



METODI STATISTICI PER LA BIOINGEGNERIA (B)

PARTE 11: REGOLARIZZAZIONE

A.A. 2025-2026

Prof. Martina Vettoretti



SCOMPOSIZIONE DELL'ERRORE DI PREDIZIONE



- ightharpoonup Relazione vera tra X_1 , X_2 , ..., X_m e Y: $Y=f(X_1,...,X_m)+\delta$, $E[\delta]=0$, $Var(\delta)=\sigma_\delta^2$
- ightharpoonup Approssimazione tramite regressione lineare multipla: $Y=eta_0+eta_1X_1+\cdots+eta_mX_m+arepsilon$
- ightharpoonup Stime dei parametri: \hat{eta}_0 , \hat{eta}_1 , ..., \hat{eta}_m
- ightharpoonup Stimatore di Y: $\hat{Y}=\hat{eta}_0+\hat{eta}_1X_1+\cdots+\hat{eta}_mX_m$
- \succ Errore quadratico di predizione del modello di regressione lineare multipla: $(Y-\widehat{Y})^2$
- Scomposizione dell'errore quadratico medio:

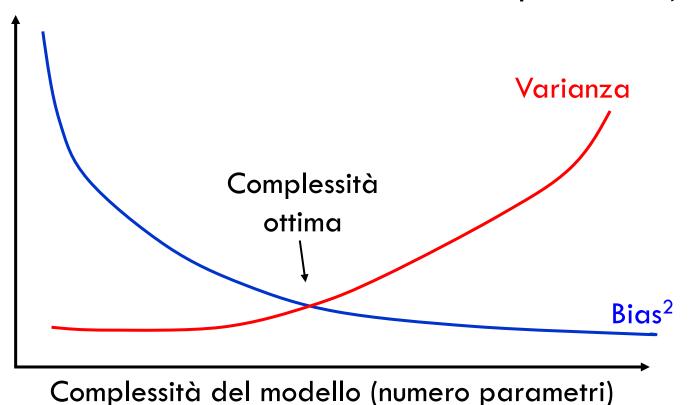
$$E[(Y-\hat{Y})^2] = Var(\delta) + (E[\hat{Y}]-Y)^2 + Var(\hat{Y})$$
 Errore irriducibile, non Bias² dello Varianza dello dipende dal modello stimatore di Y



IL COMPROMESSO TRA BIAS E VARIANZA



Figure Idealmente vorremmo minimizzare sia il bias che la varianza dello stimatore \widehat{Y} . Essi però variano in direzioni opposte al variare della complessità del modello Figure occorre trovare un compromesso (trade-off)



Per minimizzare l'errore di predizione del modello conviene limitare la complessità del modello, aumentando un po' il bias per mantenere bassa la varianza.



METODI DI SHRINKAGE



- > Introducono **bias nelle stime** dei coefficienti $\widehat{m{\beta}}$ (e quindi di $\widehat{m{Y}}$) al fine di mantenere bassa la varianza di $\widehat{m{Y}}$ (predizione meno incerta).
- \triangleright Il bias viene introdotto mediante una **regolarizzazione** della stima dei coefficienti β che penalizza valori grandi in valore assoluto delle stime dei coefficienti $\beta \rightarrow$ Le stime sono «ristrette» in valore assoluto verso lo zero.
- ➤ I metodi di shrinkage sono particolarmente utili quando abbiamo tante variabili correlate tra loro (con stima non regolarizzata potremmo avere coefficienti grandi in valore assoluto e di segno opposto per variabili tra loro correlate).



STIMA STANDARD VS STIMA REGOLARIZZATA



> Problema di ottimizzazione standard (senza regolarizzazione):

$$\widehat{\boldsymbol{\beta}}$$
: = argmin($F(\boldsymbol{\beta})$)

Quando usiamo i minimi quadrati lineari:

$$F(\boldsymbol{\beta}) := (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta})^T (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta}) = SSE$$

Problema di ottimizzazione con regolarizzazione:

$$\widehat{\boldsymbol{\beta}}_{reg,\lambda} := \underset{\boldsymbol{\beta}}{\operatorname{argmin}} (F(\boldsymbol{\beta}) + P(\lambda; \beta_{1}, \beta_{2}, \dots \beta_{m}))$$

Termine di penalità: penalizza valori grandi (in modulo) dei coefficienti

Parametro di regolarizzazione:

regola il grado di regolarizzazione.

