

Flavor Craft

GitHub

Ferdinando Tessitrice, 0512115432



Dipartimento di Informatica
Universita` degli Studi di Salerno

10/02/2024

Contents

1 Introduzione	2
1.1 Contesto	2
1.2 Idea	2
2 PEAS	3
2.1 Specifiche.	3
3 Dataset	4
3.1 Creazione Dataset	4
3.2 Data Understanding, Data Preparation & Data Cleaning... ..	5
3.3 Feature Scaling	8
3.4 Feature Selection... ..	8
4 Prima Soluzione :Algoritmi genetici	10
5 Seconda Soluzione Algoritmo A*	12
6 Terza Soluzione TF-IDF + similiarit`a del coseno	14
6.1 Word Embedding	14
6.2 TF-IDF	15
6.3 Word2vec.	16
6.4 Similarità del coseno... ..	18
6.5 KNN	19
7 Quarza Soluzione :Clustering	22
7.1 Prima Strategia	22
7.2 Seconda Strategia.	24
8 App	26
8.1 GUI	26
8.2 Input Vocale.	28
9 Riferimenti	29

Chapter 1

Introduzione

1.1 Contesto

Il progetto nasce dal problema che molte persone hanno, ovvero avere in dispensa/frigorifero degli ingredienti dimenticati oppure accumulati e non sapere cosa cucinare. Oppure banalmente ottenere nuove idee per delle ricette in base ad altre che ci piacciono.

1.2 Idea

Il progetto ha lo scopo di creare un agente intelligente capace di consigliare una o più ricette in base a quello che l'utente ha a disposizione, o che gli piace.

Chapter 2

PEAS

2.1 Specifiche

La prima cosa da fare quando si progetta un agente intelligente `e la definizione delle Specifiche Peas.

I Performance: Migliore ricetta in base agli ingredienti forniti

I Enviroment: L'ambiente nel quale il nostro agente opera `e composto da tutte le ricette del dataset con i relativi titolo, ingredienti, ratings, tempo di preparazione, istruzioni, tipologia.

L'ambiente risulta:

- **Completamente osservabile:** si ha accesso a tutte le informazioni di ogni ricetta
- **Deterministico:** L'ambiente deterministico significa che l'agente sa esattamente cosa aspettarsi in base alle informazioni fornite e agli input dell'utente. In altre parole, l'agente sa che dati gli stessi input, otterrà sempre lo stesso output, l'agente sa che se gli vengono forniti gli stessi ingredienti, troverà sempre la stessa ricetta migliore nell'ambiente statico
- **Episodico:** l'agente si trova in un singolo episodio in cui riceve gli ingredienti in input e restituisce la migliore ricetta come output, senza alcuna interazione ulteriore con il mondo esterno.
- **Statico:** L'insieme delle ricette e le relative informazioni non cambiano nel corso del tempo
- **Discreto:** L'agente ha a disposizione un numero limitato di azioni
- **Singolo Agente:** `e presente solo un agente che opera nell'ambiente

I Attuatori: L'algoritmo utilizzato per generare le possibili soluzioni e selezionare la migliore.

I Sensori: Gli input dell'utente, ovvero la lista degli ingredienti.

Chapter 3

Dataset

3.1 Creazione Dataset

Non esistendo dei dataset di ricette in italiano, ho dovuto crearlo manualmente. Per farlo ho utilizzato la tecnica del Web Scraping, ovvero una tecnica tramite la quale estrarre dati da un sito web, nel mio caso da GialloZafferano. Ho utilizzato una libreria di python chiamata recipe-scrapers insieme alla libreria BeautifulSoup. BeautifulSoup è stata utilizzata per prendere tutti i link delle ricette delle 426 pagine di GialloZafferano.

```
def searchRecipeLink():
    url_list = ["https://www.giallozafferano.it/ricette-cat/"]
    url_recipe_list = []

    count = 1
    while count < 425:
        url = "https://www.giallozafferano.it/ricette-cat/page" + str(count) + "/"
        url_list.append(url)
        count += 1

    for url in url_list:
        html = requests.get(url).content
        soup = BeautifulSoup(html, 'html.parser')
        links = soup.find_all('a', href=True)
        for link in links:
            recipe_url = link['href']
            if recipe_url.startswith("https://ricette.giallozafferano.it"):
                if ("#anchor=qz-comments-anchor" not in recipe_url):
                    if (recipe_url not in url_recipe_list):
                        url_recipe_list.append(recipe_url)

    return url_recipe_list
```

