

# Voci dal braccio della morte: un'indagine tra linguaggio ed emozioni

Erica Solinas

erica.solinas@studio.unibo.it

## Abstract

Questo progetto si propone di creare un dataset completo che raccolga informazioni personali, riassunti dei reati e ultime dichiarazioni dei 591 condannati a morte giustiziati in Texas dal 1982 al 2024. Tramite tecniche di web scraping e OCR, i dati sono stati estratti, strutturati e analizzati con un focus sulle emozioni espresse nelle ultime parole dei condannati, la cui analisi è stata effettuata attraverso le API di OpenAI.

## 1 Introduzione

Il progetto è stato concepito con l'obiettivo di applicare le tecnologie studiate durante il corso, esplorando al contempo anche ulteriori strumenti e tecniche necessari alla sua realizzazione. Ciò ci ha permesso di sviluppare una comprensione generale delle diverse soluzioni disponibili, ampliando le competenze acquisite, seppur in modo sommario.

Abbiamo scelto di concentrarci sulla realizzazione di un progetto di natura linguistica e, fin dall'inizio, abbiamo mostrato un forte interesse per il sito del TDCJ (*Texas Department of Criminal Justice*). Abbiamo voluto ascoltare le loro voci, far emergere i loro sentimenti più profondi, leggere le ragioni dietro la loro condanna e scoprire se, nei loro ultimi istanti di vita, confessavano la loro colpa o cercavano di affermare la propria innocenza.

A tale scopo, abbiamo dovuto creare il dataset da cui partire attraverso il web scraping dell'intero sito<sup>1</sup>, al fine di creare un file CSV che riproducesse una tabella simile su cui poter effettuare tutte le opportune analisi. A questo punto abbiamo deciso di realizzare due tipi di analisi: una sui dati

demografici dei condannati e una lessicale sui riassunti delle loro cause di condanna e sulle loro ultime dichiarazioni. Quest'ultime sono state il focus centrale della fase successiva, in cui abbiamo usato le API di OpenAI per estrarne i sentimenti espressi e per riconoscere se in quelle parole l'imputato si dichiarava colpevole, innocente o non lo specificava. Nell'ultima fase abbiamo preso un campione del nostro dataset e l'abbiamo annotato manualmente per valutare le performance del modello.

## 2 Fase 1: Web Scraping

Per la fase di Web Scraping sono state utilizzate librerie fondamentali come `BeautifulSoup`, `requests` e `urllib`. L'obiettivo principale era la creazione di un dataset strutturato in un file CSV contenente tutte le informazioni sui detenuti presenti sul sito del TDCJ. Questa fase si è rivelata particolarmente lunga e complessa a causa di alcune problematiche legate alla struttura HTML del sito. In particolare, è stato necessario implementare una funzione specifica per gestire due link dei condannati, poiché la loro struttura differiva dal resto della pagina HTML. Inoltre, alcune informazioni chiave, come i riassunti delle cause di condanna ("*Summary of Incident*") e il livello di istruzione ("*Education Level*"), erano presenti nella sezione "Inmate Information" del sito, ma in tanti casi erano salvate come immagini in formato `.jpg` anziché testo. Per ovviare a questo problema, è stato utilizzato un OCR (`easyocr`) per estrarre il testo dalle immagini, affiancato dalle API di OpenAI per migliorare la qualità dei dati estratti. Tuttavia, l'OCR non ha sempre restituito risultati ottimali, generando alcune imprecisioni che hanno introdotto del rumore nel dataset finale. Nonostante queste difficoltà, i dati sono stati integrati nel file CSV in formato JSON, consentendo di procedere alle fasi successive dell'analisi.

<sup>1</sup>[https://www.tdcj.texas.gov/death\\_row/dr\\_executed\\_offenders.html](https://www.tdcj.texas.gov/death_row/dr_executed_offenders.html)

### 3 Fase 2: Riconoscimento delle emozioni e del *plea status*

In questa fase centrale del progetto, il file CSV precedentemente creato è stato arricchito con due nuove colonne: *Emotions*, contenente l'analisi delle emozioni espresse, e *Plea Status*, relativa alla dichiarazione di colpa o innocenza presente nelle ultime parole dei condannati (*last statements*). L'analisi è stata effettuata utilizzando le API di OpenAI, configurate con un prompt specifico per il modello `gpt-3.5-turbo`, che fornisce una risposta strutturata in formato JSON.

