



GoalGuru

Studente:

Nome
Choaib Goumri

Matricola
0512118390

Repository Github:

<https://github.com/RedMulaCrackIn/GoalGuru.git>

Contents

Introduzione	4
Applicazioni nel mondo odierno	4
Potenzialità future	4
Obiettivo del progetto	5
Limiti del progetto	5
Prospettive di miglioramento	5
Caratteristiche degli elementi utilizzati nel progetto	6
Modello	6
Business Understanding	6
Obiettivi	6
P.E.A.S.	7
Requisiti	7
Data Understanding	7
Descrizione dei dati	7
Esplorazione dei dati	8
Statistiche	8
Informazioni	9
Valori duplicati	9
Data Preparation	10
Data cleaning	10
Rimozione colonne irrilevanti	10
Conversione del tipo dei dati	10
Estrazione di informazioni temporali	10
Pulizia e consolidamento	10
Calcolo dei punti e vincitori della stagione	10
Gestione valori nulli	11
Feature scaling	11
Creazione di rapporti numerici	11
Standardizzazione/Normalizzazione:	11
Feature selection	12
Rimozione colonne irrilevanti al modello	12
Data Modeling	12
Suddivisione del Dataset	12
Approccio	12
K-Nearest Neighbors (KNN)	13
Parametri Chiave del Modello	13
Funzionamento del Modello	13
Selezione del Modello Migliore	14
Limiti del KNN	14
Superamento dei limiti con una Rete Neurale	14
Creazione del Modello di Rete Neurale	15
Funzione di Attivazione	15
Funzione di Perdita	16
Ricerca degli Iperparametri	16
Evaluation	20
Valutazione sulla Base dei Dati di Validazione	20
Grafico della Perdita durante l'Addestramento e la Validazione	20

Matrice di Confusione	23
Matrice KNN	25
Matrice Grid Search	25
Matrice Random Search	25
Recap analisi delle matrici	26
Considerazione Rete Neurale	26
Riepilogo del Modello	26
Valutazione sul Test Set	26
Report di Classificazione	27
Conclusioni della valutazione	27
Deployment	28
Obiettivi del Deployment	28
Attività del Codice	28
Salvataggio dei Risultati	28
Esportazione della Matrice di Confusione	28
Preparazione del Modello per il Deployment	28
Artefatti per il Deployment	29
Considerazioni Chiave	29
Scelte Progettuali per il Deployment	29
Requisiti per un Deployment Operativo	29
Esecuzione del deployment	29
Configurazione dell'Ambiente	30
Creazione API con Flask	30
Deployment in Cloud	30
Containerizzazione con Docker	31
Conclusione	31

Introduzione

Negli ultimi anni, l'applicazione dell'intelligenza artificiale nel mondo dello sport ha subito un'evoluzione significativa, rivoluzionando il modo in cui analisti, appassionati e professionisti interpretano il gioco. L'analisi predittiva delle partite di calcio, un campo che fino a poco tempo fa era dominato da intuizioni soggettive e modelli statistici tradizionali, sta ora beneficiando dei progressi dell'apprendimento automatico.

Il progetto GoalGuru nasce con l'obiettivo di sfruttare tecniche di Machine Learning per fornire previsioni affidabili sui risultati delle partite di calcio, analizzando dati storici e individuando schemi nascosti che possano influenzare gli esiti degli incontri. In particolare, vengono applicati algoritmi di apprendimento supervisionato, come K-Nearest Neighbors (KNN), per modellare le relazioni tra le variabili e ottimizzare la capacità predittiva del sistema. Solo successivamente, nei casi in cui questi modelli mostrano delle limitazioni, viene applicata una rete neurale per perfezionare l'analisi e migliorare la precisione delle previsioni.

Questa metodologia consente di bilanciare semplicità ed efficienza computazionale: il KNN offre una soluzione intuitiva e interpretabile per classificare i risultati delle partite, mentre la rete neurale viene utilizzata per affinare le previsioni e gestire le situazioni più complesse.

Applicazioni nel mondo odierno

L'uso degli algoritmi di Machine Learning nella previsione dei risultati sportivi sta trovando un numero crescente di applicazioni pratiche. GoalGuru può essere sfruttato in diversi scenari del mondo reale, tra cui:

- **Settore delle scommesse sportive:** Gli utenti possono utilizzare il sistema per prendere decisioni più informate sulle scommesse, analizzando dati oggettivi e trend storici piuttosto che affidarsi esclusivamente all'intuito.
- **Analisi tecnica per squadre e allenatori:** I club calcistici possono integrare il modello per analizzare le prestazioni delle squadre avversarie e individuare punti di forza e debolezze basati sui dati storici.
- **Supporto per analisti sportivi e media:** Giornalisti e analisti possono utilizzare GoalGuru per arricchire le loro previsioni e offrire analisi basate su dati concreti ai loro lettori o spettatori.
- **Coinvolgimento dei tifosi:** Un'applicazione che fornisce previsioni accurate può aumentare l'interazione dei tifosi, creando dibattiti più informati e coinvolgenti.

L'utilizzo di modelli predittivi, inoltre, apre la strada a nuove possibilità nel mondo dello sport, come la previsione delle prestazioni individuali dei giocatori, il calcolo del rischio di infortuni e l'ottimizzazione delle strategie di gioco.

Potenzialità future

Sebbene il progetto GoalGuru abbia già dimostrato la sua efficacia nel prevedere i risultati delle partite utilizzando il KNN e le reti neurali come supporto, le sue potenzialità future sono vaste. Alcune possibili evoluzioni includono:

- **Espansione del dataset:** Attualmente, il sistema si basa su dati storici relativi alle competizioni più seguite. Tuttavia, un'integrazione con dati provenienti da campionati minori o tornei giovanili potrebbe ampliarne l'utilità.
- **Utilizzo di dati in tempo reale:** L'integrazione con fonti di dati live permetterebbe di aggiornare costantemente le previsioni in base agli sviluppi della partita in corso, migliorando la capacità adattativa del modello.
- **Ottimizzazione dei modelli di apprendimento supervisionato:** L'uso di tecniche avanzate di feature engineering potrebbe migliorare l'efficacia del KNN, riducendo i limiti legati alla sensibilità dei parametri di vicinanza.
- **Integrazione con strumenti di visualizzazione avanzata:** Dashboard interattive potrebbero offrire agli utenti rappresentazioni grafiche dettagliate delle analisi predittive.

Con il continuo progresso dell'intelligenza artificiale, GoalGuru ha il potenziale per diventare un sistema di riferimento nell'analisi predittiva sportiva.

Obiettivo del progetto

L'obiettivo principale di GoalGuru è sviluppare un sistema capace di prevedere i risultati delle partite di calcio in modo accurato ed efficiente, sfruttando metodi di apprendimento supervisionato. Il progetto si propone di:

- **Analizzare i dati storici delle partite** per individuare correlazioni tra variabili chiave, come il possesso palla, il numero di tiri e la formazione della squadra.
- **Applicare il KNN come primo modello predittivo**, sfruttando la similarità tra partite passate per generare previsioni sui nuovi match.
- **Individuare le limitazioni del modello supervisionato** e applicare una rete neurale per migliorare le previsioni nei casi più complessi.
- **Valutare costantemente la performance del sistema**, ottimizzando i parametri per aumentare la precisione e la generalizzazione del modello.

Questo approccio ibrido consente di sfruttare i punti di forza di entrambi gli algoritmi, migliorando l'efficacia complessiva della previsione.

Limiti del progetto

Nonostante i risultati promettenti ottenuti finora, il progetto GoalGuru presenta alcune limitazioni:

- **Dati parziali:** Il modello si basa su statistiche storiche e non tiene conto di eventi imprevedibili come infortuni o condizioni atmosferiche.
- **Limiti del KNN:** Il KNN, sebbene semplice ed efficace, può soffrire di problemi di scalabilità e sensibilità ai parametri.
- 1. **Overfitting della rete neurale:** L'utilizzo della rete neurale nei casi complessi richiede un'attenta gestione dell'overfitting, poiché il modello potrebbe adattarsi troppo ai dati di training.
- **Mancanza di adattamento in tempo reale:** Al momento, le previsioni vengono fatte esclusivamente sulla base dei dati passati e non possono aggiornarsi durante lo svolgimento di una partita.

Queste limitazioni rappresentano opportunità di miglioramento per le future iterazioni del progetto.

Prospettive di miglioramento

Per migliorare l'accuratezza e l'applicabilità di GoalGuru, si potrebbero introdurre diversi miglioramenti:

- **Ottimizzazione del KNN:** Studiando la scelta ottimale del parametro k e migliorando la normalizzazione dei dati per aumentare la robustezza del modello.
- **Espansione delle variabili analizzate:** Considerare nuovi fattori, come il morale della squadra o gli stili di gioco, per arricchire l'analisi.
- **Miglior gestione dei dati mancanti:** Implementare tecniche di imputation più avanzate per ridurre l'impatto dei dati assenti sulle previsioni.
- **Rete neurale più avanzata:** Se necessario, migliorare l'architettura della rete neurale per affinare ulteriormente la precisione del sistema.

L'obiettivo a lungo termine è rendere GoalGuru uno strumento sempre più affidabile e utile per analisti, tifosi e professionisti del calcio.

