif intercept] = "Signif",
  "Output = " & FORMAT ([LR intercept], "#,0.00#") & " + "
    & FORMAT ( [LR slope1], "#,0.00#" ) & " * Input + \epsilon",
```

info@kubisco.com













```
IF (
    [Signif slope1] = "Signif",
    "Output = " & FORMAT ( [LR slope1], "#,0.00#" ) & " * Input + ε",
    IF (
        [Signif intercept] = "Signif",
        "Output = " & FORMAT ( [LR intercept], "#,0.00#" ) & " + ε",
        "No statistical significance"
    )
)
```

In figura 2, tuttavia, i p-value sono talmente piccoli da portare al non rifiuto dell'ipotesi nulla in qualunque condizione di selezione delle soglie (valore selezionato della colonna *Criteria[Criteria]*). Per evidenziare la cosa, mostriamo qui di seguito la figura 2 con una modifica del valore soglia in un valore più restrittivo (dal 5% al 1%), come si vede sia l'intercetta che la pendenza stimate rimangono significative. Valori meno restrittivi, ovviamente, porterebbero a nessuna modifica ancora a maggior ragione.

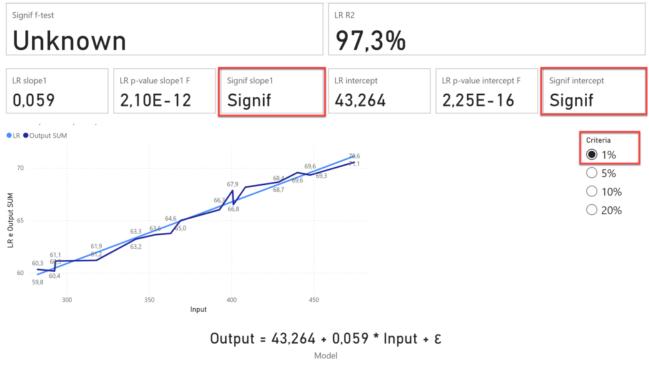


Figura 6



www.kubisco.com



info@kubisco.com



Pag. 9 di 13















Per permettere di apprezzare la dinamicità del report, come preannunciato ad inizio articolo, è stato creato un parametro in Power Query che permette di modificare la tabella *SampleData* in modo da riempirla con dei dati fake creati da noi e che portano ad una stima di intercetta con p-value più grandi, che permettono di osservare il cambiamento di significatività sulla base della selezione della soglia di non rifiuto. A tale scopo, si proceda come in figura 7.

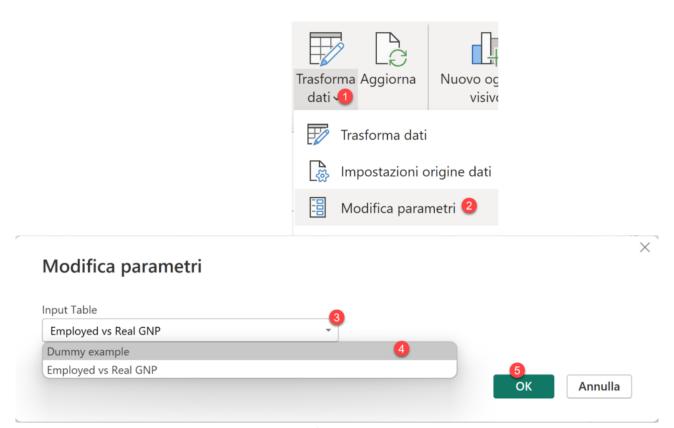


Figura 7

Si clicchi, infine, su "Applica modifiche" in Power BI Desktop. Ecco come si presenta adesso la tabella SampleData:

Input 🔻	Output 🔽	Additional Generic Column
1,000	3,000	
2,000	6,000	
3,000	7,000	
4,000	10,000	
5,000	11,000	

Figura 8



www.kubisco.com



info@kubisco.com



Pag. 10 di 13











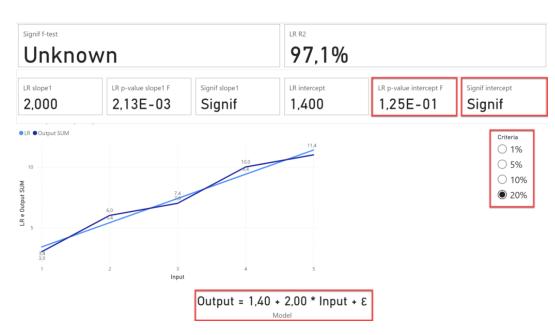


Ed ecco come si presenta il report di figura 2:



Figura 9

Come si vede in figura 9, il *p-value* dell'intercetta è 0,125 cioè il 12,5%. Questo valore è non significativo (rifiuto ipotesi nulla) quando la soglia di non rifiuto è al 1% come mostrato in figura. Si noti che, in questo caso, l'intercetta non compare nell'equazione del modello. Ecco cosa succede se seleziona una soglia che porta al non rifiuto (20%).





www.kubisco.com



info@kubisco.com



Pag. **11** di **13**















Figura 10

Se ci fosse interesse per l'uso della R-visual in Power BI Desktop fatecelo sapere nei commenti e faremo un articolo al riguardo, come anticipazione, ecco un paio di schermate.

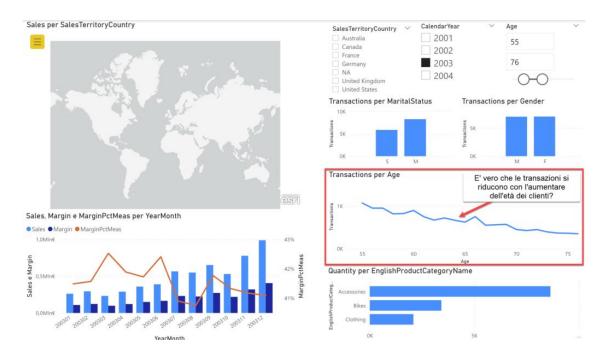
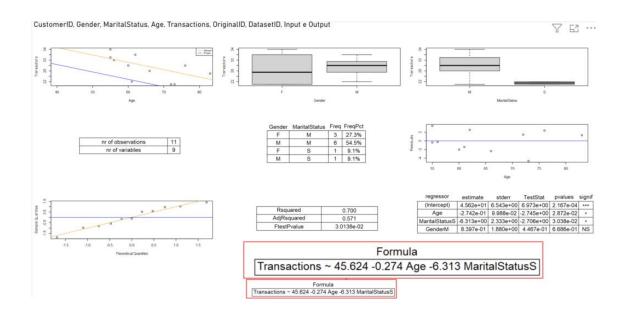


Figura 11





www.kubisco.com



info@kubisco.com



Pag. 12 di 13















Figura 12

Conclusioni

È possibile calcolare i p-value dell'intercetta e delle pendenze stimate da LINESTX. Purtroppo, ad oggi, non è possibile farlo per il test F. Speriamo Microsoft crei la funzione DAX F.DIST.RT che permetterebbe di farlo. Ciò porterebbe alla possibilità di fare tutto in DAX. Tuttavia, restano più comodi i software statistici che possono essere usati dentro Power BI Desktop. Ciò che porta al desideri odi fare tutto in DAX è la bellezza dei report Power BI rispetto a quelli in R e la flessibilità del data model. Si veda questo articolo per apprezzare la possibilità di usare slicer in Power BI Desktop e osservare l'intercetta e la pendenza cambiare in real-time nel report.

file .pbixDownload

Autore del Post



₩kubisco

<u>Francesco Bergamaschi e David Bianconi</u>

See author's posts