Dopodichè per ogni pagina di una ricetta sono state usate le funzioni di recipe scraper per ottenere le feature da inserire nel dataset quali: titolo ricetta, tempo preparazione,

valutazione, categoria/tipologia, istruzioni.

```
def createRecipesList(url_recipe_list):
    recipes = []
    count = 0
    for url in url_recipe_list:
        scraper = scrape_me(url)

        print(str(count) + " " + scraper.title())
        count += 1
        ingredients = cleanIngredients(scraper.ingredients())
        #ingredients = scraper.ingredients()

        recipe = {'title': scraper.title(), 'ingredients': ingredients, 'time': scraper.total_time(),
                  'ratings': scraper.ratings(), 'type': scraper.category(), 'instructions': scraper.instructions()}

        # Aggiungi le informazioni della ricetta alla lista
        recipes.append(recipe)

    return recipes
```

Poi tutte le ricette sono state salvate in un file csv

```
def saveRecipe2CSV(recipes):
    import csv

    # Apri il file in modalità di scrittura
    with open('ricetteEsatte.csv', mode='w', newline='') as csv_file:
        writer = csv.writer(csv_file)
        writer.writerow(['title', 'ingredients', 'time', 'ratings', 'type', 'instructions'])

    # Scrivi le informazioni delle ricette nel file
    for recipe in recipes:
        row = [
            recipe['title'],
            recipe['ingredients'],
            recipe['time'],
            recipe['ratings'],
            recipe['type'],
            recipe['instructions']
        ]

        writer.writerow(row)
    print("dataset cvs creato\n")
```

3.2 Data Understanding, Data Preparation & Data Cleaning

Ho creato un programma per l'esplorazione del dataset in modo tale da effettuare il conteggio delle ricette, del tempo medio per la preparazione, del numero di ricette per tipologia, del numero ingredienti totali, del numero di ingredienti più usati, di quelli di meno usati

Un primo problema si è presentato con la lista degli ingredienti. Se prendiamo ad esempio la ricetta “Tiramisu” i suoi ingredienti sono:

Mascarpone 750 g

Uova (freschissime, circa 5 medie) 260 g

Savoardi 250 g

Zucchero 120 g

Caffè (della moka, zuccherato a piacere) 300 g

A noi interessano solo i nomi degli ingredienti, in quanto saranno utilizzati per trovare ricette simili dati degli ingredienti in input. Quindi andavano cancellate tutte le grammature e altre informazioni non necessarie come “freschissime, circa 5 medie”.

Per fare questo è stato creato un programma in python che andava a rimuovere le grammature, numeri, informazioni nelle parentesi e informazioni di contorno in generale. Inoltre bisognava anche andare a sostituire alcuni ingredienti composti usando solo l'ingrediente principale ad esempio: Estratto di vaniglia, doveva diventare vaniglia.

Poi sono stati rimossi tutti i caratteri non necessari come: [,], (,) /, ecc

Sono stati rimossi tutti quei dettagli come “piccola”. E se c'è piccolo o piccoli. Ho rimosso tutte le varianti manualmente,

- Tutti gli ingredienti sono stati portati in lower case.

Fatto ciò la lista di ingredienti pulita è stata salvata come stringa ed inserita al posto della lista di ingredienti precedente alla pulizia. Alla fine del processo veniva creato un nuovo dataset.

La fase di pulizia del dataset però non finiva qui.

Grazie al conteggio per ogni singolo ingrediente dell'esplorazione dataset aggiornata. Mi ha permesso di pulire ulteriormente il dataset da ricette superflue. Difatti ho eliminato circa un centinaio di ricette che avevano ingredienti unici, troppo specifici non mi avrebbero aiutato nell'identificazione di ricette simili, in quanto uniche nel loro genere (Per il codice vedere removeRow.py in RecipeScraper)

```
# Carica il dataset
df = pd.read_csv('ricette.csv', encoding='iso-8859-1')

df.dropna(subset=['ingredients'], inplace=True)

print(df.shape)
print(df.head)
```

[6037 rows x 6 columns]

Il nostro dataset ha poco più di 6000 ricette ed è un file CSV.

```
(6037, 6)
<class 'pandas.core.frame.DataFrame'>
Int64Index: 6037 entries, 0 to 6038
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype
---  -
0   title                  6037 non-null   object
1   ingredients             6037 non-null   object
2   time                   6037 non-null   int64
3   ratings                6037 non-null   float64
4   type                   6037 non-null   object
5   instructions            6036 non-null   object
dtypes: float64(1), int64(1), object(4)
memory usage: 330.1+ KB
None
```

Esempio prime righe. . .