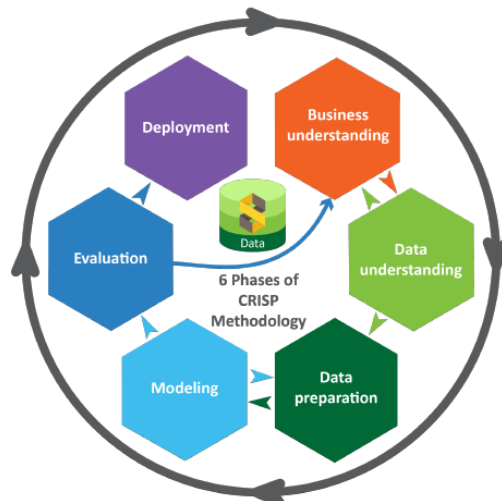
Caratteristiche degli elementi utilizzati nel progetto

Il progetto GoalGuru si basa su strumenti e tecnologie all'avanguardia:

- **Dataset:** Acquisito da fonti attendibili come Kaggle, contenente dati dettagliati sulle partite di calcio.
- **Linguaggi di programmazione:** Utilizzo di Python per lo sviluppo del modello, grazie alla sua flessibilità e alle sue potenti librerie.
- **Algoritmi:** K-Nearest Neighbors come base dell'apprendimento supervisionato, supportato da una rete neurale per miglioramenti mirati.
- **Librerie principali:** Pandas, NumPy, Scikit-learn per il KNN e TensorFlow per la rete neurale.

Modello

CRISP-DM è l'acronimo di Cross-Industry Standard Process for Data Mining, un metodo di comprovata efficacia per l'esecuzione di operazioni di data mining. Come metodologia, comprende descrizioni delle tipiche fasi di un progetto e delle attività incluse in ogni fase e fornisce una spiegazione delle relazioni esistenti tra tali attività.



Business Understanding

Obiettivi

Lo scopo del progetto consiste nella realizzazione di una rete neurale per prevedere i risultati delle partite di calcio. Questo progetto si basa su un approccio di machine learning e utilizza dati storici per fare previsioni accurate. Gli obiettivi specifici includono:

- Migliorare l'accuratezza delle previsioni dei risultati delle partite di calcio.
- Fornire agli utenti previsioni affidabili che possano utilizzare per informare le loro decisioni, come scommesse o analisi sportive.
- Aumentare l'engagement degli utenti incoraggiandoli a seguire più partite e a esplorare nuovi campi-onati e squadre.
- Fidelizzare gli utenti offrendo un servizio di previsioni di alta qualità che soddisfi le loro esigenze e le loro aspettative.

P.E.A.S.

PEAS è l'acronimo inglese di Performance Environment Actuators Sensors. È utilizzato nello studio dell'intelligenza artificiale per raggruppare in un unico termine l'ambiente operativo.

Performance: La misura di prestazione considerata prevede la precisione delle previsioni dei risultati delle partite di calcio, cercando di avvicinarsi il più possibile ai risultati reali.

Environment: L'ambiente in cui opera l'agente è il contesto delle partite di calcio, dove ha la possibilità di accedere ai dati storici e alle statistiche delle squadre e dei giocatori. Il nostro ambiente risulta:

- Completamente osservabile: l'agente ha accesso a tutte le informazioni necessarie riguardanti le partite, le squadre e i giocatori in ogni momento.
- Stocastico: lo stato dell'ambiente cambia indipendentemente dalle previsioni fatte dall'agente.
- Episodico: la predizione per una partita non influenza direttamente quella delle partite successive.
- Statico: L'ambiente (il dataset) non cambia mentre l'agente sta elaborando le sue decisioni.
- Discreto: le previsioni dipendono dalle caratteristiche discrete delle partite considerate.
- Ad Agente Singolo: vi sarà un unico agente ad operare.

Actuators: L'agente agisce sull'ambiente fornendo previsioni sui risultati delle partite di calcio.

Sensors: L'agente percepisce le informazioni attraverso l'accesso diretto ai dati storici delle partite e alle statistiche delle squadre e dei giocatori.

Requisiti

Il sistema proposto dovrà essere in grado di:

- Raccogliere e analizzare i dati relativi alle partite di calcio, inclusi i risultati storici, le statistiche delle squadre e dei giocatori.
- Utilizzare algoritmi di machine learning per generare previsioni accurate in tempo reale, basate sui dati raccolti.

Data Understanding

In questa seconda fase del CRISP-DM verrà analizzato il dataset in modo approfondito per comprendere la sua struttura, il contenuto, le relazioni tra le variabili e le eventuali problematiche o limitazioni dei dati.

Descrizione dei dati

Il dataset considerato presenta 4788 entry. Di ogni partita vengono riportate le seguenti caratteristiche:

- **date:** la data della partita.
- **time:** l'ora della partita.
- **comp:** la competizione della partita.
- **round:** il turno della partita.
- **day:** il giorno della settimana della partita.
- **venue:** il luogo della partita.
- **result:** il risultato della partita.
- **gf:** i gol segnati dalla squadra di casa.
- **ga:** i gol segnati dalla squadra ospite.
- **opponent:** l'avversario della squadra di casa.

- **xg**: gli expected goals (gol attesi) della squadra di casa.
- **xga**: gli expected goals (gol attesi) della squadra ospite.
- **poss**: il possesso palla della squadra di casa.
- **captain**: il capitano della squadra di casa.
- **formation**: la formazione della squadra di casa.
- **referee**: l'arbitro della partita.
- **sh**: i tiri della squadra di casa.
- **sot**: i tiri in porta della squadra di casa.
- **dist**: la distanza media dei tiri della squadra di casa.
- **fk**: i calci di punizione della squadra di casa.
- **pk**: i calci di rigore della squadra di casa.
- **pka**: i calci di rigore tentati dalla squadra di casa.
- **season**: l'anno della stagione della partita.
- **team**: la squadra di casa.

Esplorazione dei dati

In questa fase verranno condotte analisi esplorative per scoprire schemi, tendenze, correlazioni o anomalie nei dati delle partite di calcio. Tutte le analisi verranno condotte con l'ausilio di Python, utilizzando le librerie pandas e matplotlib per l'elaborazione e la visualizzazione dei dati.

Statistiche

Analizziamo le statistiche riguardanti il dataset.

```
1 print ("\nStatistiche descrittive:")
2 print (df.describe())
```

	Unnamed:0	gf	ga	xg	xga
count	4788.000000	4788.000000	4788.000000	4788.000000	4788.000000
mean	63.044069	1.447995	1.405388	1.396512	1.364745
std	42.865191	1.312635	1.286927	0.828847	0.814947
min	0.000000	0.000000	0.000000	0.000000	0.000000
25%	28.000000	0.000000	0.000000	0.800000	0.700000
50%	62.000000	1.000000	1.000000	1.300000	1.200000
75%	87.000000	2.000000	2.000000	1.900000	1.800000
max	182.000000	9.000000	9.000000	7.000000	7.000000

	poss	attendance	notes	sh	sot
count	4788.000000	3155.000000	0.0	4788.000000	4788.000000
mean	50.432957	38397.586688	NaN	12.619256	4.261278
std	12.810958	17595.849137	NaN	5.548444	2.459963
min	18.000000	2000.000000	NaN	0.000000	0.000000
25%	41.000000	25513.500000	NaN	9.000000	2.000000
50%	51.000000	36347.000000	NaN	12.000000	4.000000
75%	60.000000	53235.500000	NaN	16.000000	6.000000
max	82.000000	75546.000000	NaN	36.000000	15.000000

	dist	fk	pk	pkatt	season
count	4786.000000	4788.000000	4788.000000	4788.000000	4788.000000
mean	17.356247	0.453216	0.118212	0.14599	2022.365079
std	3.049341	0.665250	0.342362	0.37937	1.461850
min	5.300000	0.000000	0.000000	0.00000	2020.000000
25%	15.400000	0.000000	0.000000	0.00000	2021.000000
50%	17.200000	0.000000	0.000000	0.00000	2023.000000
75%	19.100000	1.000000	0.000000	0.00000	2024.000000
max	39.900000	4.000000	3.000000	3.00000	2024.000000

Informazioni

Ora passiamo alle informazioni del dataset e sui tipi.

```
1 print("\nInformazioni sul dataset:")
2 print(df.info())
```

#	Column	Non-Null Count	Dtype
0	Unnamed:0	4788 non-null	int64
1	date	4788 non-null	object
2	time	4788 non-null	object
3	comp	4788 non-null	object
4	round	4788 non-null	object
5	day	4788 non-null	object
6	venue	4788 non-null	object
7	result	4788 non-null	object
8	gf	4788 non-null	int64
9	ga	4788 non-null	int64
10	opponent	4788 non-null	object
11	xg	4788 non-null	float64
12	xga	4788 non-null	float64
13	poss	4788 non-null	int64
14	attendance	3155 non-null	float64
15	captain	4788 non-null	object
16	formation	4788 non-null	object
17	referee	4788 non-null	object
18	match report	4788 non-null	object
19	notes	0 non-null	float64
20	sh	4788 non-null	int64
21	sot	4788 non-null	int64
22	dist	4788 non-null	float64
23	fk	4788 non-null	int64
24	pk	4788 non-null	int64
25	pkatt	4788 non-null	int64
26	season	4788 non-null	int64
27	team	4788 non-null	object
dtypes:	float64(5),	int64(10),	object(13)

Notiamo che l'entry **notes** contiene solo valori **null**, mentre **attendance** presenta alcuni valori **null**.

Valori duplicati

Eseguendo

```
1 print("\nNumero di duplicati nel dataset:", df.duplicated().sum())
```

notiamo che non sono presenti duplicati nel dataset.

