



A.D. 1308
unipg

DIPARTIMENTO
DI INGEGNERIA

Progetto di
Signal Processing and Optimization for Big Data

Corso di Laurea in Ingegneria Informatica e Robotica

Curriculum Data Science

DIPARTIMENTO DI INGEGNERIA

docente

Prof. Paolo BANELLI

Confronto tra Algoritmi di Regressione Lasso Sviluppati da Zero

363433 **Gian Marco Ferri** gianmarco.ferri@studenti.unipg.it

Indice

1	Introduzione	2
1.1	Contesto Generale	2
2	Analisi Teorica	3
2.1	LASSO	3
2.2	Soft-Thresholding	3
2.3	ADMM	4
2.4	ADMM Distribuito	6
3	Implementazione	8
3.1	Dataset	8
4	Comparazioni	9

1 Introduzione

Questo progetto affronta la risoluzione del problema della regressione con regolarizzazione Lasso tramite l'implementazione e il confronto di diversi algoritmi. In particolare, sono stati sviluppati in Matlab tre approcci distinti: ISTA, ADMM e una versione simulata di ADMM distribuito su più agenti.

Nella parte finale della documentazione vengono presentati dei confronti tra i tre algoritmi in termini di tempi di calcolo, numero di iterazioni e grafici che illustrano l'andamento della convergenza durante le iterazioni.

1.1 Contesto Generale

La regressione LASSO (Least Absolute Shrinkage and Selection Operator) è una tecnica fondamentale nel campo della modellistica statistica e del machine learning.

Questa metodologia permette di affrontare problemi come l'overfitting e la selezione delle variabili, offrendo una strategia efficace per sviluppare modelli affidabili e generalizzabili.

Con la crescente complessità degli algoritmi, il rischio di overfitting, ossia la tendenza di un modello a modellare troppo fedelmente i dati di addestramento includendo anche il rumore, è diventato sempre più rilevante. Il metodo di regolarizzazione LASSO affronta tale problema introducendo una penalizzazione sulla somma dei valori assoluti dei coefficienti delle variabili nel modello. L'intensità di questa penalità è regolata da un iperparametro.

La caratteristica distintiva della regressione LASSO è proprio l'utilizzo della regolarizzazione di tipo $L1$, che incentiva la sparsità nei coefficienti stimati. In pratica, ciò comporta che LASSO tende ad azzerare molti coefficienti, selezionando automaticamente solo le variabili più significative per la previsione.

Questa sua proprietà la differenzia da altre tecniche come la regressione Ridge, che usa una penalità $L2$, e la rende particolarmente utile quando si ha a che fare con dataset che hanno un numero elevato di feature.

2 Analisi Teorica

In questo capitolo si propone un'analisi dettagliata del problema dal punto di vista sia teorico che matematico, con l'obiettivo di approfondire la comprensione degli algoritmi che verranno implementati in seguito.

Viene presentata una descrizione esaustiva della regressione LASSO, introducendo il concetto della regolarizzazione L1 e discutendo le soluzioni teoriche relative al problema di ottimizzazione collegato.

Si prendono inoltre in considerazione le principali varianti dell'algoritmo, ovvero ISTA, ADMM e ADMM distribuito, illustrandone le basi teoriche necessarie per comprendere appieno le applicazioni pratiche che saranno affrontate nei capitoli successivi.

2.1 LASSO

Il LASSO (Least Absolute Shrinkage and Selection Operator) è un metodo di regressione che applica una penalizzazione di tipo $L1$ alla funzione di ottimizzazione. Tale penalità, proporzionale alla somma dei valori assoluti dei coefficienti, induce sparsità nel vettore dei parametri, portando molti di essi ad annullarsi. Questo meccanismo consente non solo di selezionare automaticamente le variabili più rilevanti, ma anche di ottenere un modello più semplice, riducendo efficacemente il fenomeno dell'overfitting.

Tale risultato è ottenuto risolvendo un problema di ottimizzazione che coinvolge la minimizzazione di una funzione di costo, la quale comprende sia il termine di regressione che la penalità $L1$.

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - A\mathbf{x}\|_2^2 + \alpha \|\mathbf{x}\|_1$$

L'espressione riportata descrive il problema di ottimizzazione relativo alla regressione LASSO. Nell'equazione, \mathbf{x} indica il vettore dei coefficienti del modello, \mathbf{y} rappresenta il vettore delle variabili dipendenti, A corrisponde alla matrice delle variabili indipendenti (features) e α è l'iperparametro che controlla l'intensità della regolarizzazione.