Per l'estrazione delle emozioni, il modello è stato istruito a identificare tre emozioni principali per ogni dichiarazione. Questa scelta è stata dettata dalla necessità di garantire coerenza e comparabilità tra i risultati, mantenendo tuttavia una certa libertà generativa per il modello. Infatti, abbiamo deliberatamente evitato di vincolarlo a un elenco predefinito di emozioni, così da consentirgli di esprimere il massimo delle sue capacità interpretative. Nel prompt, abbiamo fornito un elenco di esempi come punto di riferimento, senza imporre limitazioni rigide, se non il numero fisso di tre emozioni. Abbiamo anche sperimentato la possibilità di consentire al modello di identificare fino a un massimo di tre emozioni, ma questa configurazione ha prodotto risultati meno coerenti e soddisfacenti.

Per quanto riguarda il riconoscimento dello *Plea Status*, il modello è stato istruito a individuare se l'imputato si dichiarasse colpevole (*guilty*), innocente (*not guilty*) o se non facesse riferimento alla propria colpa (*unspecified*). Questo approccio ha permesso di estrarre informazioni significative dalle dichiarazioni, integrandole nel dataset per le analisi successive.

### 4 Fase 3: Analisi dei dati demografici e lessicali

Per condurre queste analisi è stato sviluppato un notebook dedicato, progettato per presentare i risultati in modo chiaro e intuitivo, anche grazie all'utilizzo di numerosi grafici. L'analisi si articola in due sezioni principali:

#### 4.1 Analisi dei dati demografici

Questa analisi fornisce una panoramica dettagliata delle caratteristiche anagrafiche e demografiche presenti nel dataset, rappresentando un primo passo fondamentale per comprendere il profilo dei

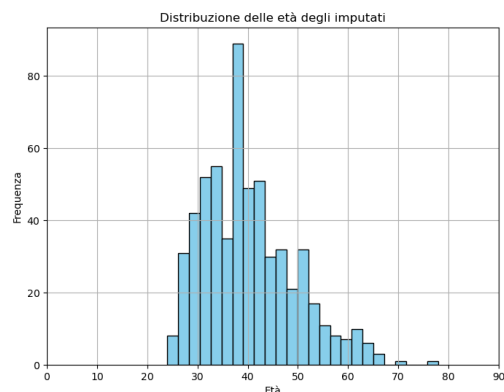


Figure 1: *Età media degli imputati al momento della morte*

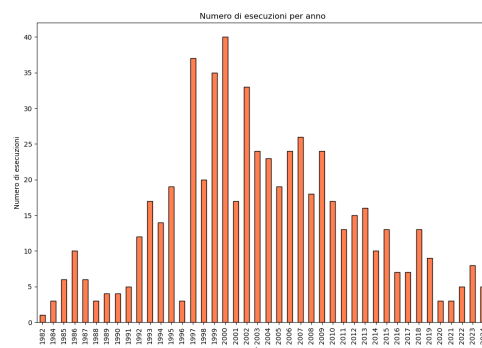


Figure 2: *N° di esecuzioni per anno*

condannati a morte inclusi nel nostro studio. I principali indicatori calcolati includono:

1. **Età media dei condannati a morte al momento dell'esecuzione**, pari a 40,2 anni. La distribuzione dell'età è stata rappresentata graficamente nella Figura 1, che evidenzia le tendenze relative all'età dei condannati.
2. **Anno con maggiori/minori esecuzioni**, insieme alla **media di esecuzioni in un anno**, fornisce una panoramica dell'andamento temporale delle esecuzioni dal 1982 al 2024. La media annuale è di 14 esecuzioni. Come illustrato nella Figura 2, il 1982 registra il numero minimo di esecuzioni, con una sola condanna portata a termine, mentre il 2000 si distingue come l'anno con il numero massimo, pari a 40 esecuzioni.
3. **Distribuzione di frequenza delle etnie tra i condannati a morte**: l'etnia più rappresentata nel dataset è "White", con 263

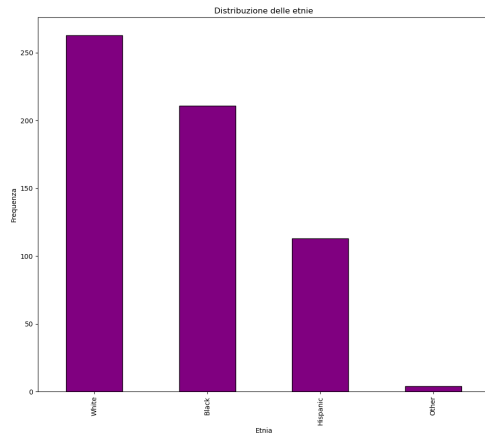


Figure 3: *Distribuzione di frequenza delle etnie*

casi, seguita da “Black” (211), “Hispanic” (113) e “Other” (4). La distribuzione è visualizzata nella Figura 3, che evidenzia chiaramente la predominanza delle prime tre categorie.