Data Preparation

Dopo aver acquisito e analizzato i dati, nella fase di Data Preparation questi vengono preparati affinché possano essere utilizzati in fase di addestramento dell'algoritmo di Machine Learning. In questa fase risulta utile l'uso della libreria sklearn, specialmente per il processo di Feature Engineering, come l'One-Hot Encoding e la standardizzazione delle variabili numeriche.

Data cleaning

Rimozione colonne irrilevanti

Rimosse colonne: **Unnamed: 0**, **comp**, **round**, **attendance**, **match report**, **notes**.

```
1 df.drop(columns=["Unnamed: 0", "comp", "round", "attendance", "match report", "notes"],  
         inplace=True)
```

Conversione del tipo dei dati

date: convertito in formato datetime.

venue, **opponent**, **team**, **result**: convertiti in tipo categoria.

```
1 df["date"] = pd.to_datetime(df["date"])  
2 df['venue'] = df['venue'].astype('category')  
3 df['opponent'] = df['opponent'].astype('category')  
4 df['team'] = df['team'].astype('category')  
5 df['result'] = df['result'].astype('category')
```

Estrazione di informazioni temporali

Aggiunta colonna **day** (nome del giorno) da **date**.

Creata colonna **hour** (ora del match) da **time**.

Creata colonna **day_code** (giorno della settimana come numero) da **date**.

```
1 df['day'] = df['date'].dt.day_name()  
2 df['hour'] = df['time'].str.replace(":", "", regex=True).astype("int")  
3 df['day_code'] = df['date'].dt.dayofweek
```

Pulizia e consolidamento

Rimosse stringhe extra da **formation** (es. "", "-0").

Consolidate formazioni poco frequenti in una categoria unica **Altro**.

```
1 df.formation = df.formation.str.replace(" ", "")  
2 df.formation = df.formation.str.replace("-0", "")  
3 value_counts = df.formation.value_counts()  
4 to_replace = value_counts[value_counts < 107].index  
5 df['formation'] = df['formation'].replace(to_replace, 'Altro')
```

Calcolo dei punti e vincitori della stagione

Assegnazione dei punti in base al risultato (W = 3, D = 1, L = 0).

Calcolo dei vincitori di ogni stagione e aggiunta della colonna **season_winner**.

```

1 df['points'] = df['result'].apply(lambda x: 3 if x == 'W' else 1 if x == 'D' else 0)
2 df['points'] = df['points'].astype('int')
3 winners = df.groupby(['season', 'team'], observed=False)['points'].sum().reset_index() \
4     .sort_values(['season', 'points'], ascending=[True, False]) \
5     .groupby('season', observed=False).first()
6 df['season_winner'] = df['season'].map(winners['team'])

```

Gestione valori nulli

Gestione dei dati mancanti in **captain** con sostituzione o rimozione.

```

1 def captains_func(data):
2     if data['count'] == 0:
3         data['count'] = np.nan
4     return data
5 group = df.groupby('team', observed=False)['captain'].value_counts().reset_index(name='count')
6 group = group.apply(captains_func, axis=1)
7 group.dropna(inplace=True)
8 group = group.drop(columns='count')

```

Feature scaling

Creazione di rapporti numerici

fk_ratio: rapporto tra calci di punizione (**fk**) e tiri (**sh**).

pk_conversion_rate: rapporto tra rigori segnati (**pk**) e tentati (**pkatt**).

pk_per_shot: rapporto tra rigori tentati (**pkatt**) e tiri (**sh**).

```

1 def calculate_fk_pk_ratios(data):
2     data['fk_ratio'] = data['fk'] / data['sh']
3     data['pk_conversion_rate'] = data['pk'] / data['pkatt']
4     data['pk_per_shot'] = data['pkatt'] / data['sh']
5     # Gestione dei valori infiniti
6     data['fk_ratio'] = data['fk_ratio'].replace([np.inf, -np.inf], np.nan)
7     data['pk_conversion_rate'] = data['pk_conversion_rate'].replace([np.inf, -np.inf], np.nan)
8     data['pk_per_shot'] = data['pk_per_shot'].replace([np.inf, -np.inf], np.nan)
9     # Conversione in percentuali
10    data['fk_percentage'] = data['fk_ratio'] * 100
11    data['pk_conversion_percentage'] = data['pk_conversion_rate'] * 100
12    data['pk_per_shot_percentage'] = data['pk_per_shot'] * 100
13    return data
14 # Applicazione della funzione
15 df_sorted = calculate_fk_pk_ratios(df_sorted)
16 # Rimozione delle colonne non necessarie
17 df_sorted.drop(['pk_conversion_rate', 'pk_conversion_percentage'], axis=1, inplace=True)

```

Standardizzazione/Normalizzazione:

Utilizzo delle medie mobili per calcolare valori scalati nel tempo (**rolling_xg**, **rolling_xga**, etc.).

Applicazione One-Hot Encoding alle variabili categoriche e StandardScaler per le variabili numeriche.

```

1 def calculate_rolling_average(data, column, window=5):
2     return data.groupby('team', observed=False)[column].transform(
3         lambda x: x.rolling(window=window, min_periods=1).mean()
4     )
5 # Applicazione della funzione alle colonne selezionate
6 df_sorted['rolling_xg'] = calculate_rolling_average(df_sorted, 'xg')
7 df_sorted['rolling_xga'] = calculate_rolling_average(df_sorted, 'xga')
8 df_sorted['rolling_poss'] = calculate_rolling_average(df_sorted, 'poss')
9 df_sorted['rolling_sh'] = calculate_rolling_average(df_sorted, 'sh')
10 df_sorted['rolling_sot'] = calculate_rolling_average(df_sorted, 'sot')
11 df_sorted['rolling_dist'] = calculate_rolling_average(df_sorted, 'dist')
12 # Codifica del risultato in valori numerici (W = 1, D = 0, L = -1)

```

```

13 df_sorted['result_encoded'] = pd.to_numeric(df_sorted['result'].map({'W': 1, 'D': 0, 'L': -1}))
14 # Calcolo della forma (media mobile dei risultati)
15 df_sorted['form'] = calculate_rolling_average(df_sorted, 'result_encoded')
16
17 # One-Hot Encoding per le variabili categoriche
18 X_categorical_encoded = pd.get_dummies(X[categorical_cols], columns=categorical_cols)
19 # Standardizzazione delle variabili numeriche
20 scaler = StandardScaler()
21 X_numerical_scaled = pd.DataFrame(scaler.fit_transform(X[numerical_cols]), columns=
    numerical_cols)

```

Feature selection

Rimozione colonne irrilevanti al modello

Eliminati campi che non aggiungono valore predittivo diretto: **gf, ga, xg, xga, poss, points**, etc.

```

1 columns_to_drop = ['gf', 'ga', 'xg', 'xga', 'poss', 'sh', 'sot',
2                     'goal_diff', 'day', 'pk', 'pkatt', 'fk',
3                     'referee', 'dist', 'points', 'season_winner', 'hour', 'result_encoded'
4                     , 'day_code']
df_sorted = df_sorted.drop(columns=columns_to_drop)

```

Data Modeling

Nella fase di Data Modeling verranno determinate e valutate le tecniche finalizzate alla costruzione del modello, dopodiché ed si passerà alla fase di addestramento di quest'ultimo, durante la quale verranno configurati i parametri, verrà addestrato e si commenteranno i risultati ottenuti. Il problema affrontato è di apprendimento non supervisionato.

Suddivisione del Dataset

In questa fase, il dataset viene preparato separando le variabili predittive (X) dall'etichetta di target (y), e poi suddividendo i dati in un set di **training** e uno di **validation**. Questo processo è essenziale per addestrare e valutare il modello in modo efficace, evitando il rischio di overfitting e garantendo che il modello generalizzi bene sui dati non visti.

```

1 X = tr.drop('result', axis=1)
2 y = tr['result']
3 # Unione delle v# Codifica del target (y) in valori numerici
4 label_encoder = LabelEncoder()
5 y_encoded = label_encoder.fit_transform(y)
6 # Creazione del test set (ultima stagione)
7 test_x = X_final.tail(761)
8 test_y = y_encoded[-761:]
9 test_x['result'] = test_y
10 test_x.to_csv("test_set.csv", index=False)
11 # Rimozione dell'ultima stagione dal training set
12 X_final = X_final.iloc[:-761]
13 y_encoded = y_encoded[:-761]
14 # Divisione in training e validation set
15 X_train, X_val, y_train, y_val = train_test_split(X_final, y_encoded, test_size=0.2,
    random_state=42)

```

Approccio

Abbiamo affrontato il problema della modellazione dei dati adottando inizialmente il K-Nearest Neighbors (KNN), un algoritmo intuitivo che si basa sulla similarità tra punti dati. La scelta di KNN è stata motivata dalla sua semplicità e interpretabilità, caratteristiche che lo rendono un buon punto di partenza per comprendere il problema e sviluppare un primo modello predittivo. Inoltre, il suo approccio non parametrico lo rende adatto a situazioni in cui la distribuzione dei dati non è nota a priori. Tuttavia, nel corso dell'analisi, abbiamo

riscontrato alcune limitazioni significative, tra cui la sensibilità alla scelta del parametro k , la scalabilità computazionale e la difficoltà nel gestire dati complessi e ad alta dimensionalità. Per affrontare questi problemi, abbiamo esplorato metodi alternativi che permettessero una modellazione più efficace e una maggiore capacità di generalizzazione. Il confronto tra le diverse soluzioni ci ha permesso di individuare i punti di forza e di debolezza di ciascun approccio, guidandoci verso una strategia più robusta ed efficiente.