2.2 Soft-Thresholding

Nell'algoritmo ISTA (Iterative Shrinkage Thresholding Algorithm), la regressione LASSO viene risolta utilizzando un approccio iterativo basato sul Gradient Descent.

In questo contesto, il Gradient Descent viene impiegato per affrontare il problema di ottimizzazione precedentemente illustrato, aggiornando progressivamente i pesi del modello nella direzione opposta al gradiente della funzione costo, fino al raggiungimento della convergenza.

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \gamma \nabla \left[\frac{1}{2} \|\mathbf{y} - A\mathbf{x}^{(k)}\|_2^2 + \alpha \|\mathbf{x}^{(k)}\|_1 \right]$$

Poiché la norma L_1 non è differenziabile, si ricorre al concetto di sub-gradiente attraverso l'impiego dell'operatore di Soft-Thresholding.

$$[\text{sub}\nabla f(x_0)]_i = \begin{cases} \frac{\partial f(x_0)}{\partial x_i}, & \text{se differenziabile in } x_0 \\ \left[\frac{\partial f(x_0^-)}{\partial x_i}, \frac{\partial f(x_0^+)}{\partial x_i} \right], & \text{se non differenziabile in } x_0 \end{cases}$$

Nel contesto della regressione LASSO, questo si traduce nella seguente espressione:

$$\begin{aligned} \text{sub}\nabla\Phi &= \begin{cases} \gamma\alpha + (u_j - x_j) & \text{se } u_j > 0 \\ -\gamma\alpha + (u_j - x_j) & \text{se } u_j < 0 \\ [-\gamma\alpha + (u_j - x_j), \gamma\alpha + (u_j - x_j)] & \text{se } u_j = 0 \end{cases} \\ \Rightarrow &\begin{cases} u_j = x_j - \gamma\alpha & \text{se } x_j > \gamma\alpha \\ u_j = x_j + \gamma\alpha & \text{se } x_j < -\gamma\alpha \\ u_j = 0 & \text{se } -\gamma\alpha < x_j < \gamma\alpha \end{cases} \end{aligned}$$

Questa operazione, nota come Soft-Thresholding, fornisce una soluzione analitica per l'aggiornamento dei coefficienti nel caso di penalità L_1 , dove il gradiente classico non è applicabile.

Il processo iterativo prosegue fino a quando la norma della differenza tra i coefficienti di due iterazioni successive diventa inferiore a una soglia di tolleranza prestabilita, indicando il raggiungimento della convergenza.

2.3 ADMM

L'Alternating Direction Method of Multipliers (ADMM) è un algoritmo di ottimizzazione primal-dual che risolve problemi complessi scomponendoli in una serie di sottoproblemi più semplici e gestibili, che vengono risolti in modo iterativo attraverso l'aggiornamento alternato delle variabili primali e duali.

Per applicare l'ADMM al problema LASSO, è necessario riformularlo introducendo una variabile ausiliaria (o *slack variable*) \mathbf{z} . Questa scomposizione permette di separare l'ottimizzazione della funzione di loss originale dalla penalità L_1 , facilitando la risoluzione attraverso sottoproblemi indipendenti.

La formulazione risultante è la seguente:

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{z}} \quad & \frac{1}{2} \|A\mathbf{x} - \mathbf{y}\|_2^2 + \alpha \|\mathbf{z}\|_1 \\ \text{s.t.} \quad & \mathbf{x} - \mathbf{z} = \mathbf{0} \end{aligned}$$

Il framework ADMM scompone il problema di ottimizzazione in tre distinti sottoproblemi, le cui soluzioni vengono calcolate in modo alternato e iterativo fino al soddisfacimento di un criterio di convergenza.