4. **Distribuzione di frequenza dei distretti di appartenenza dei condannati a morte**, analizzata per identificare la loro provenienza. La Figura 4 riporta i primi dieci distretti più rappresentati, con Harris al primo posto con 135 condannati, seguito da Dallas (65), Bexar (46), Tarrant (45) e Nueces (17). Numerosi altri distretti presentano un’unica occorrenza, confermando una concentrazione geografica in specifiche aree del Texas.

5. **Livello di istruzione medio** dei condannati a morte, calcolato sulla base dei dati disponibili. La colonna relativa al “Education Level” presenta alcune anomalie dovute a errori generati dall’utilizzo dell’OCR, come evidenziato nella Figura 5. Tuttavia, notiamo che il valore più ricorrente è 12, che corrisponde al “12th grade” nel sistema scolastico statunitense, equivalente al quinto anno delle scuole superiori italiane.

## 4.2 Analisi dei dati lessicali

Questa sezione analizza i dati testuali relativi alle ultime dichiarazioni dei condannati a morte (*Last Statement*) e al riassunto del reato commesso (*Summary of Incident*). Per entrambe le categorie, i testi sono stati sottoposti a un processo di normalizzazione e tokenizzazione, con l’obiettivo

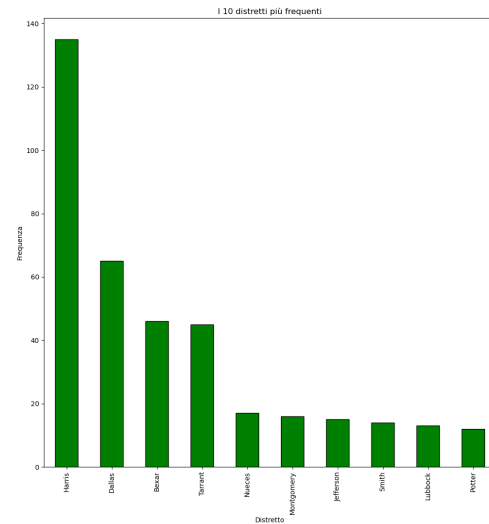


Figure 4: *I 10 distretti più frequenti*

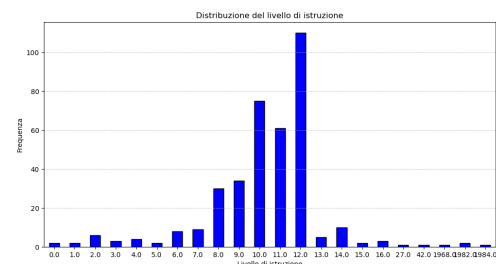


Figure 5: *Livello di istruzione medio*

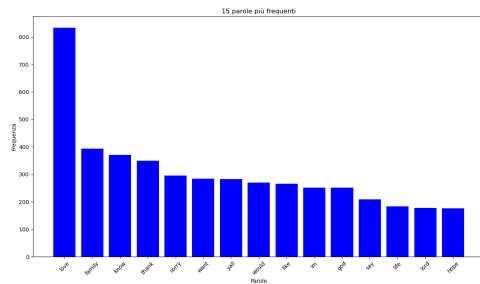


Figure 6: Le 15 parole più frequenti della sezione "Last Statement"

di calcolare i token e i bigrammi più frequenti. I risultati ottenuti sono stati rappresentati graficamente e memorizzati in file CSV, contenenti tutti i token e i bigrammi e le relative occorrenze ordinati in ordine decrescente di frequenza.