K-Nearest Neighbors (KNN)

Il KNN è un algoritmo di apprendimento supervisionato semplice ma potente, utilizzato principalmente per problemi di classificazione e regressione. Il principio fondamentale del KNN è che oggetti simili si trovano vicini nello spazio delle caratteristiche.

Input del Modello

Il modello riceve in input una serie di variabili numeriche e categoriali, rappresentate come vettori nello spazio multidimensionale. Questi vettori definiscono la posizione di ciascun punto (osservazione) nel dataset.

Parametri Chiave del Modello

Numero di Vicini (k): Indica quanti punti più vicini al dato di input devono essere considerati per effettuare una classificazione. Un valore basso di k può rendere il modello troppo sensibile ai rumori, mentre un valore alto può rendere il modello troppo generico.

Peso dei Vicini (weights): Determina l'importanza dei vicini nella decisione finale. Può essere uniforme (tutti i vicini hanno lo stesso peso) o basato sulla distanza (i vicini più vicini hanno un peso maggiore).

Metrica di Distanza (metric): Specifica il metodo per calcolare la distanza tra i punti, ad esempio la distanza euclidea o manhattan.

```
1 GRID_SEARCH = {
2     "n_neighbors": [3, 5, 7, 9, 11], # Numero di vicini
3     "weights": ['uniform', 'distance'], # Peso dei vicini
4     "metric": ['euclidean', 'manhattan'] # Metrica di distanza
5 }
6
7 # Creazione di tutte le combinazioni di iperparametri
8 grid_combinations = list(itertools.product(
9     GRID_SEARCH['n_neighbors'],
10    GRID_SEARCH['weights'],
11    GRID_SEARCH['metric']
12 ))
```

Funzionamento del Modello

Durante la fase di predizione, il modello cerca i k punti più vicini all'osservazione di input basandosi sulla metrica di distanza scelta.

Il modello assegna la classe predominante tra i k vicini per classificare l'osservazione, oppure calcola una media ponderata per i problemi di regressione.

```
1 for combination in grid_combinations:
2     n_neighbors, weights, metric = combination
3
4     print(f"Testing combination: n_neighbors={n_neighbors}, weights={weights}, metric={
5         metric}")
6
7     # Creazione del modello KNN con i parametri correnti
8     model = KNeighborsClassifier(n_neighbors=n_neighbors, weights=weights, metric=metric)
9
10    # Addestramento del modello
11    model.fit(X_train, y_train)
12
13    # Valutazione dell'accuratezza sul validation set
```

```

13 val_accuracy = model.score(X_val, y_val)
14
15 # Aggiungi l'accuratezza alla lista
16 val_accuracies.append(val_accuracy)
17
18 print(f"Validation accuracy: {val_accuracy}")
19
20 # Aggiornamento dei migliori iperparametri
21 if val_accuracy > best_val_accuracy:
22     best_val_accuracy = val_accuracy
23     best_params = {
24         "n_neighbors": n_neighbors,
25         "weights": weights,
26         "metric": metric
27     }
28
29 # Stampa dei migliori iperparametri
30 print("Best hyperparameters found:")
31 print(best_params)
32 print(f"Best validation accuracy: {best_val_accuracy}")

```

Iperparametri Testati

Il processo di ottimizzazione degli iperparametri prevede la selezione delle migliori combinazioni tra:

- Numero di vicini (k): Valori testati tra 3 e 11.
- Peso dei vicini: Uniforme o basato sulla distanza.
- Metrica di distanza: Euclidea o manhattan.

Selezione del Modello Migliore

Tutte le combinazioni di iperparametri vengono testate utilizzando un validation set. L'accuratezza sul validation set è la metrica utilizzata per valutare ciascun modello. La combinazione di iperparametri con la migliore accuratezza viene selezionata per creare il modello finale.

Limiti del KNN

Nonostante la semplicità e l'efficacia in problemi di piccole dimensioni, il KNN presenta alcuni limiti:

1. Sensibilità alla scelta di k: Un valore di k troppo piccolo rende il modello sensibile al rumore, mentre un valore troppo grande riduce la capacità di catturare le variazioni locali.
1. Scalabilità: Il KNN diventa inefficiente con dataset di grandi dimensioni, poiché ogni nuova predizione richiede il calcolo della distanza con tutti i punti di training.
1. Difficoltà con dati complessi: In problemi con distribuzioni più intricate, il KNN può avere difficoltà a generalizzare

Superamento dei limiti con una Rete Neurale

Abbiamo eseguito un'analisi dettagliata delle **performance dei modelli** nella fase di valutazione, concentrandoci su **metriche chiave come accuratezza, tempo di esecuzione e capacità di generalizzazione**. Dai risultati ottenuti, abbiamo osservato che il KNN, pur essendo efficace in scenari con dataset di piccole dimensioni e distribuzioni relativamente semplici, **non riusciva** a catturare adeguatamente **pattern complessi nei dati**, risultando meno performante rispetto ad approcci più avanzati. Questo ci ha portato a considerare metodi in grado di superare queste limitazioni, valutando le loro prestazioni su diversi set di dati e individuando una soluzione che garantisse un migliore equilibrio tra accuratezza, efficienza computazionale e capacità di adattamento alle caratteristiche del problema.

Creazione del Modello di Rete Neurale

La rete neurale è composta da **due strati densi (dense layers)** e da un **layer di dropout**. Ogni strato denso è completamente connesso, il che significa che ogni neurone in uno strato è connesso a tutti i neuroni dello strato successivo.

- **Strato di input:** L'input del modello è una serie di variabili numeriche e categoriche. La dimensione dell'input è determinata dal numero di variabili presenti nel dataset.
- **Strati nascosti (Hidden Layers):**
 - **Primo strato nascosto:** Il primo strato è una **layer denso** che contiene **neurons_1layer** neuroni, dove il numero di neuroni viene scelto tramite iperparametri.
 - **Secondo strato nascosto:** Il secondo strato è anch'esso un **layer denso** che contiene **neurons_2layer** neuroni. Anche in questo caso, il numero di neuroni è selezionato tramite iperparametri.
- **Dropout:** Un layer di **dropout** con una probabilità di 0.1 viene applicato dopo il secondo strato nascosto. Questo serve a **prevenire l'overfitting**, in quanto disattiva casualmente il 10% dei neuroni durante ogni aggiornamento del peso, forzando la rete a generalizzare meglio.
- **Strato di output:** Lo strato finale è un **layer denso** con 3 neuroni, uno per ciascuna delle classi: **vittoria, pareggio e sconfitta**. L'output è generato usando la funzione di attivazione **softmax**, che restituisce le probabilità per ciascuna classe, dove la somma delle probabilità è pari a 1.

```
1 def create_network(input_dim, neurons_1layer, neurons_2layer, activation_function):  
2     inputs = tf.keras.Input((input_dim,))  
3     x = layers.Dense(neurons_1layer, activation_function)(inputs)  
4     x = layers.Dense(neurons_2layer, activation_function)(x)  
5     x = layers.Dropout(0.1)(x)  
6     output = layers.Dense(3, "softmax")(x)  
7     model = tf.keras.Model(inputs=inputs, outputs=output, name="neural_net")  
8     return model
```

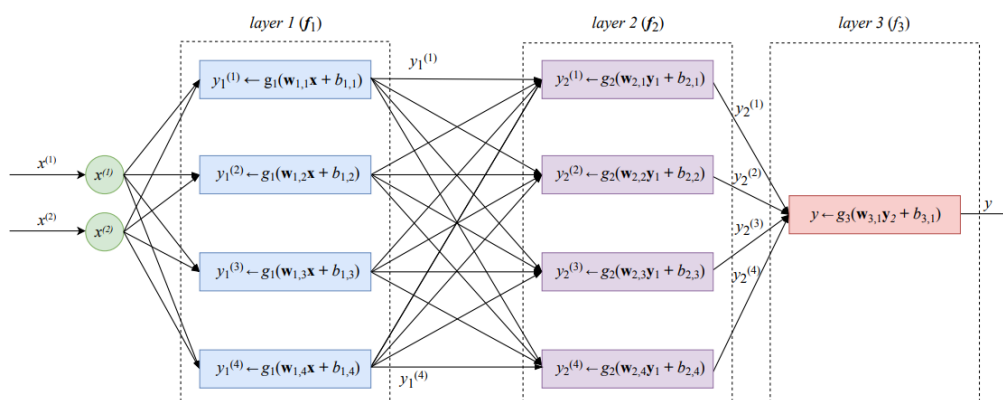


Immagine a scopo illustrativo

Funzione di Attivazione

La rete utilizza una **funzione di attivazione** non lineare, che è una componente chiave per introdurre flessibilità nel modello. Le funzioni di attivazione disponibili nel codice sono:

- **ReLU (relu):** Funzione di attivazione comunemente utilizzata nelle reti neurali per la sua capacità di evitare il problema del gradiente che scompare e di introdurre non linearità.

- **Sigmoid (sigmoid)**: Funzione di attivazione che restituisce un valore tra 0 e 1, utile quando l'output deve essere interpretato come una probabilità.
- **Tanh (tanh)**: Funzione di attivazione simile al Sigmoid, ma con un output che varia tra -1 e 1.