Di seguito viene presentata la formulazione standard in forma scalata dell'algoritmo:

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \underset{\mathbf{x}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|A\mathbf{x} - \mathbf{y}\|_2^2 + \frac{\rho}{2} \|\mathbf{x} - \mathbf{z}^{(k)} + \mathbf{u}^{(k)}\|_2^2 \right\} \\ \mathbf{z}^{(k+1)} &= \underset{\mathbf{z}}{\operatorname{argmin}} \left\{ \alpha \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\mathbf{x}^{(k+1)} - \mathbf{z} + \mathbf{u}^{(k)}\|_2^2 \right\} \\ \mathbf{u}^{(k+1)} &= \mathbf{u}^{(0)} + \sum_{i=1}^{k+1} \|\mathbf{x}^{(i)} - \mathbf{z}^{(i)}\|_2^2 \end{aligned}$$

Il primo step ammette una soluzione in forma chiusa, ottenuta ponendo a zero la derivata della funzione obiettivo:

$$\begin{aligned} \nabla_{\mathbf{x}} \left(\frac{1}{2} \|A\mathbf{x} - \mathbf{y}\|_2^2 + \frac{\rho}{2} \|\mathbf{x} - \mathbf{z}^{(k)} + \mathbf{u}^{(k)}\|_2^2 \right) &= \mathbf{0} \\ A^T (A\mathbf{x} - \mathbf{y}) + \rho(\mathbf{x} - \mathbf{z}^{(k)} + \mathbf{u}^{(k)}) &= \mathbf{0} \\ (A^T A + \rho I) \mathbf{x} &= A^T \mathbf{y} + \rho(\mathbf{z}^{(k)} - \mathbf{u}^{(k)}) \\ \rightarrow \mathbf{x}^{(k+1)} &= (A^T A + \rho I)^{-1} (A^T \mathbf{y} + \rho(\mathbf{z}^{(k)} - \mathbf{u}^{(k)})) \end{aligned}$$

Per il secondo sottoproblema, la funzione obiettivo non è differenziabile a causa della norma $L1$. La sua soluzione, tuttavia, può essere ottenuta in forma chiusa imponendo le condizioni di sub-gradiente, le quali, per la norma $L1$, si traducono nell'applicazione dell'operatore di soft-thresholding.

In conclusione, la formulazione del problema LASSO mediante ADMM in forma scalata è la seguente:

$$\begin{aligned} \mathbf{x}^{(k+1)} &= (A^T A + \rho I)^{-1} (A^T \mathbf{y} + \rho(\mathbf{z}^{(k)} - \mathbf{u}^{(k)})) \\ \mathbf{z}^{(k+1)} &= S_{\frac{\alpha}{\rho}}(\mathbf{x}^{(k+1)} + \mathbf{u}^{(k)}) \\ \mathbf{u}^{(k+1)} &= \mathbf{u}^{(k)} + \mathbf{x}^{(k+1)} - \mathbf{z}^{(k+1)} \end{aligned}$$

I parametri ρ e α rivestono un ruolo critico nell'algoritmo, poiché la loro scelta influenza fondamentalemente il bilanciamento tra velocità di convergenza e robustezza della soluzione.

Nonostante questa sensibilità ai parametri, l'utilizzo dell'ADMM per il problema LASSO offre notevoli vantaggi.

Il suo vantaggio principale deriva dalla scomposizione del problema originale in sottoproblemi più semplici e gestibili, spesso risolvibili in forma chiusa. Questa architettura garantisce all'algoritmo una notevole stabilità e consente di raggiungere la soluzione ottimale in un numero di iterazioni generalmente contenuto.

Un ulteriore vantaggio decisivo è la sua capacità di essere parallelizzato e implementato in modo distribuito, che lo rende particolarmente adatto per scalare su problemi di grandi dimensioni.

2.4 ADMM Distribuito

La versione distribuita dell'ADMM estende il framework classico per gestire problemi su larga scala distribuendo il carico computazionale tra diversi nodi (o agenti). L'ottimizzazione del consenso viene implementata partizionando il dataset tra gli agenti, ciascuno dei quali mantiene una copia locale delle variabili di ottimizzazione. In questo modo, il problema globale viene scomposto rispetto ai dati e risolto in maniera collaborativa attraverso lo scambio di informazioni tra i nodi.

Questo approccio è possibile quando la funzione obiettivo presenta una struttura additiva rispetto ai dati, il che permette di scomporre il problema di ottimizzazione in sottoproblemi indipendenti che possono essere risolti in parallelo.

La formulazione standard del problema di ottimizzazione del consenso, in cui N agenti devono minimizzare una funzione obiettivo globale scomponibile, è la seguente:

$$\begin{aligned} \min_{(x_1, \dots, x_N), \mathbf{z}} \quad & \sum_{i=1}^N f_i(\mathbf{x}) + g(\mathbf{z}) \\ \text{s.t.} \quad & \mathbf{x}_i - \mathbf{z} = \mathbf{0}, \quad i = 1, \dots, N \end{aligned}$$

Il problema Lasso è intrinsecamente distribuibile grazie alla natura separabile della sua funzione obiettivo. Infatti, sia le osservazioni \mathbf{y} che la matrice A possono essere partizionati in N blocchi, rendendo possibile la scomposizione del problema originale in sottoproblemi indipendenti.