- **Sezione "Last Statement"**

Nell'analisi delle ultime dichiarazioni, abbiamo individuato frasi non rilevanti dal punto di vista lessicale, poiché si trattava di formule standardizzate che indicavano il rifiuto dell'imputato di rilasciare un'ultima dichiarazione. Queste frasi, se incluse, avrebbero potuto distorcere i risultati dell'analisi. Poiché nel dataset non era presente una formulazione unica per tali dichiarazioni, abbiamo effettuato una ricerca per identificare e raccogliere tutte le varianti utilizzate, escludendole successivamente dall'elaborazione dei dati. La Figura 6 mostra le 15 parole più utilizzate nelle ultime dichiarazioni dei condannati a morte. Tra queste, emergono termini fortemente emotivi come *love* (con ben 833 occorrenze si afferma come il più frequente in assoluto), *sorry* e *thank*, che riflettono sentimenti di affetto, rimorso e gratitudine. Parole come *family*, *lord* e *god* evidenziano un legame con i valori della famiglia e della religione, suggerendo che, nei momenti finali, i pensieri dei condannati si rivolgono principalmente a coloro che hanno rappresentato un sostegno emotivo o spirituale durante la loro vita. Questi dati indicano un quadro emotivo complesso, in cui prevalgono messaggi di amore e ringraziamento, ma anche di speranza, accompagnati da riflessioni sul rimorso e sulla fede.

La Figura 7 mostra i 10 bigrammi più frequenti nelle ultime dichiarazioni dei condan-

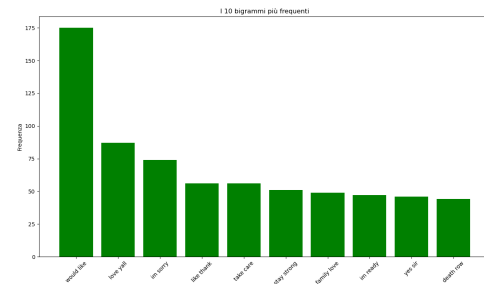


Figure 7: I 10 bigrammi più frequenti della sezione "Last Statement"

nati a morte. Il bigramma più ricorrente è *would like*, con 175 occorrenze, seguito da espressioni significative come *take care*, *stay strong*, *love family*, dunque messaggi di incoraggiamento rivolti ai familiari e agli amici, mostrando una preoccupazione per i propri cari anche negli ultimi istanti di vita. Tra le espressioni più ricorrenti troviamo anche *im sorry* e *im ready*, che suggeriscono rassegnazione per il destino imminente, ma anche un profondo rimorso per il dolore causato.

- **Sezione "Summary of Incident"** In questa sezione, abbiamo analizzato i riassunti dei reati commessi dai condannati, applicando le stesse tecniche utilizzate nella sezione precedente, ossia la normalizzazione, la tokenizzazione e il calcolo delle frequenze dei termini. La Figura 8 illustra le 15 parole più frequenti nei riassunti, evidenziando chiaramente la natura dei crimini descritti. Tra i termini più ricorrenti, spiccano quelli legati alla violenza, come *shot* (la parola più frequente con 342 occorrenze), *murder* e *robbery*, nonché quelli che fanno riferimento agli esiti tragici dei crimini, come *death* e *victim*. Inoltre, appaiono parole che indicano i contesti in cui i crimini sono stati commessi o scoperti, come *home*, *car* e *found*. Infine, termini come *police* e *convicted* sottolineano l'aspetto investigativo e giudiziario dei reati descritti.

Tra i bigrammi più ricorrenti mostrati nel grafico della Figura 9, *year old* appare come il più frequente (146 occorrenze), evidenziando l'età delle vittime o dei colpevoli. Un altro bigramma significativo è *capital murder*, che suggerisce l'entità dei crimini commessi, ovvero omicidi capitali, a loro

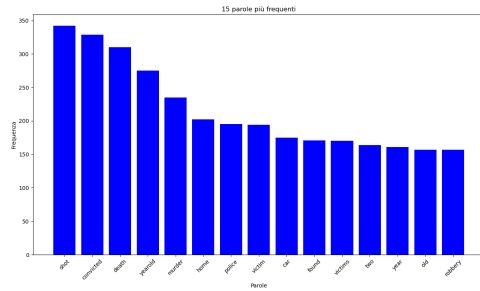


Figure 8: Le 15 parole più frequenti della sezione "Summary of Incident"

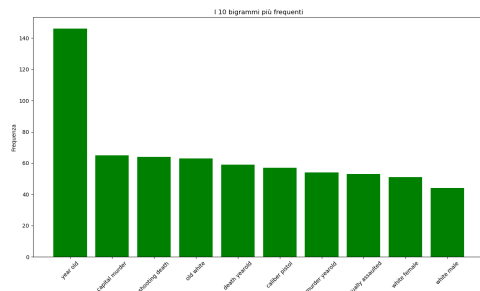


Figure 9: I 10 bigrammi più frequenti della sezione "Summary of Incident"

volta seguiti da termini come *shooting death*, *caliber pistol* e *sexually assaulted*, che si riferiscono a modalità e strumenti utilizzati durante l'atto. Questi bigrammi, nel loro complesso, rivelano dettagli chiave sui crimini commessi, sulle caratteristiche delle vittime e sulle modalità con cui sono stati perpetrati.