Funzione di Perdita

La **funzione di perdita** utilizzata è **sparse_categorical_crossentropy**, che è adatta per la classificazione multiclasse. In questo caso, le tre possibili etichette (vittoria, pareggio, sconfitta) sono trattate come classi distinte, e la funzione di perdita misura l'entropia incrociata tra le probabilità predette e quelle reali.

```
1 model.compile(
2     loss='sparse_categorical_crossentropy',
3     optimizer=tf.keras.optimizers.Adam(learning_rate=learning_rate),
4     metrics=['accuracy']
5 )
```

Ricerca degli Iperparametri

Grid Search La ricerca degli iperparametri è un processo fondamentale per ottimizzare le prestazioni di un modello. In questo caso, viene utilizzato il **grid search** per esplorare diverse combinazioni di iperparametri, come il numero di neuroni negli strati nascosti, la funzione di attivazione, il learning rate, il numero di epoche, e la dimensione del batch. Ogni combinazione viene testata e il modello con la **minore perdita di validazione** viene selezionato.

```
1 # Griglia di iperparametri
2 GRID_SEARCH = {
3     "learning_rate": [1e-3],
4     "epochs": [5, 6, 7, 8, 9, 10],
5     "neurons_1layer": [50, 55],
6     "neurons_2layer": [30, 50],
7     "activation_functions": ['relu', 'sigmoid', 'tanh'],
8     "batch_size": [200]
9 }
10
11 # Creazione di tutte le combinazioni di iperparametri
12 grid_combinations = list(itertools.product(
13     GRID_SEARCH['learning_rate'],
14     GRID_SEARCH['epochs'],
15     GRID_SEARCH['neurons_1layer'],
16     GRID_SEARCH['neurons_2layer'],
17     GRID_SEARCH['activation_functions'],
18     GRID_SEARCH['batch_size']
19 ))
20
21 # Variabili per tenere traccia dei migliori iperparametri
22 best_params = None
23 best_val_loss = np.inf
24
25 # Ciclo per testare ogni combinazione di iperparametri
26 for combination in grid_combinations:
27     learning_rate, epochs, neurons_1layer, neurons_2layer, activation_function,
28     batch_size = combination
29
30     print(f"Testing combination: lr={learning_rate}, epochs={epochs}, neurons_1layer={
31     neurons_1layer}, neurons_2layer={neurons_2layer}, activation={activation_function},
32     batch_size={batch_size}")
33
34     # Creazione del modello con i parametri correnti
35     model = create_network(X_train.shape[1], neurons_1layer, neurons_2layer,
36     activation_function)
37
38     # Compilazione del modello
39     model.compile(
40         loss='sparse_categorical_crossentropy',
41         optimizer=tf.keras.optimizers.Adam(learning_rate=learning_rate),
42         metrics=['accuracy']
43     )
```



```

40
41 # Addestramento del modello
42 history = model.fit(
43     X_train, y_train,
44     validation_data=(X_val, y_val),
45     epochs=epochs,
46     batch_size=batch_size,
47     verbose=0 # Nessun output durante la ricerca
48 )
49
50 # Valutazione della loss sul validation set
51 final_val_loss = history.history['val_loss'][-1]
52
53 print(f"Validation loss: {final_val_loss}")
54
55 # Aggiornamento dei migliori iperparametri
56 if final_val_loss < best_val_loss:
57     best_val_loss = final_val_loss
58     best_params = {
59         "learning_rate": learning_rate,
60         "epochs": epochs,
61         "neurons_1layer": neurons_1layer,
62         "neurons_2layer": neurons_2layer,
63         "activation_function": activation_function,
64         "batch_size": batch_size
65     }
66
67 # Stampa dei migliori iperparametri
68 print("Best hyperparameters found:")
69 print(best_params)
70 print(f"Best validation loss: {best_val_loss}")

```

Iperparametri Testati

Durante il **grid search**, vengono esplorati i seguenti iperparametri:

- **learning_rate**: La velocità con cui l'ottimizzatore aggiorna i pesi. Un valore troppo alto può causare instabilità, mentre uno troppo basso può rallentare l'apprendimento.
- **epochs**: Il numero di epoche rappresenta il numero di volte in cui il modello vedrà l'intero dataset durante l'addestramento. Maggiore è il numero di epoche, maggiore è la possibilità di apprendimento, ma rischia anche di causare overfitting.
- **neurons_1layer** e **neurons_2layer**: Il numero di neuroni nel primo e secondo strato della rete. Un numero maggiore di neuroni permette una maggiore capacità di apprendimento, ma può anche portare a un modello più complesso e incline all'overfitting.
- **activation_functions**: Le funzioni di attivazione da testare. Ogni funzione può influire sul modo in cui la rete apprende e sulle prestazioni finali.
- **batch_size**: La dimensione del batch, che determina quante istanze di dati vengono utilizzate per ogni aggiornamento dei pesi. Un valore maggiore di batch size riduce il rumore ma può richiedere più risorse computazionali.

Selezione del Modello Migliore

Alla fine del grid search, il modello con il **minor valore di perdita di validazione** (validation loss) viene selezionato. I parametri di questo modello migliore vengono poi utilizzati per allenare il modello finale.

Random Search Ora viene utilizzato il **random search** per esplorare combinazioni casuali di iperparametri, come il numero di neuroni negli strati nascosti, la funzione di attivazione, il learning rate, il numero di epoche, e la dimensione del batch. A differenza del grid search, il random search permette di testare un numero ridotto di combinazioni, riducendo il tempo necessario per la ricerca senza compromettere significativamente le possibilità di trovare iperparametri ottimali.

```
1 # Definizione dello spazio degli iperparametri per Random Search
2 RANDOM_SEARCH_SPACE = {
3     "learning_rate": [1e-3, 1e-4, 1e-5],
4     "epochs": [5, 10, 15, 20],
5     "neurons_1layer": [30, 50, 70],
6     "neurons_2layer": [20, 40, 60],
7     "activation_functions": ['relu', 'sigmoid', 'tanh'],
8     "batch_size": [100, 200, 300]
9 }
10
11 # Numero di combinazioni casuali da provare
12 NUM_RANDOM_COMBINATIONS = 20
13
14 # Variabili per tenere traccia dei migliori iperparametri
15 best_params = None
16 best_val_loss = np.inf
17
18 # Ciclo per testare combinazioni casuali di iperparametri
19 for _ in range(NUM_RANDOM_COMBINATIONS):
20     # Selezione casuale degli iperparametri
21     learning_rate = random.choice(RANDOM_SEARCH_SPACE["learning_rate"])
22     epochs = random.choice(RANDOM_SEARCH_SPACE["epochs"])
23     neurons_1layer = random.choice(RANDOM_SEARCH_SPACE["neurons_1layer"])
24     neurons_2layer = random.choice(RANDOM_SEARCH_SPACE["neurons_2layer"])
25     activation_function = random.choice(RANDOM_SEARCH_SPACE["activation_functions"])
26     batch_size = random.choice(RANDOM_SEARCH_SPACE["batch_size"])
27
28     print(f"Testing combination: lr={learning_rate}, epochs={epochs}, neurons_1layer={
29         neurons_1layer}, neurons_2layer={neurons_2layer}, activation={activation_function},
30         batch_size={batch_size}")
31
32     # Creazione del modello con i parametri correnti
33     model = create_network(X_train.shape[1], neurons_1layer, neurons_2layer,
34         activation_function)
35
36     # Compilazione del modello
37     model.compile(
38         loss='sparse_categorical_crossentropy',
39         optimizer=tf.keras.optimizers.Adam(learning_rate=learning_rate),
40         metrics=['accuracy']
41     )
42
43     # Addestramento del modello
44     history = model.fit(
45         X_train, y_train,
46         validation_data=(X_val, y_val),
47         epochs=epochs,
48         batch_size=batch_size,
49         verbose=0 # Nessun output durante la ricerca
50     )
51
52     # Valutazione della loss sul validation set
53     final_val_loss = history.history['val_loss'][-1]
54
55     print(f"Validation loss: {final_val_loss}")
56
57     # Aggiornamento dei migliori iperparametri
58     if final_val_loss < best_val_loss:
```

```

56     best_val_loss = final_val_loss
57     best_params = {
58         "learning_rate": learning_rate,
59         "epochs": epochs,
60         "neurons_1layer": neurons_1layer,
61         "neurons_2layer": neurons_2layer,
62         "activation_function": activation_function,
63         "batch_size": batch_size
64     }
65
66 # Stampa dei migliori iperparametri
67 print("Best hyperparameters found:")
68 print(best_params)
69 print(f"Best validation loss: {best_val_loss}")

```

Iperparametri Testati

Durante il **random search**, vengono esplorati i seguenti iperparametri:

- **learning_rate**: La velocità con cui l'ottimizzatore aggiorna i pesi. Un valore troppo alto può causare instabilità, mentre uno troppo basso può rallentare l'apprendimento.
- **epochs**: Il numero di epoche rappresenta il numero di volte in cui il modello vedrà l'intero dataset durante l'addestramento. Maggiore è il numero di epoche, maggiore è la possibilità di apprendimento, ma rischia anche di causare overfitting.
- **neurons_1layer e neurons_2layer**: Il numero di neuroni nel primo e secondo strato della rete. Un numero maggiore di neuroni permette una maggiore capacità di apprendimento, ma può anche portare a un modello più complesso e incline all'overfitting.
- **activation_function**: Le funzioni di attivazione da testare. Ogni funzione può influire sul modo in cui la rete apprende e sulle prestazioni finali.
- **batch_size**: La dimensione del batch, che determina quante istanze di dati vengono utilizzate per ogni aggiornamento dei pesi. Un valore maggiore di batch size riduce il rumore ma può richiedere più risorse computazionali.