Titolo	Ingredienti	tempo	Valutazione	Tipologia	Istruzioni
Tiramisù	mascarpone,uova, 46 savoiardi,zucchero, caffè,cacao		4.2	Dolci	Per preparare il tiramisù preparate il caffè...
Spaghetti alla carbonara	spaghetti,guanciale,25 uova,pecorino,pepe		4.2	Primi piatti	Per preparare gli spaghetti alla carbonara...
ecc...					

Conteggio ingredienti `e stato fatto sia con la libreria nltk sia con la libreria csv leggerle righe:

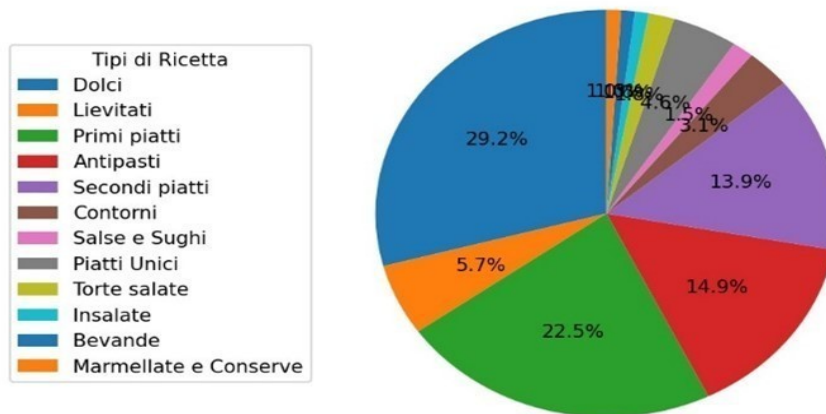
```
vocabulary = nltk.FreqDist()
for ingredients in df['ingredients']:
    ingredients = ingredients.split(',')
    vocabulary.update(ingredients)

for word, frequency in vocabulary.most_common(200):
    print(f'{word}: {frequency}')
    fdist = nltk.FreqDist(ingredients)

common_words = []
for word, _ in vocabulary.most_common(250):
    common_words.append(word)
print(common_words)
```

```
sale: 4385
olio extravergine d'oliva: 3043
pepe: 2954
farina: 2461
uova: 2432
zucchero: 2125
burro: 1929
acqua: 1531
latte: 1371
aglio: 1324
cipolle: 1182
limone: 1106
lievito: 1095
```

Conteggio ricette per Tipologia :



Numero di ricette per tipo:

Dolci: 1761
 Lievitati: 342
 Primi piatti: 1356
 Antipasti: 898
 Secondi piatti: 840
 Contorni: 189
 Salse e Sughi: 90
 Piatti Unici: 277
 Torte salate: 109
 Insalate: 58
 Bevande: 58
 Marmellate e Conserve: 61

3.3 Feature Scaling

Non c'erano particolarmente discordantie quindi non c'è stato bisogno di scalare e normalizzare.

3.4 Feature Selection

Come già detto, il dataset delle ricette è stato creato e quindi di seguito vengono spiegati i motivi di ogni caratteristica:

- Titolo: è utile per mostrare all'utente la ricetta consigliata.
- Ingredienti: è la feature dominante in quanto ci permette di capire se due ricette sono simili.
- Rating: utile per capire se una determinata ricetta piace agli Utenti o meno (non verrà usata nella soluzione finale).

- Prep_time: `e utile per capire il tempo di preparazione della ricetta.
- Categoria: `e utile per capire la tipologia della ricetta.

Chapter 4

Prima Soluzione :Algoritmi genetici

Ho provato a realizzare un agente intelligente attraverso gli Algoritmi Genetici.
(per il codice vedere cartella GA)

Un GA è una procedura ad alto livello (meta-euristica) ispirata alla genetica per definire un algoritmo di ricerca. Evolve una popolazione di individui (soluzioni candidate) producendo di volta in volta soluzioni sempre migliori rispetto ad una funzione obiettivo, fino a raggiungere l'ottimo o un'altra condizione di terminazione. La creazione di nuove generazioni di individui avviene applicando degli operatori genetici, precisamente selezione, crossover e mutazione.

Nel mio caso le caratteristiche del GA erano le seguenti: Codifica individuo (ricetta) era formata da un array di 8 interi che potevano assumere il valore solo di 0 o 1. L'intero individuo andava ad individuare l'indice di una riga del dataset e quindi di una ricetta.