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}, \quad A\mathbf{x} = \begin{bmatrix} A_1 \\ \vdots \\ A_N \end{bmatrix} \mathbf{x} = \begin{bmatrix} A_1\mathbf{x} \\ \vdots \\ A_N\mathbf{x} \end{bmatrix}$$

Andando a riformulare il problema nella forma di un'ottimizzazione al consenso, si ottiene la seguente rappresentazione:

$$\begin{aligned} \min_{(x_1, \dots, x_N), \mathbf{z}} \quad & \sum_{i=1}^N \|\mathbf{y}_i - A_i \mathbf{x}_i\|_2^2 + \lambda \|\mathbf{z}\|_1 \\ \text{s.t.} \quad & \mathbf{x}_i - \mathbf{z} = \mathbf{0}, \quad i = 1, \dots, N \end{aligned}$$

La soluzione viene ottenuta applicando la versione scalata dell'ADMM, la cui struttura è la seguente:

$$\begin{aligned} \mathbf{x}_i^{(k+1)} &= \underset{\mathbf{x}_i}{\operatorname{argmin}} \left\{ \|\mathbf{y}_i - A_i \mathbf{x}_i\|_2^2 + \frac{\rho}{2} \|\mathbf{x}_i - \mathbf{z}^{(k)} + \mathbf{u}_i^{(k)}\|_2^2 \right\} \\ \mathbf{z}^{(k+1)} &= \underset{\mathbf{z}}{\operatorname{argmin}} \left\{ g(\mathbf{z}) + \frac{\rho}{2} \sum_{i=1}^N \|\mathbf{x}_i^{(k+1)} - \mathbf{z} + \mathbf{u}_i^{(k)}\|_2^2 \right\} \quad i = 1, \dots, N \\ \mathbf{u}_i^{(k+1)} &= \mathbf{u}_i^{(k)} + (\mathbf{x}_i^{(k+1)} - \mathbf{z}^{(k+1)}) \quad i = 1, \dots, N \end{aligned}$$

Il primo step si compone di N sottoproblemi indipendenti, ciascuno assegnabile a un agente distinto. Come nella formulazione centralizzata, questi ammettono una soluzione in forma chiusa e possono essere quindi calcolati in parallelo.

Il secondo step richiede la raccolta di tutte le variabili primali e duali locali per il calcolo di un valore di consenso globale. Questa operazione, non distribuibile, viene affidata a un fusion center.

Il terzo step può essere eseguito in modo completamente autonomo e parallelo da ciascun agente sulla base dei nuovi valori di \mathbf{x}_i e \mathbf{z} .

Di conseguenza, le iterazioni dell'algoritmo per ogni agente i e per il fusion center sono definite dalle seguenti equazioni:

$$\begin{aligned} \mathbf{x}_i^{(k+1)} &= (A_i^T A_i + \frac{\rho}{2} I)^{-1} (A_i^T \mathbf{y}_i + \frac{\rho}{2} (\mathbf{z}_i^{(k)} - \mathbf{u}_i^{(k)})) \\ \mathbf{z}^{(k+1)} &= \operatorname{Prox}_{\frac{\lambda}{N\rho} \|\cdot\|_1} \{ \hat{\mathbf{x}}^{(k+1)} + \hat{\mathbf{u}}^{(k)} \} = S_{\frac{\lambda}{N\rho}} (\hat{\mathbf{x}}^{(k+1)} + \hat{\mathbf{u}}^{(k)}) \\ \mathbf{u}_i^{(k+1)} &= \mathbf{u}_i^{(k)} + \mathbf{x}_i^{(k+1)} - \mathbf{z}^{(k+1)} \quad i = 1, \dots, N \end{aligned}$$

dove $\hat{\mathbf{x}}^{(k+1)}$ e $\hat{\mathbf{u}}^{(k)}$ sono le medie delle variabili primarie e duali al passo $k + 1$ e k rispettivamente.

Questa tecnica è particolarmente adatta per l'elaborazione di dataset di grandi dimensioni e per scenari in cui i dati sono intrinsecamente distribuiti su più nodi computazionali.