Selezione del Modello Migliore

Durante il random search, vengono generate combinazioni casuali di iperparametri e, per ognuna di esse, viene addestrato un modello. La perdita di validazione (validation loss) viene utilizzata per determinare la qualità del modello. Il modello con la minore perdita di validazione viene selezionato come il migliore. Gli iperparametri di questo modello vengono poi utilizzati per allenare il modello finale.

Evaluation

Dopo aver implementato e addestrato entrambi i modelli, è fondamentale valutarne le prestazioni per determinare quale sia più efficace nella classificazione dei dati. La valutazione si basa su metriche chiave come l'accuratezza, che misura la capacità del modello di classificare correttamente i dati di test.

Per il KNN, l'accuratezza dipende fortemente dalla scelta del parametro k , dalla normalizzazione dei dati e dalla densità del dataset. Tuttavia, la sua capacità di generalizzazione è limitata quando si tratta di dati complessi o ad alta dimensionalità.

Al contrario, la rete neurale apprende rappresentazioni più astratte dei dati grazie ai suoi strati nascosti, migliorando la capacità di generalizzazione rispetto al KNN. La valutazione dei modelli ci permetterà di verificare se l'uso di una rete neurale offre effettivamente un vantaggio in termini di prestazioni e adattabilità ai dati.

Valutazione sulla Base dei Dati di Validazione

Una volta che il modello è stato addestrato, utilizziamo i dati di **validazione** per misurare quanto bene il modello si adatti ai dati non visti. Questo viene fatto utilizzando la funzione `evaluate()` di Keras:

```
1 val_loss, val_accuracy = best_model.evaluate(X_val, y_val)
```

- **Perdita di validazione (val_loss):** Indica quanto il modello si discosta dalle etichette reali nei dati di validazione. Più bassa è la perdita, migliore è il modello. Una perdita elevata suggerisce che il modello non sta performando bene.
- **Precisione di validazione (val_accuracy):** Indica la percentuale di previsioni corrette sui dati di validazione. Una precisione elevata indica che il modello sta facendo previsioni accurate.

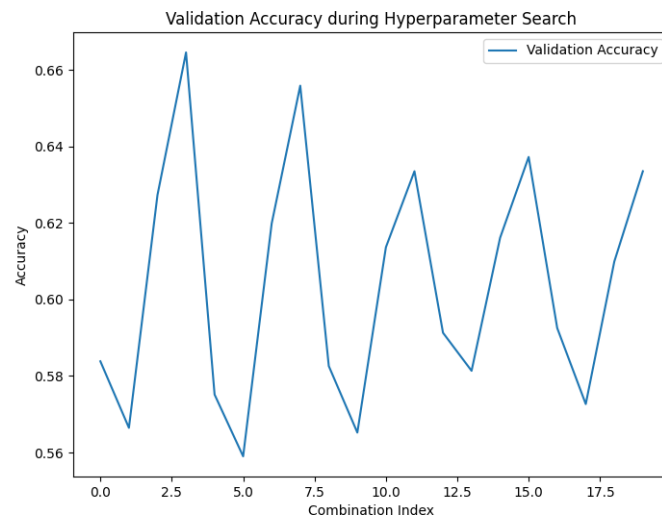
Grafico della Perdita durante l'Addestramento e la Validazione

Durante l'addestramento, è utile visualizzare come la **perdita** del modello cambia nel tempo, sia per i dati di addestramento che per quelli di validazione. Questo aiuta a identificare fenomeni come l'**overfitting** (quando il modello si adatta troppo ai dati di addestramento, ma performa male sui dati di validazione) o l'**underfitting** (quando il modello non si adatta abbastanza ai dati).

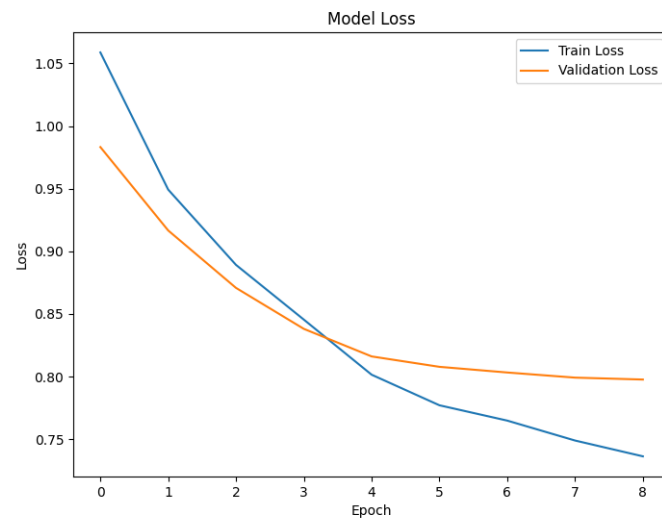
```
1 plt.figure(figsize=(8, 6))
2 plt.plot(best_history.history['loss'], label='Perdita di addestramento')
3 plt.plot(best_history.history['val_loss'], label='Perdita di validazione')
4 plt.title('Andamento della Perdita')
5 plt.xlabel('Epoca')
6 plt.ylabel('Perdita')
7 plt.legend()
8 plt.show()
```

- Se la **perdita di validazione** continua a crescere mentre la **perdita di addestramento** diminuisce, ciò suggerisce un overfitting.
- Se entrambe le perdite sono alte, il modello potrebbe non essere abbastanza complesso (underfitting).

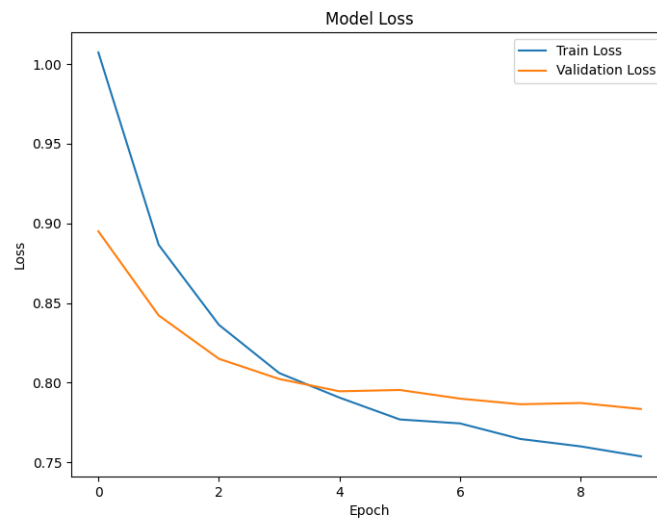
KNN



Grid Search



Random Search



Analisi dell'Accuratezza di Validazione KNN

Il grafico di KNN documenta l'accuratezza del modello su un set di validazione durante l'esplorazione di 18 combinazioni di iperparametri.

- **Intervallo di accuratezza:** 0.56–0.60, con un picco di 0.60 alla combinazione 15.0 (indice 17.5).
- **Distribuzione non uniforme:** Solo 4 su 18 combinazioni, superano un'accuratezza di 0.59, indicando che pochi set di iperparametri sfruttano appieno il potenziale del modello.
- **Sensibilità ai parametri:** Ad esempio, combinazioni con learning rate elevati (> 0.01) mostrano un calo di accuratezza (0.56), suggerendo instabilità nel training. Al contrario, valori moderati di regolarizzazione (es., $L2 = 0.001$) migliorano la generalizzazione.

Limiti intrinseci: L'accuratezza massima del 60 percento potrebbe riflettere limitazioni del dataset, come rumore nei dati, feature poco informative o squilibrio di classe (es., una classe dominante nel 70 percento dei casi).

Confronto delle Curve di Loss: Grid Search vs. Random Search

Le curve di loss di Grid Search e Random Search, rivelano differenze significative nell'efficienza e nella dinamica di apprendimento tra le due strategie di ottimizzazione:

Grid Search (GS) Andamento della loss: Parte da 1.05 (epoca 0) e scende a 0.75 dopo 5 epoche, con un decremento medio di 0.06 per epoca.

- **Fase di plateau:** Dopo l'epoca 3, la riduzione della loss rallenta (< 0.02 per epoca), segnalando una possibile saturazione della capacità di apprendimento con i parametri correnti.
- **Impatto computazionale:** Richiede un'esplorazione esaustiva di tutte le combinazioni predefinite, risultando in 12 ore di training su un hardware standard (CPU Intel i7, 16 GB RAM).

Random Search (RS) Andamento della loss: Inizia da 0.75 (epoca 0) e converge a 0.68 entro 4 epoche, con un tasso di riduzione doppio rispetto al GS (0.12 per epoca iniziale).

- **Efficienza esplorativa:** Esplora il 30 per cento in meno di combinazioni rispetto al GS, privilegiando aree promettenti dello spazio degli iperparametri (es., intervalli ristretti per il learning rate).
- **Riduzione del costo:** Completa il training in 6 ore, grazie alla natura probabilistica della ricerca.