Nota: nella formulazione classica non si regolarizza l'intercetta, β_0 .



REGOLARIZZAZIONE L2 O REGRESSIONE RIDGE (1/2)



➤ Regolarizzazione L2: il termine di penalità è il quadrato della norma 2 del vettore dei coefficienti di regressione (intercetta esclusa) → regressione Ridge

$$\widehat{\boldsymbol{\beta}}_{Ridge,\lambda} := \underset{\boldsymbol{\beta}}{\operatorname{argmin}} (F(\boldsymbol{\beta}) + \lambda \cdot \sum_{j=1}^{m} \beta_{j}^{2})$$

$$F(\boldsymbol{\beta}) := (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta})^{T} (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta}) = SSE$$

- > Si penalizzano valori elevati (in modulo) dei coefficienti di regressione.
- Effetto: le stime dei parametri vengono «schiacciate» verso lo 0. Tuttavia nessun coefficiente viene posto a 0 → complessità del modello invariata



EFFETTO DEL PARAMETRO DI REGOLARIZZAZIONE



- $\triangleright \lambda$ è uno scalare positivo ($\lambda \ge 0$).
- ightharpoonup Più λ è grande, maggiore è la penalità imposta su valori grandi dei coefficienti.
 - lacktriangle Quando $\lambda o \infty$, $\hat{eta}_i o 0$
 - lacktriangle Quando $\lambda=0$ ightarrow le stime \hat{eta}_j saranno identiche a quelle del modello non regolarizzato
- \triangleright Il valore ottimale di λ va trovato per ogni modello.
- \triangleright Valori come λ , che non sono parametri del modello di Y (non sono coefficienti β), ma hanno un impatto sulla forma finale del modello si chiamano **iperparametri**.



STIMATORE RIDGE



 \succ Fissato λ si può dimostrare che la soluzione al problema di regressione Ridge è:

$$\widehat{\boldsymbol{\beta}}_{Rid,ge,\lambda} = (\boldsymbol{X}^T\boldsymbol{X} + \lambda \cdot \boldsymbol{Q})^{-1}\boldsymbol{X}^T\boldsymbol{Y}$$

dove Q è la matrice identità m+1 x m+1 con uno zero sulla diagonale in corrispondenza dell'intercetta (parametro non penalizzato).

Esempio: $\beta = [\beta_0, \beta_1, \beta_2]$ con β_0 intercetta $\Rightarrow Q = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

ESEMPIO: EFFETTO DELLA REGOLARIZZAZIONE L2



Modello per la predizione del diametro della componente acetabolare della protesi all'anca. Variabili indipendenti: altezza, girovita, lunghezza piede, età, sesso, patologia (2 variabili dummy: frattura e necrosi).

Coefficienti	$\lambda = 0$	$\lambda = 1$	$\lambda = 10$	$\lambda = 100$
Intercetta	46.04	46.53	48.86	51.97
Altezza	6.91	6.19	3.58	1.02
Girovita	5.61	5.22	3.26	0.85
Lunghezza piede	4.35	4.27	3.10	1.02
Età	-0.56	-0.61	-0.71	-0.36
Sesso	-0.49	-0.23	0.95	1.44
Frattura	-0.26	-0.27	-0.28	-0.18
Necrosi	0.11	0.10	0.03	-0.01



REGOLARIZZAZIONE L1 O REGRESSIONE LASSO



➤ **Regolarizzazione L1**: il termine di penalità è la norma 1 del vettore dei coefficienti di regressione (intercetta esclusa) → **regressione LASSO** (Least Absolute Shrinkage and Selection Operator).

$$\widehat{\boldsymbol{\beta}}_{LASSO,\lambda} := \underset{\boldsymbol{\beta}}{\operatorname{argmin}} (F(\boldsymbol{\beta}) + \lambda \cdot \sum_{j=1}^{m} |\beta_j|)$$
$$F(\boldsymbol{\beta}) := (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta})^T (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta}) = SSE$$

- Effetto: le stime di tutti i parametri vengono «schiacciate» verso lo 0 e, per valori sufficientemente grandi di λ , <u>i coefficienti più piccoli vengono posti a 0 \rightarrow si riduce la complessità del modello</u>.
- \triangleright Il valore ottimo per l'iperparametro λ va trovato.

ESEMPIO: EFFETTO DELLA REGOLARIZZAZIONE L1



Modello per la predizione del diametro della componente acetabolare della protesi all'anca. Variabili indipendenti: altezza, girovita, lunghezza piede, età, sesso, patologia (2 variabili dummy: frattura e necrosi).