Size popolazione: 5

Selezione ho usato la tecnica dei State GA, invece di generare nuove generazioni si va a migliorare costantemente quella attuale. Si selezionano i due migliori individui come genitori

Crossover: il crossover (Single/Two-Point) avviene tra i due genitori selezionati al punto precedente. Il nuovo individuo andrà a sostituire il peggiore della popolazione. Nel caso in cui si va a creare un individuo già esistente nella popolazione allora si applica la mutazione (bit flip)

Stopping condition terminiamo l'evoluzione quando si superano le 500 iterazioni o quando si è raggiunta una certa fitness.

La funzione di fitness è data dal numero di ingredienti in comune con quelli inseriti moltiplicati per i likes il tutto diviso per il tempo di preparazione.

Dopo varie prove e modifiche agli operatori genetici ci siamo resi conto che l'algoritmo era un po' troppo casuale. Molte volte trovava le due migliori ricette in una 50ina di generazioni altre volte arrivava fino a 500 iterazioni trovando due ricette con fitness molto basse.

Inserisci gli ingredienti separati da virgole o spazi:

farina, acqua, sale, lievito, pomodoro, mozzarella, basilico

Soluzione trovata nella generazione 50

Fitness: 46.666668

Individuo 1 :Kebab turco. Fintness: 8.333333

Individuo 2 :Pescado frito spagnolo. Fintness: 10.0

Individuo 3 :Parmigiana di zucchini. Fintness: 11.333333

Individuo 4 :Ratatouille. Fintness: 6.666665

Individuo 5 :Pizza Margherita. Fintness: 46.666668

Individuo 6 :Riso e fagioli. Fintness: 6.0

Individuo 7 :Bistecca alla Fiorentina. Fintness: 7.5

Individuo 8 :Panzanella. Fintness: 10.666667

Individuo 9 :Gnocchi al pomodoro. Fintness: 15.0

Individuo 10 :Churros. Fintness: 15.2

Ricetta più adatta: Pizza Margherita

ingredienti: [farina, acqua, sale, lievito, pomodoro, mozzarella, basilico]

Altra ricetta: Churros

ingredienti: [farina, acqua, olio, zucchero, cannella]

Così ho pensato che la codifica non era adatta o probabilmente non era adatto al nostro problema un algoritmo Genetico.

Chapter 5

Seconda Soluzione : Algoritmo A*

Quindi ho optato per un algoritmo di ricerca informata, il più utilizzato tra gli Algoritmi di ricerca informata è l'A*. (Dataset Vecchio utilizzato, title, ingredients, likes, prep time, Solo 300 ricette)

Per il codice vedere la cartella A Star

L'algoritmo A* è un algoritmo di ricerca che utilizza una funzione di valutazione per determinare la distanza tra un nodo corrente e un nodo obiettivo. In questo caso, la funzione di valutazione è "evaluate recipe" che valuta una singola ricetta sulla base di quanti ingredienti corrispondono a quelli forniti in input e sulla base del numero di likes della ricetta.

La funzione **"evaluate recipe"** valuta una singola ricetta in base agli ingredienti Forniti in input. La funzione utilizza un ciclo for per verificare se ciascun ingrediente della ricetta è presente negli IngredientForniti in input. Se un ingrediente è presente, la valutazione della ricetta viene incrementata di 1. Inoltre, la valutazione della ricetta viene incrementata anche dal numero di "likes" della ricetta diviso 100, per normalizzare i "likes"

```
def evaluate_recipe(recipe, input_ingredients):
    score = 0
    for ingredient in recipe.ingredients:
        if ingredient in input_ingredients:
            score += 1
    score += recipe.likes / 100 # dividere per 100 per normalizzare i likes
    return score
```

La funzione **"find best recipe"** utilizza l'algoritmo A* per trovare la migliore ricetta possibile. Inizia creando due liste "open list" e "closed list". "open list" contiene tutte le ricette del dataset all'inizio, mentre "closed list" è vuota.

All'interno del ciclo while, l'algoritmo seleziona la ricetta con la valutazione più alta dalla lista "open list" e la sposta nella lista "closed list". Successivamente, verifica se tutti gli ingredienti della ricetta corrente sono presenti Negli ingredienti forniti in input. Se lo sono, la ricetta viene aggiunta alla lista "best recipes".