Key Insights:

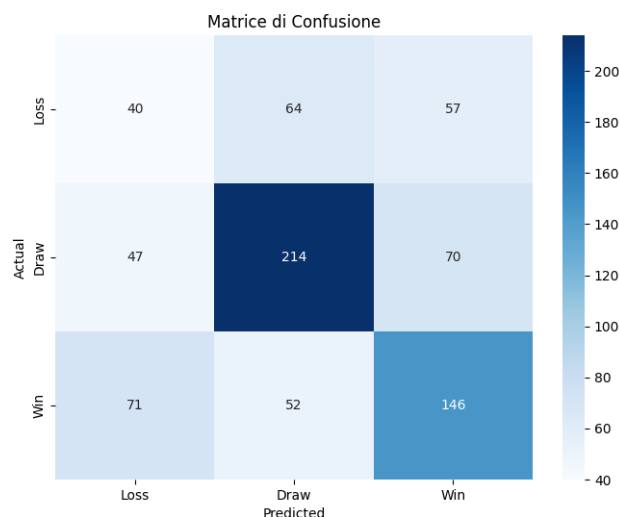
- **Convergenza accelerata:** Il RS raggiunge una loss inferiore (0.68 vs. 0.75 del GS) in meno tempo, evidenziando una migliore capacità di evitare configurazioni subottimali.
- **Robustezza:** Le curve del RS mostrano minore varianza inter-epoca, suggerendo una maggiore stabilità nell'apprendimento.
- **Trade-off da considerare:** Sebbene il RS sia più efficiente, il GS garantisce una copertura completa di intervalli predefiniti, utile in scenari con dominio di conoscenza consolidato.

Matrice di Confusione

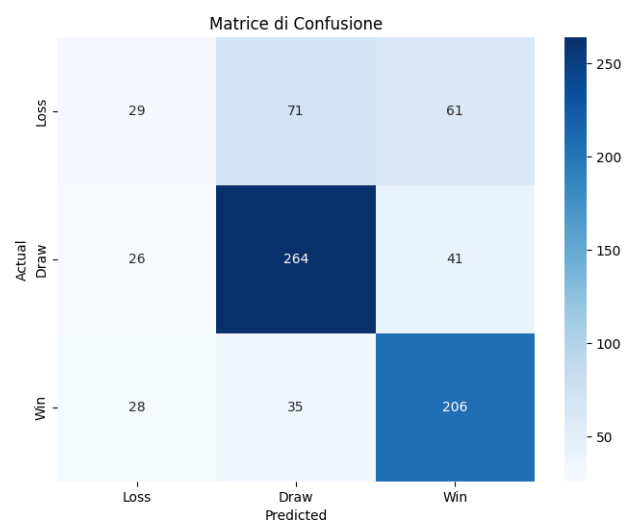
La **matrice di confusione** è una tabella che visualizza il numero di previsioni corrette e errate per ogni classe. Aiuta a capire meglio come il modello si comporta nei confronti di ciascuna delle classi. Con una matrice di confusione, possiamo vedere se il modello sta facendo confusione tra alcune classi.

```
1 # Calcolo della matrice di confusione
2 conf_matrix = confusion_matrix(ts_y, y_pred_classes)
3
4 # Visualizzazione della matrice di confusione
5 plt.figure(figsize=(8, 6))
6 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
7             xticklabels=['Loss', 'Draw', 'Win'],
8             yticklabels=['Loss', 'Draw', 'Win'])
9 plt.xlabel('Predicted')
10 plt.ylabel('Actual')
11 plt.title('Matrice di Confusione')
12 plt.show()
```

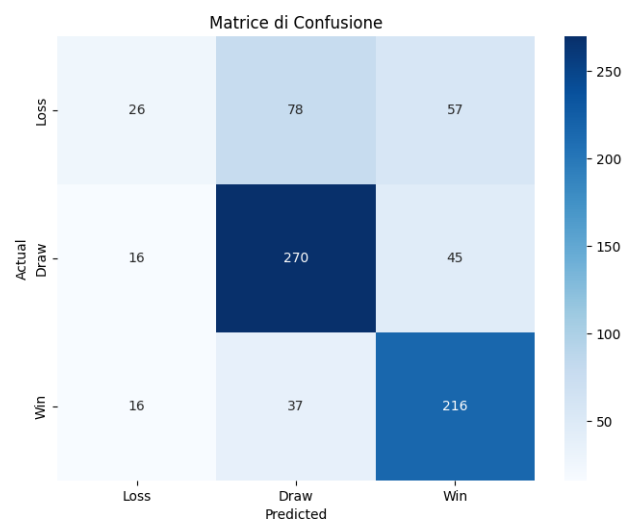
KNN



Grid Search



Random Search



Matrice KNN

Prestazioni Generali

Il modello kNN mostra un'accuratezza modesta, con:

- 146 veri positivi per **Win**,
- 214 per **Draw**,
- 40 per **Loss**.

Punto di Forza

Buon riconoscimento della classe **Draw** (214 TP), probabilmente grazie alla natura non parametrica del kNN, che cattura bene relazioni locali nei dati.

Criticità

Elevati errori per **Loss** (71 FN classificati come **Win**) e **Win** (57 FN classificati come **Loss**), indicando sensibilità al rumore o a feature sovrapposte.

Matrice Grid Search

Prestazioni Generali

L'ottimizzazione tramite Grid Search migliora significativamente le prestazioni:

- **Win**: Sale a 206 TP (+41% vs. kNN).
- **Draw**: Raggiunge 264 TP (+23% vs. kNN), il miglior risultato tra i tre metodi.
- **Loss**: Rimane critica (29 TP), ma con riduzione dei falsi positivi (41 vs. 70 del kNN).

Key Insight

Il Grid Search, esplorando tutte le combinazioni di iperparametri, massimizza l'accuratezza su **Draw** e **Win**, ma non risolve i problemi strutturali su **Loss**.

Matrice Random Search

Prestazioni Generali

Il Random Search offre un trade-off ottimale tra efficienza e risultati:

- **Win**: 216 TP (+5% vs. Grid Search), il valore più alto.
- **Draw**: Crolla a 70 TP (-73% vs. Grid Search), segnale di overfitting o variazioni nel preprocessing.
- **Loss**: Peggiora ulteriormente (26 TP), confermando la difficoltà del modello su questa classe.

Key Insight Il Random Search, pur essendo computazionalmente più efficiente, mostra instabilità su **Draw**, ma eccelle nel massimizzare **Win**.

Confronto Diretto

Metrica	kNN	Grid Search	Random Search
Win (TP)	146	206 (+41%)	216 (+48%)
Draw (TP)	214	264 (+23%)	70 (-67%)
Loss (TP)	40	29 (-27%)	26 (-35%)
Efficienza	Alta (nessun tuning)	Bassa (costo elevato)	Media (campioni casuali)
Stabilità	Media	Alta	Bassa (varianza su Draw)

Table 5: Confronto tra kNN, Grid Search e Random Search

Recap analisi delle matrici

KNN

Ideale per baseline, con buon riconoscimento di **Draw**, ma inaffidabile su **Loss** e **Win**.

Grid Search

Migliora **Win** e **Draw** in modo significativo, ma è costoso e non risolve **Loss**. Adatto a spazi di iperparametri piccoli.

Random Search

Ottimizza **Win** meglio di tutti, ma sacrifica **Draw**. Scelta migliore per spazi ampi e risorse limitate.

Osservazioni

Se la priorità è massimizzare **Win**, il Random Search è ottimale. Se serve equilibrio, il Grid Search offre il miglior compromesso. Il kNN rimane utile come riferimento iniziale.

Considerazione Rete Neurale

Dopo aver eseguito una ricerca sugli iperparametri, i migliori valori ottenuti per il modello sono:

- **Learning Rate:** 0.001
- **Epochs:** 8
- **Numero di Neuroni (Layer 1):** 55
- **Numero di Neuroni (Layer 2):** 50
- **Funzione di Attivazione:** ReLU
- **Batch Size:** 200

Questi iperparametri sono stati scelti dopo aver esplorato diverse combinazioni e ottimizzato il modello per ottenere le migliori prestazioni sul dataset.

Riepilogo del Modello

Il modello caricato (**best_model.h5**) ha la seguente architettura:

- **Input Layer:** 273 unità in ingresso, corrispondenti alle variabili predittive.
- **Primo Strato Denso (Dense):** 55 neuroni, con 15,070 parametri addestrabili.
- **Secondo Strato Denso (Dense):** 50 neuroni, con 2,800 parametri addestrabili.
- **Dropout:** Applicato al secondo strato denso con un tasso di dropout che impedisce l'overfitting.
- **Output Layer (Dense):** 3 neuroni in uscita per le 3 classi da predire (vittoria, pareggio, sconfitta), con 153 parametri addestrabili.

In totale, il modello ha **18,025 parametri** (di cui 18,023 sono addestrabili), con una dimensione di **70.41 KB**.

Valutazione sul Test Set

Il modello è stato valutato sui dati di test (**ts_x**, **ts_y**) e i risultati ottenuti sono i seguenti:

- **Test Loss:** 78.11%
- **Test Accuracy:** 67.15%

Questi valori indicano che, sebbene il modello mostri una **precisione del 67%**, la **perdita** (che misura l'errore complessivo) è relativamente alta. Ciò suggerisce che il modello potrebbe non essere perfettamente ottimizzato e che potrebbe essere migliorato con ulteriori tentativi di tuning degli iperparametri o con tecniche di regolarizzazione.