Coefficienti	$\lambda = 0$	$\lambda = 0.01$	$\lambda = 0.1$	$\lambda = 0.2$	$\lambda = 1$
Intercetta	46.04	46.52	49.23	52.41	53.86
Altezza	6.91	6.45	3.98	0.62	0
Girovita	5.61	5.16	2.29	0	0
Lunghezza piede	4.35	3.93	2.23	0.26	0
Età	-0.56	-0.48	0	0	0
Sesso	-0.49	-0.15	0.94	2.26	0
Frattura	-0.26	-0.24	0	0	0
Necrosi	0.11	0.05	0	0	0



REGOLARIZZAZIONE ELASTIC NET (1/2)



Regolarizzazione Elastic Net: termine di penalità dato da una combinazione lineare dei termini di penalità delle regolarizzazioni L1 e L2.

$$\widehat{\boldsymbol{\beta}}_{ENet,\lambda,\alpha} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} (F(\boldsymbol{\beta}) + \lambda \cdot \sum_{j=1}^{\infty} [(1 - \alpha) \cdot \beta_j^2 + \alpha \cdot |\beta_j|])$$

$$F(\boldsymbol{\beta}) = (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta})^T (\boldsymbol{Y} - \boldsymbol{X} \cdot \boldsymbol{\beta}) = SSE$$

- \triangleright Iperparametri: λ , α
 - λ regola il grado di regolarizzazione
 - $\alpha \in [0,1]$ indica quanto prevale il termine di penalità L1 su quello L2.
- > I valori ottimi di entrambi gli iperparametri vanno trovati.



REGOLARIZZAZIONE ELASTIC NET (1/2)



$$\widehat{\boldsymbol{\beta}}_{ENet,\lambda,\alpha} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} (F(\boldsymbol{\beta}) + \lambda \cdot \sum_{j=1}^{m} [(1 - \alpha) \cdot \beta_j^2 + \alpha \cdot |\beta_j|])$$

- \triangleright Cosa succede per $\alpha=0$?
- \triangleright E per $\alpha = 1$?



QUALE REGOLARIZZAZIONE ?



- ➤ Ridge vs LASSO:
 - Ridge non riduce la complessità del modello, mentre LASSO lo fa.
 - Se ipotizziamo che il dataset contenga pochi predittori forti e molte variabili «rumore» (non associate all'outcome) → preferiamo LASSO
 - Se ipotizziamo non esserci variabili indipendenti dall'outcome > preferiamo Ridge
- Elastic net consente di avere sia i vantaggi della Ridge che della LASSO
 - Tuttavia richiede 2 iperparametri da stimare → la ricerca degli iperparametri ottimi può diventare computazionalmente onerosa



TUNING DEGLI IPERPARAMETRI

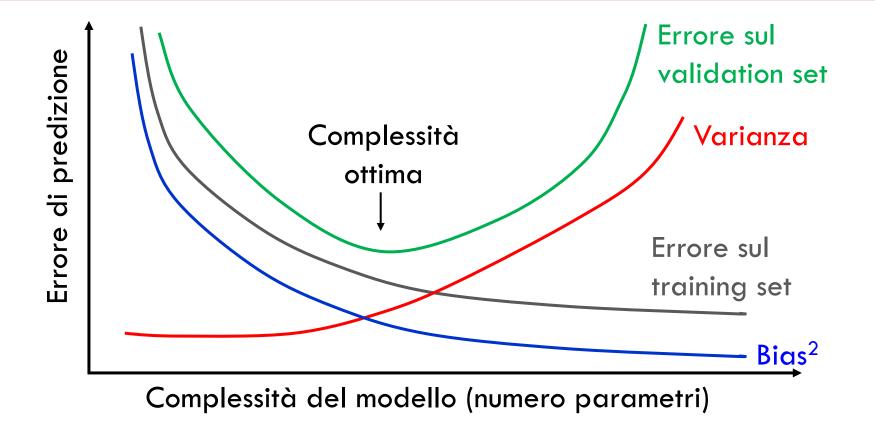


- > Non è possibile stabilire a priori il valore degli iperparametri.
- > Non esistono stimatori per calcolare il valore degli iperparametri.
- \triangleright Ricerca empirica del valore ottimo degli iperparametri: si testano diversi valori per λ (ed eventualmente α) e si seleziona il valore (o la combinazione di valori) che minimizza l'errore di predizione del modello (MSE o RMSE) su un nuovo set di dati detto <u>validation set</u>.
- \triangleright Training set: set di dati utilizzati per stimare i coefficienti β del modello.
- Validation set: nuovo set di dati, non utilizzato per stimare i coefficienti del modello β, ma utilizzato per misurare l'errore di predizione del modello al variare degli iperparametri.