3 Implementazione

Gli algoritmi descritti in precedenza sono stati implementati all'interno della classe `LassoReg`. Al momento dell'istanziatura, è possibile specificare i seguenti parametri: l'algoritmo di ottimizzazione per la fase di training, lo step-size, la tolleranza del criterio di convergenza, il numero massimo di iterazioni e il coefficiente di penalizzazione.

3.1 Dataset

Il dataset utilizzato in questo studio è il California Housing Dataset, un insieme di dati pubblici ampiamente utilizzato nella letteratura di machine learning. Questo dataset contiene informazioni statistiche aggregate sul mercato immobiliare californiano, dove ogni osservazione rappresenta un blocco residenziale. Le feature disponibili includono sia variabili demografiche che proprietà fisiche delle abitazioni, in particolare:

- **longitude**: longitudine del blocco;
- **latitude**: latitudine del blocco;
- **housing_median_age**: età mediana delle abitazioni nel blocco;
- **total_rooms**: numero totale di stanze nel blocco;
- **total_bedrooms**: numero totale di camere da letto nel blocco;
- **population**: popolazione nel blocco;
- **households**: numero di nuclei familiari nel blocco;
- **median_income**: reddito mediano delle famiglie nel blocco;
- **median_house_value**: valore mediano delle abitazioni nel blocco;

Durante la fase di preprocessing, i valori mancanti nella variabile *total_bedrooms* sono stati sostituiti con la mediana. Tutte le feature numeriche sono state normalizzate nell'intervallo $[0, 1]$. La variabile categorica *ocean_proximity* è stata esclusa dall'analisi.

4 Comparazioni

I parametri configurati per il confronto tra i tre algoritmi sono i seguenti:

- iterazioni massime = 50000;
- step-size = 0.01;
- l1-penalty = 1;
- tolerance = $1e-4$;
- agenti = 8 (ADMM distribuito);

La Tabella 1 riassume il confronto delle prestazioni ottenute.

Tabella 1: Comparazione algoritmi

	R2	time (s)	iterazioni
ISTA	0.5588	12.876	50000
ADMM	0.5839	0.0010	4
ADMM-Dist	0.5726	0.0376	170

È importante sottolineare che i valori riportati sono stati ottenuti effettuando una suddivisione casuale del dataset in training set e test set. Pertanto, esecuzioni successive potrebbero produrre risultati leggermente differenti.

A seguire viene illustrato l'andamento dei criteri di convergenza durante le iterazioni degli algoritmi, insieme a dei grafici che confrontano le previsioni con i valori reali.

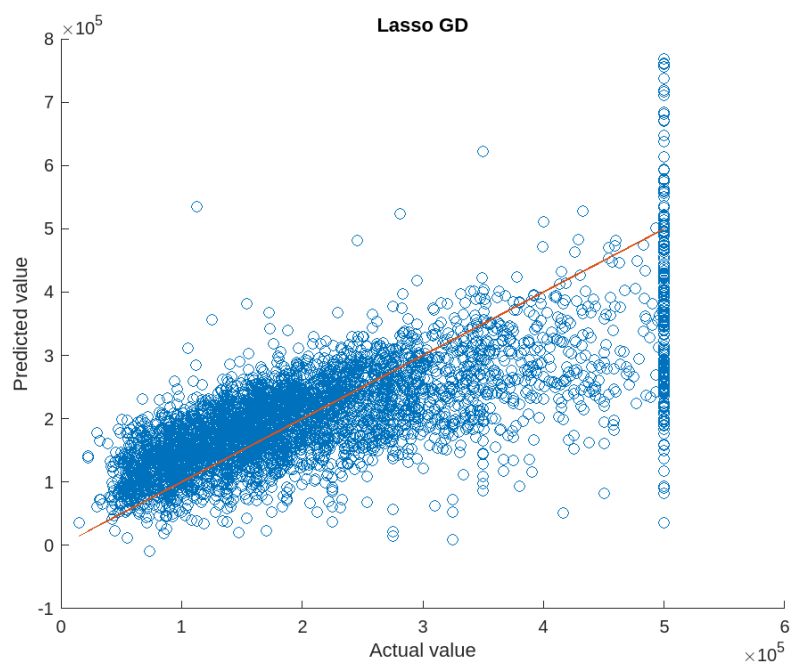


Figura 1: Predizioni ottenute con Soft-Thresholding.