Tuttavia, ci teniamo a sottolineare che anche in questa sezione è stato utilizzato un OCR per estrarre i dati, il quale ha introdotto del rumore, con la conseguente possibilità che i risultati ottenuti non siano completamente accurati.

## 5 Fase 4: Valutazione finale

Per valutare in modo significativo i risultati ottenuti, abbiamo annotato manualmente un campione casuale del dataset, pari al 10% (59 frasi), con l'obiettivo di creare un Gold Standard per il confronto con le predizioni del modello. L'annotazione delle emozioni nel Gold Standard è stata effettuata assegnando fino a tre emozioni per ciascuna dichiarazione, selezionando quelle ritenute più rappresentative. Le prestazioni del modello sono state misurate utilizzando due met-

riche principali: la percentuale di corrispondenza tra emozioni predette e quelle annotate (*Emotion Match*) e il tasso di correttezza nel riconoscimento del *plea status* (*Correct Plea*). I risultati di queste analisi sono stati inclusi in una colonna aggiuntiva del file CSV. Dall'analisi emerge che il tasso di corrispondenza tra le emozioni predette e quelle del Gold Standard è pari al **68,64%**, mentre la percentuale di correttezza nel riconoscimento del *plea status* raggiunge l'**81,36%**. Tra le emozioni predette e quelle annotate nel Gold Standard, "love" risulta essere la più frequente in entrambi i casi. Tuttavia, si osservano alcune differenze nella distribuzione delle emozioni più comuni, come evidenziato nelle tabelle sottostanti.

Predicted Emotion	Frequency
love	28
hope	26
remorse	20
gratitude	20
acceptance	13
forgiveness	12
regret	11
silence	10
defiance	6
anger	5

Table 1: 10 emozioni più frequenti predette dal modello

Gold Emotion	Frequency
love	26
gratitude	17
hope	16
remorse	13
regret	8
forgiveness	7
anger	5
defiance	4
resignation	3
disappointment	3

Table 2: 10 emozioni più frequenti nel Gold Standard.

Per quanto riguarda il riconoscimento del *plea status*, i risultati sono abbastanza simili a quelli del Gold Standard: notiamo comunque una certa difficoltà nel riconoscere "guilty" e "not guilty", con un numero inferiore di predizioni rispetto al Gold Standard come mostra la tabella sottostante.

<b>Plea Status</b>	<b>Predicted</b>	<b>Gold Standard</b>
unspecified	51	40
guilty	6	15
not guilty	2	4

Table 3: Confronto tra il Plea Status Predetto e il Gold Standard

## 6 Conclusioni

Questo progetto ha rappresentato un'importante esperienza di apprendimento, sia dal punto di vista tecnico che umano. Ha dimostrato l'applicabilità di tecniche avanzate come il web scraping, l'OCR e l'analisi del linguaggio naturale per estrarre e analizzare informazioni da un ampio dataset riguardante i condannati a morte in Texas.

Nonostante alcune difficoltà, in particolare con l'OCR nell'estrazione del testo dalle immagini, il progetto ha portato alla creazione di un dataset completo e ben strutturato. La valutazione finale del modello, confrontata con un campione annotato manualmente, ha evidenziato risultati promettenti in termini di accuratezza nel riconoscimento delle emozioni e dello stato di colpa. Tuttavia, esistono margini di miglioramento, ad esempio mediante l'uso di modelli alternativi o l'ottimizzazione dei prompt per l'estrazione dei dati. Sebbene il percorso sia ancora lungo, siamo certi che questo rappresenta un buon punto di partenza.

## References

Beautiful Soup documentation, <https://www.crummy.com/software/BeautifulSoup/bs4/doc>

Requests documentation, <https://docs.python-requests.org/en/v2.0.0/>

Matplotlib documentation, <https://matplotlib.org/stable/users/index.html>

Numpy documentation, <https://numpy.org/doc/>

easyOCR documentation, <https://www.jaided.ai/easyocr/documentation/>

OpenAI API documentation, <https://platform.openai.com/docs/api-reference/authentication>

Pandas documentation, <https://pandas.pydata.org/docs/>