ERRORE DI PREDIZIONE SU DATI DI VALIDAZIONE





Dobbiamo stimare i valori degli iperparametri che minimizzano l'errore di predizione su un set di dati di validazione, non utilizzato per stimare i parametri $oldsymbol{eta}$ del modello.



OTTIMIZZAZIONE DEGLI IPERPARAMETRI MEDIANTE VALIDATION SET (1/2)



- > Dividiamo i dati a disposizione in due set:
 - Training set: set di dati impiegato per stimare i coefficienti del modello di regressione.
 - Validation set: set di dati impiegato per valutare la miglior combinazione di iperparametri.
- > Scegliamo una griglia di possibili valori per ciascun iperparametro:

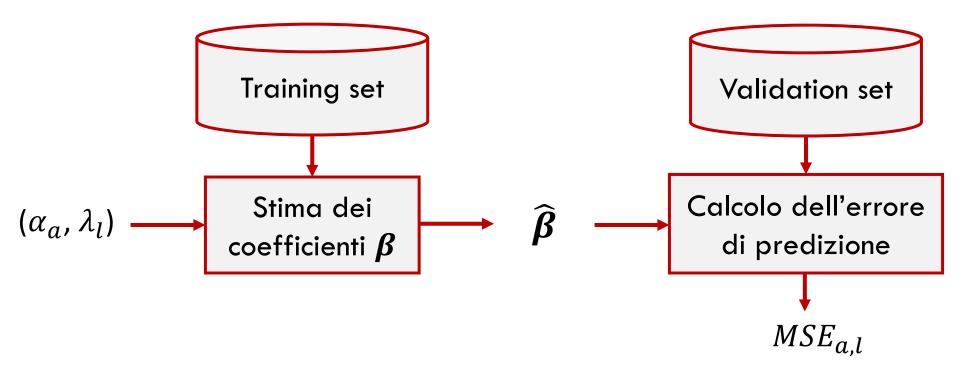
$$\alpha_1, \alpha_2, \dots, \alpha_A$$
 $\lambda_1, \lambda_2, \dots, \lambda_L$

A x L possibili coppie di
valori per gli iperparametri



OTTIMIZZAZIONE DEGLI IPERPARAMETRI MEDIANTE VALIDATION SET (2/2)





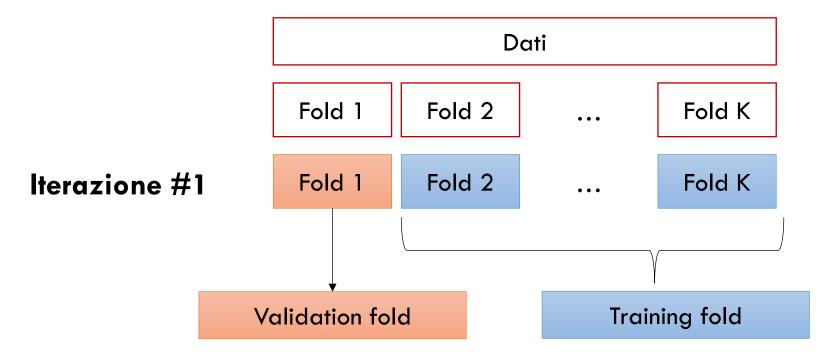
- \succ Scelgo gli iperparametri ($lpha_{opt}$, λ_{opt}) per cui $MSE_{a,l}$ è minimo.
- ightharpoonup Uso ($lpha_{opt}$, λ_{opt}) per stimare i coefficienti $m{\beta}$ del modello finale su training+validation set.
- Limitazione: siamo sicuri che il particolare split dei dati che facciamo per ricavare il training set e il validation set non influenzi il risultato?