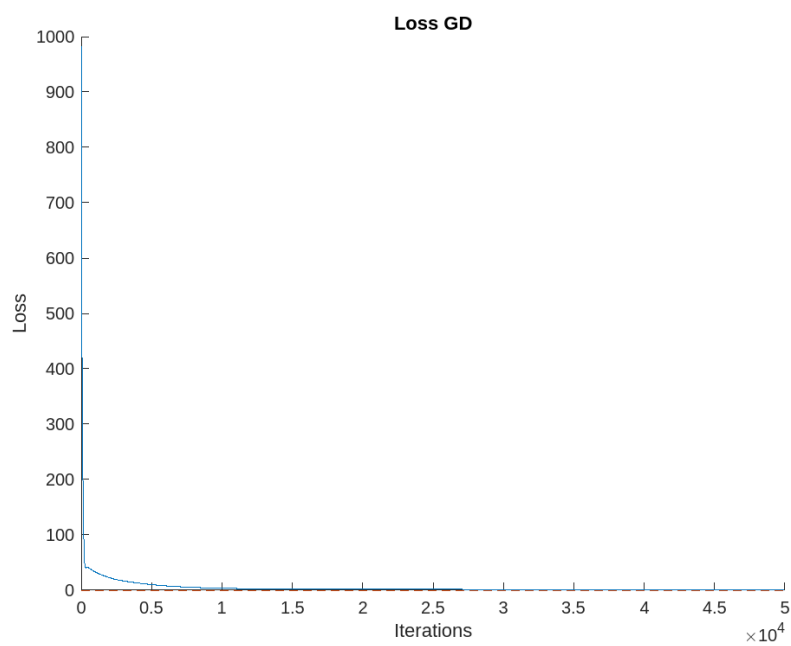


Figura 2: Convergenza del Soft-Thresholding.

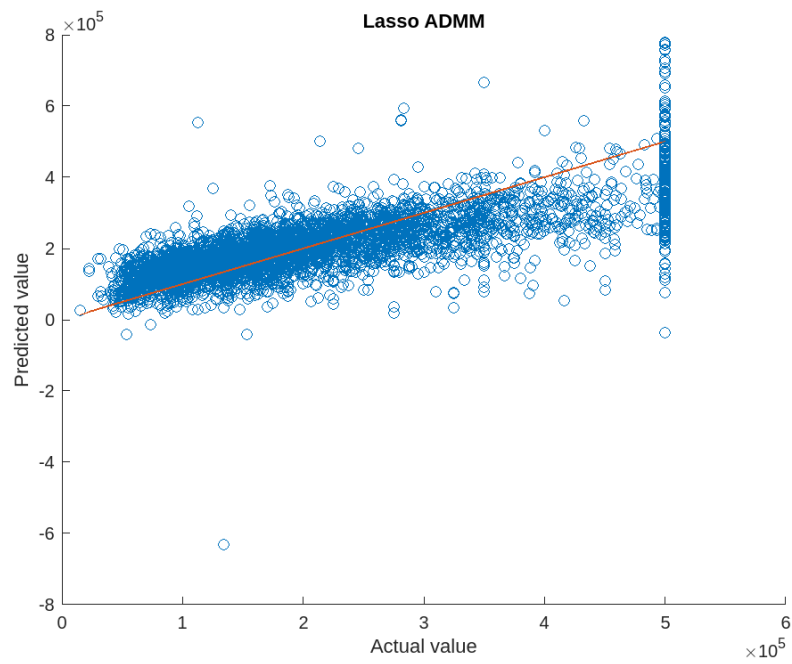


Figura 3: Predizioni ottenute con ADMM.