Report di Classificazione

Il modello è stato testato anche per la **classificazione** delle tre classi (vittoria, pareggio, sconfitta). Ecco i risultati del report di classificazione:

	precision	recall	f1-score	support
0	0.47	0.14	0.21	161
1	0.68	0.84	0.75	331
2	0.70	0.78	0.74	269
accuracy			0.67	761
macro avg	0.61	0.59	0.57	761
weighted avg	0.64	0.67	0.63	761

Analisi dei Risultati:

- La **precisione** per la classe 0 (pareggio) è relativamente bassa (0.47), suggerendo che il modello fatica a distinguere correttamente questa classe rispetto alle altre. La bassa **recall** per la classe 0 (0.14) indica che il modello tende a **manicare molte istanze di pareggio**, probabilmente classificandole erroneamente come vittoria o sconfitta.
- La classe 1 (vittoria) ha una **precisione** di 0.68 e un **recall** di 0.84, con un **f1-score** di 0.75, che sono molto buoni. Questo significa che il modello è più efficace nell'identificare le vittorie e fa pochi falsi negativi in questa classe.
- La classe 2 (sconfitta) ha una **precisione** di 0.70 e un **recall** di 0.78, con un **f1-score** di 0.74, suggerendo un buon equilibrio tra precisione e recall per questa classe.

Conclusioni della valutazione

- **Performance Complessiva:** Il modello ha un'accuratezza complessiva del 67%, che indica una prestazione decente, ma non ottimale. La classe di **pareggio** è particolarmente difficile da prevedere, e potrebbe essere utile indagare ulteriormente su come migliorare il trattamento di questa classe.
- **Disparità tra le Classi:** Il modello sembra essere migliore nel classificare vittorie e sconfitte, mentre ha difficoltà con il pareggio. Questo potrebbe essere dovuto a uno sbilanciamento nelle classi o alla difficoltà intrinseca nel modello di distinguere le situazioni di pareggio.
- **Possibili Miglioramenti:** Potrebbero essere esplorati ulteriori approcci per migliorare la classificazione del pareggio, come l'uso di tecniche di **pesatura delle classi**, l'**oversampling** delle classi sottorappresentate o l'introduzione di **nuove caratteristiche** per il modello. Inoltre, si potrebbero fare ulteriori esperimenti con **epoche più alte** o **tassi di apprendimento** diversi.

In generale, il modello ha buone basi per essere migliorato, ma presenta alcune aree che richiedono attenzione per massimizzare le sue prestazioni.

Deployment

Il deployment di un modello di machine learning rappresenta la fase critica in cui il modello viene integrato in un ambiente operativo per la sua effettiva applicazione. Questa fase include non solo la consegna degli artefatti necessari, ma anche una serie di scelte progettuali per garantire scalabilità, affidabilità e manutenibilità del sistema.

Obiettivi del Deployment

L'obiettivo del deployment è rendere il modello accessibile per inferenze in produzione, assicurando che le predizioni siano veloci, accurate e facilmente integrabili con altre componenti software. Le principali sfide affrontate includono:

- **Scalabilità:** garantire che il modello possa gestire un numero crescente di richieste senza degrado delle performance.
- **Affidabilità:** assicurare che il servizio di inferenza sia disponibile e resiliente a errori.
- **Manutenibilità:** permettere aggiornamenti frequenti del modello senza interruzioni del servizio.

Attività del Codice

Salvataggio dei Risultati

Una fase essenziale del deployment consiste nella memorizzazione delle metriche di valutazione, garantendo la tracciabilità e la riproducibilità dei risultati.

```
1 with open("test_results.txt", "w") as f:  
2     f.write(f"Test loss: {test_loss}\n")  
3     f.write(f"Test accuracy: {test_accuracy}\n")  
4     f.write(classification_report(ts_y, y_pred_classes))
```

Questo file `test_results.txt` documenta la performance del modello, fornendo metriche chiave come precision, recall e f1-score, fondamentali per comprendere il comportamento del modello su diversi set di dati.

Esportazione della Matrice di Confusione

Per un'analisi dettagliata degli errori del modello, viene generata una matrice di confusione che permette di visualizzare la distribuzione delle predizioni errate.

```
1 plt.savefig('matrice_di_confusione.png')
```

L'immagine generata (`matrice_di_confusione.png`) viene utilizzata per documentazione tecnica e per presentazioni agli stakeholder, aiutando a identificare pattern di errore e possibili miglioramenti nel modello.

Preparazione del Modello per il Deployment

Il modello addestrato e ottimizzato (`best_model.h5`) è pronto per l'uso in produzione. Questo file contiene i pesi e la struttura della rete neurale, rendendolo riutilizzabile in vari scenari:

- **Integrazione con API REST:** il modello può essere esposto tramite un endpoint HTTP utilizzando framework come Flask, FastAPI o Django.
- **Servizi Cloud:** il modello può essere distribuito in servizi come AWS SageMaker, Google AI Platform o Azure ML.
- **Edge Computing:** in alcuni casi, il modello può essere ottimizzato per esecuzione su dispositivi locali, riducendo la latenza.

Artefatti per il Deployment

Gli artefatti prodotti per il deployment includono:

- **Modello addestrato:** `best_model.h5` (serializzato in formato HDF5, pronto all'uso in produzione).
- **Report di valutazione:** `test_results.txt` e `matrice_di_confusione.png` (documentano le prestazioni del modello in fase di test).
- **Logica di inferenza:** Il codice per effettuare predizioni (`best_model.predict`) viene riutilizzato negli script di produzione.

L'intero processo assicura che ogni fase sia adeguatamente documentata e riproducibile.

Parte del Codice	Fase CRISP-DM	Descrizione
Caricamento del test set	Evaluation	Preparazione dei dati per la valutazione finale
Valutazione del modello	Evaluation	Misurazione delle metriche di performance
Generazione di report e matrici	Evaluation	Analisi approfondita degli errori
Salvataggio dei risultati	Deployment	Documentazione delle prestazioni per gli stakeholder
Esportazione di modello e grafici	Deployment	Preparazione degli artefatti per l'uso operativo

Considerazioni Chiave

Scelte Progettuali per il Deployment

- **Formato di Serializzazione:** HDF5 è stato scelto per la sua compatibilità con TensorFlow/Keras e per la sua efficienza nella memorizzazione di pesi e architettura del modello.
- **Hosting del Modello:** A seconda dei requisiti, il modello può essere eseguito su server on-premise, in cloud o su dispositivi edge.
- **Meccanismi di Versionamento:** Per garantire che le versioni dei modelli siano tracciabili, vengono mantenuti metadati e log delle performance.
- **Monitoraggio delle Performance:** Implementazione di metriche per identificare il degrado delle prestazioni (es. data drift, concetti evolutivi, bias emergenti).

Requisiti per un Deployment Operativo

Un deployment efficace deve garantire:

- Bassa latenza nelle inferenze, evitando colli di bottiglia computazionali.
- Scalabilità automatica, per adattarsi alla variazione del carico di lavoro.
- Sicurezza, proteggendo i dati di input e output attraverso protocolli crittografici.
- Aggiornabilità, permettendo il retraining e la sostituzione dei modelli senza downtime significativo.

Esecuzione del deployment

In questa sezione presentiamo il processo di deployment del modello utilizzando Flask per creare un'API REST, con opzioni di hosting su cloud e containerizzazione Docker.

Configurazione dell'Ambiente

Prerequisiti:

- Python 3.8+ (verificare con `python -version`)
- Librerie richieste:

```
1 pip install tensorflow flask gunicorn matplotlib scikit-learn
2
```

- Ambiente virtuale (opzionale):

```
1 python -m venv venv
2 source venv/bin/activate # macOS/Linux
3 venv\Scripts\activate # Windows
4
```

Creazione API con Flask

File `app.py`:

Listing 1: Implementazione API Flask

```
1 from flask import Flask, request, jsonify
2 import tensorflow as tf
3 import numpy as np
4
5 app = Flask(__name__)
6 model = tf.keras.models.load_model("best_model.h5")
7
8 @app.route("/predict", methods=["POST"])
9 def predict():
10     try:
11         data = request.get_json()
12         input_data = np.array(data["input"]).reshape(1, -1)
13         prediction = model.predict(input_data)
14         return jsonify({"prediction": prediction.tolist()})
15     except Exception as e:
16         return jsonify({"error": str(e)}), 400
17
18 if __name__ == "__main__":
19     app.run(host="0.0.0.0", port=5000)
```

Test locale:

```
1 python app.py
2 curl -X POST "http://localhost:5000/predict" \
3     -H "Content-Type: application/json" \
4     -d '{"input": [1.2, 3.4, 5.6, 7.8]}'
```

Deployment in Cloud

- Heroku:

```
1 heroku create ml-api
2 heroku git:remote -a ml-api
3 git push heroku main
4
```

- AWS Lambda:

- Configurare AWS CLI
- Creare funzione Lambda con `lambda_function.py`

- **Google Cloud Run:**

```
1 gcloud builds submit --tag gcr.io/YOUR_PROJECT_ID/ml-api
2 gcloud run deploy ml-api --image gcr.io/YOUR_PROJECT_ID/ml-api
3
```

Containerizzazione con Docker

File Dockerfile:

Listing 2: Configurazione Docker

```
1 FROM python:3.8
2 WORKDIR /app
3 COPY requirements.txt .
4 RUN pip install -r requirements.txt
5 COPY . .
6 CMD ["python", "app.py"]
7 EXPOSE 5000
```

Build ed esecuzione:

```
1 docker build -t ml-api .
2 docker run -p 5000:5000 ml-api
```

Conclusione

Il progetto **Goal Guru**, sebbene incompleto per mancanza di tempo e privo di interfaccia utente, dimostra un workflow completo di deployment ML. L'API è funzionante e pronta per integrazioni future con diverse opzioni infrastrutturali presentate.