K-FOLD CROSS-VALIDATION (1/2)



- > Dividiamo i dati in K sottoinsiemi, detti fold. Realizziamo K iterazioni.
- Iterazione 1:
 - Fold 2-K \rightarrow Dati da usare per il training del modello (stima dei parametri β)
 - Fold 1 \rightarrow Set di validazione per calcolare le performance del modello (MSE)

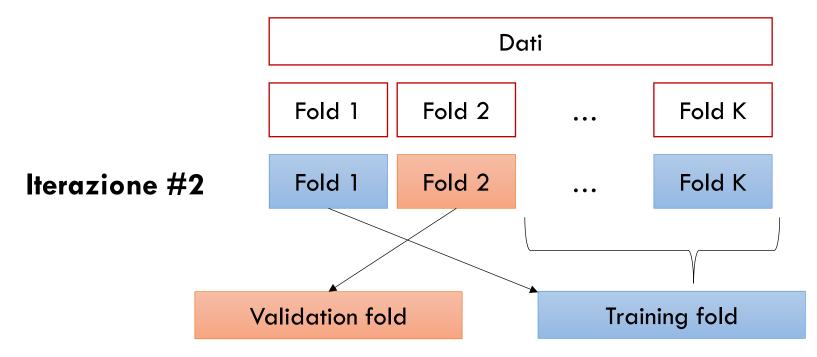




K-FOLD CROSS-VALIDATION (2/2)



- Dividiamo i dati in K sottoinsiemi, detti fold. Realizziamo K iterazioni.
- Iterazione 2:
 - Fold 1, 3-K \rightarrow Dati da usare per il training del modello (stima dei parametri β)
 - Fold 2 \rightarrow Set di validazione per calcolare le performance del modello (MSE)



In generale alla iesima iterazione
usiamo la fold i-esima
come set di dati per la
validazione, l'unione
delle altre fold come
set di dati per il
training del modello.



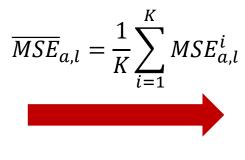
K-FOLD CROSS-VALIDATION PER L'OTTIMIZZAZIONE DEGLI IPERPARAMETRI



 $N = A \times L$ diverse combinazioni di iperparametri candidate.

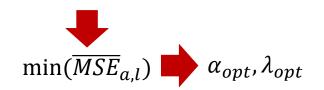
- Ad ogni iterazione della K-fold cross-validation, alleniamo sulle fold di training gli N modelli corrispondenti alle N combinazioni di iperparametri. Testiamo ciascun modello sulla fold di validazione, calcolando l'MSE.
- 2. Calcoliamo la media dell'MSE sulle fold di validazione per ciascuna combinazione di iperparametri.

MSE sulla fold i-esima					
λ_1 λ_2				λ_{L}	
α_1	MSE ⁱ _{1,1}	MSE ⁱ _{1,2}		MSE ⁱ _{1,L}	
α_2	MSE ⁱ _{2,1}	MSE ⁱ _{2,2}	•••	MSE ⁱ _{2,L}	
α_{A}	MSE ⁱ _{A,1}	MSE ⁱ _{A,2}		MSE ⁱ _{A,L}	



	MSE medio sulle K fold				
	λ_1	λ_2	•••	λ_{L}	
α_1	$\overline{\text{MSE}}_{1,1}$	$\overline{\text{MSE}}_{1,2}$	•••	$\overline{\text{MSE}}_{1,L}$	
α_2	$\overline{\text{MSE}}_{2,1}$	$\overline{\text{MSE}}_{2,2}$	•••	$\overline{\text{MSE}}_{2,L}$	
•••					
α_{A}	$\overline{\text{MSE}}_{A,1}$	$\overline{\text{MSE}}_{A,2}$		$\overline{\text{MSE}}_{A,L}$	

3. Combinazione ottima di iperparametri: quella che minimizza l'MSE medio sulle fold di validazione.





NOTE



- \triangleright Una volta scelti i valori degli iperparametri ottimi, α_{opt} , λ_{opt} , si utilizzano questi valori per stimare i coefficienti β del modello finale sull'intero set di dati a disposizione (tutte le fold insieme).
- > Scelta del numero di fold K:
 - Valori tipici di K sono 5 o 10.
 - Se K è pari al numero di osservazioni nel dataset n → leave-one-out cross-validation.
 - Approccio computazionalmente oneroso quando n è grande.
- Occorre fare attenzione alla scelta della griglia di valori degli iperparametri da testare.