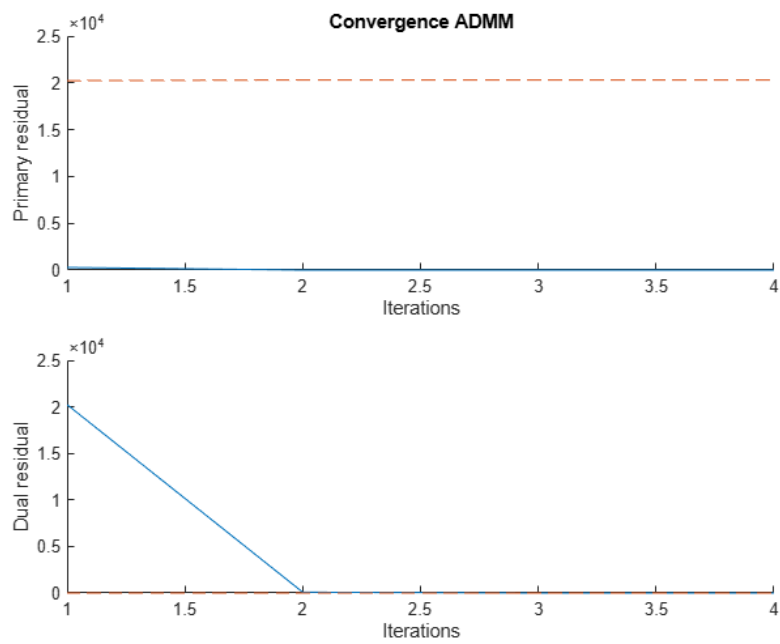


Figura 4: Convergenza di ADMM.

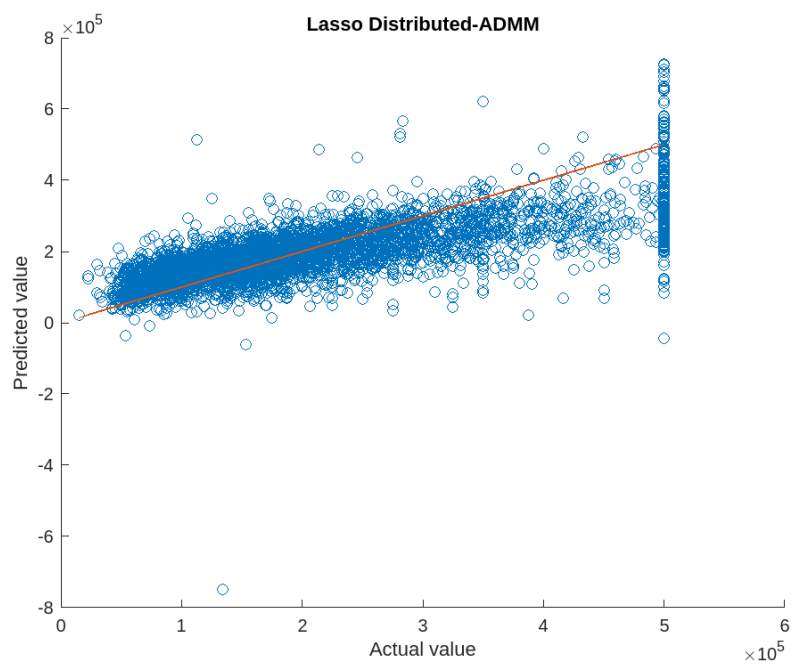


Figura 5: Predizioni ottenute con ADMM distribuito.

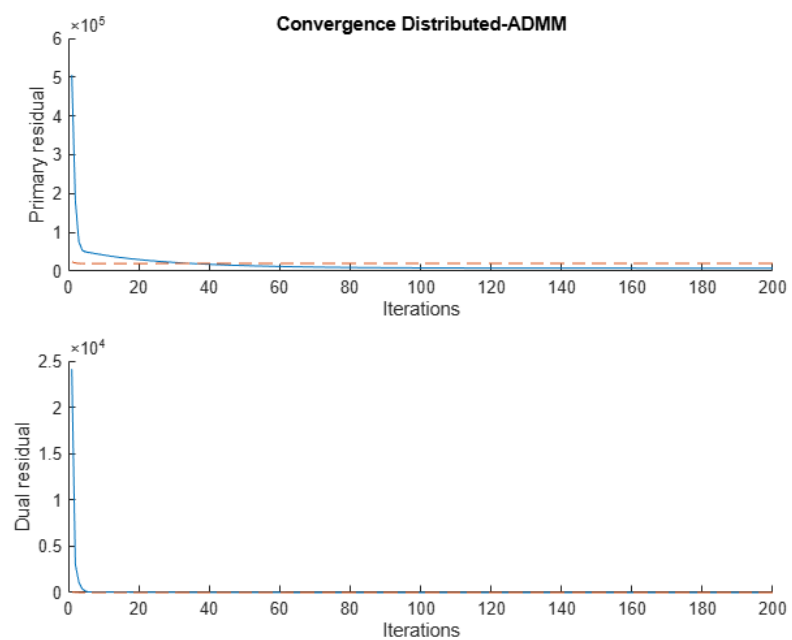


Figura 6: Convergenza di ADMM distribuito.