SCELTA DELLA GRIGLIA DEGLI IPERPARAMETRI



- \blacktriangleright Sappiamo che $lpha_{opt}$ è compreso tra 0 e 1
 - \rightarrow possibile griglia per α : {0, 0.1, 0.2, 0.3,...,0.9, 1}
- \triangleright Sappiamo che $\lambda_{opt}>0$, ma non ne conosciamo l'ordine di grandezza
 - \rightarrow possibile griglia iniziale per λ : {10-p, 10-p+1, ...,1, 10, ..., 10p}
 - Se al primo giro di ottimizzazione finisco in uno degli **estremi della griglia** dei $\lambda \rightarrow$ non ho davvero trovato l'ottimo! Espando la griglia dei λ .

• Quando λ_{opt} cade all'interno della griglia dei λ possiamo fermarci qui o decidere di testare una griglia più fitta di valori di λ attorno al valore ottimo trovato.

$\lambda_{\text{min,new}}$	•••	$\lambda_{\text{min,old}}$	•••	λ_{\max}	$\lambda_{\text{min,new}} < \lambda_{\text{opt}} < \lambda_{\text{max}} \rightarrow \mathbf{OK}$
----------------------------	-----	----------------------------	-----	------------------	--



ESEMPIO



- Modello per la predizione del diametro della componente acetabolare della protesi all'anca. Variabili indipendenti: altezza, girovita, lunghezza piede, età, sesso, patologia (2 variabili dummy: frattura e necrosi).
- Regolarizzazione elastic net con ottimizzazione degli iperparametri mediante 5-fold cross-validation.
- \triangleright Griglia di valori per α : {0.1, 0.4, 0.7, 1}
- \triangleright Griglia di valori per λ : {10⁻¹⁰,10⁻⁷,10⁻⁴,10⁻¹,10²,10⁵,10⁸}



ESEMPIO: RISULTATI



Minimo dell'MSE medio sulle fold di validazione

$$\alpha = \frac{\text{MSE_CV} = }{2.1617} = \frac{2.1617}{2.1617} = \frac{2.1615}{3.6503} = \frac{7.4811}{7.4811} = \frac{7.4811}{7.4811} = \frac{7.4811}{7.4811} = \frac{7.4811}{2.1617} = \frac{2.1617}{2.1617} = \frac{2.1616}{2.1617} = \frac{3.3916}{2.9100} = \frac{7.4811}{7.4811} = \frac{7.4811}{7.48$$

$$\alpha_{opt} = 0.1$$

$$\alpha_{opt} = 0.1$$

$$\lambda_{opt} = 10^{-4}$$



ESEMPIO: RISULTATI



No regolarizzazione	$lpha=0$. 1, $\lambda=10^{-4}$
46.04	46.52
6.91	6.88
5.61	5.59
4.35	4.35
-0.56	-0.56
-0.49	-0.48
-0.27	-0.27
0.11	0.11
	46.04 6.91 5.61 4.35 -0.56 -0.49 -0.27

In questo caso il modello regolarizzato si discosta poco dal modello non regolarizzato. In generale però la regolarizzazione può impattare in maniera importante i risultati, soprattutto per modelli con tanti predittori.



INFLUENZA DELLA SCALA DELLE VARIABILI



- $\succ \lambda$ rappresenta il **grado di regolarizzazione**, che idealmente vorremmo essere **uniforme** per tutte le variabili indipendenti.
- Attenzione: se le variabili hanno scale diverse, il grado di regolarizzazione non è lo stesso per tutte le variabili!
 - A parità di impatto su Y, le variabili che assumono valori più piccoli tenderanno ad avere coefficienti più grandi → maggiore penalizzazione.
- Per rendere uniforme il grado di regolarizzazione per tutte le variabili, prima di applicare la regressione regolarizzata, è buona norma normalizzare le variabili indipendenti per riportarle ad avere la stessa scala.
- Approcci di normalizzazione più diffusi:
 - Standardizzazione
 - Min-max scaling



STANDARDIZZAZIONE



$$Z = \frac{X - \bar{X}}{S_X}$$

- X: variabile originale
- Z: variabile standardizzata
- lacktriangle $ar{X}$: media campionaria di X
- S_X : deviazione standard campionaria di X

> Dopo la standardizzazione, tutte le variabili avranno media 0 e varianza 1.



MIN-MAX SCALING



$$X' = \frac{X - X_{min}}{X_{max} - X_{min}}$$

- X: variabile originale
- X': variabile trasformata
- X_{min}: valore minimo della variabile originale
- X_{max}: valore massimo della variabile originale

- > Dopo min-max scaling, tutte le variabili avranno range tra 0 e 1.
- P Questo approccio è preferibile quando abbiamo molte variabili qualitative che entreranno nel modello come variabili binarie di valori 0 o 1.