L'algoritmo termina quando la lista "open list" è vuota e restituisce la lista "best recipes" contenente le migliori ricette

```
# funzione di ricerca A*
def find_best_recipe(recipes, input_ingredients):
    open_list = []
    closed_list = []
    best_recipes = []
    best_score = 0
    alternative_recipes = []
    missing_ingredient = []

    open_list = recipes

    while open_list:
        # scegli la ricetta con la valutazione più alta
        current_recipe = max(open_list, key=lambda x: evaluate_recipe(x, input_ingredients))
        open_list.remove(current_recipe)
        closed_list.append(current_recipe)

        score = evaluate_recipe(current_recipe, input_ingredients)
        if score != 0:
            if score == len(current_recipe.ingredients):
                best_recipes.append(current_recipe) # aggiungo la ricetta corrente alla lista delle migliori ricette
            elif len(current_recipe.ingredients) - score == 1:
                otherPossibleRecipe(alternative_recipes, missing_ingredient, current_recipe, input_ingredients)

    return best_recipes, alternative_recipes, missing_ingredient
```

Il programma funzionava abbastanza bene riusciva a dare in output le ricette che ottimizzavano gli ingredienti forniti in input ed anche le ricette a cui mancava un solo ingrediente di quelli a disposizione.

```
Inserisci gli ingredienti separati da una virgola: caffè,uova,zucchero,panna,savoiardi,maionese,pomodoro,formaggio,gnocchi,mozzarella

-----RICETTA MIGLIORE: Tiramisù
-----Ingredienti usati: ['uova', 'zucchero', 'caffè', 'panna', 'savoiardi']

-----POSSIBILI RICETTE:

--- Gnocchi al pomodoro ---
-----Ricetta 1  Ingredienti richiesti: ['gnocchi', 'pomodoro', 'basilico', 'mozzarella']
-----Ricetta 1  Ingrediente mancante da comprare: basilico
--- Gnocchi alla sorrentina ---
-----Ricetta 2  Ingredienti richiesti: ['gnocchi', 'pomodoro', 'mozzarella', 'basilico']
-----Ricetta 2  Ingrediente mancante da comprare: basilico
```

Il problema era che era troppo troppo lento, già con un dataset piccolo, inoltre, analizzando meglio il codice ho visto che probabilmente era una semplice ricerca e non un algoritmo A Star, in quanto non utilizza $h(n)$.

Chapter 6

Terza Soluzione TF-IDF + similarità del coseno

Fatte un paio di ricerche sul web ho compreso che era un problema facente parte del NLP, una branca dell'AI. Nello specifico un Problema di Text Classification, ovvero identificare gli argomenti di un testo nel mio caso gli ingredienti delle ricette e la loro correlazione.

Per poter risolvere i problemi di Text Classification si pu`o utilizzare il word embedding

6.1 Word Embedding

La semplice definizione di word embedding `e la conversione di testo in numeri. Per far s` che il computer capisca il linguaggio, convertiamo il testo in forma vettoriale in modo che i computer possano sviluppare connessioni tra i vettori e le parole e capire cò che diciamo.

Con il word embedding, risolviamo i problemi relativi al Natural Language Processing. Le rappresentazioni Vettoriali dei tokens, chiamate Word Embeddings, sono apprese da modelli NLP, definiti come modelli spaziali vettoriali

Le singole parole sono analizzate e rappresentate atomicamente come unità singole in **Bag-of-Words** (BOW): questo 'bagaglio' di parole, o 'set', sono una rappresentazione sparsa delle parole e della loro presenza, indipendentemente dall'ordine sintattico in cui appaiono in una data frase.

L'idea `e molto semplice: ogni documento del nostro corpus viene rappresentato contando quante volte ogni parola appare in esso. Ad esempio se prendiamo due ricette e quindi due liste di ingredienti (documenti)

R1 Tiramis`u : "mascarpone,uova,savoiardi,zucchero,caff`e,cacao"

R2 Spaghetti alla Carbonara : "spaghetti,guanciale,uova,pecorino,pepe"

	mascarpone	uova	savoiardi	zucchero	caff`e	cacao	spaghetti	guanciale	pecorino	pepe
R1	1	1	1	1	1	1	0	0	0	0
R2	0	1	0	0	0	0	1	1	1	1

Il modello TF-IDF viene addestrato sull'intero dataset di ingredienti utilizzando il metodo fit().

Il modello viene quindi utilizzato per trasformare questi ingredienti in una rappresentazione numerica utilizzando il metodo transform(). Questa rappresentazione viene chiamata "codifica TF-IDF".

Per Misurare la somiglianza tra le liste degli ingredienti e gli ingredienti dati in input ho utilizzato la Similarità del coseno, funzione di Sklearn con le stesse somiglianze e per i dati testuali viene utilizzata per trovare la somiglianza di testi nel documento. Viene usata la funzione ottieni Consigli per classificare i punteggi e generare un dataframe, tramite la libreria pandas, contenente titolo, gli ingredienti il punteggio delle ricette consigliate.

La similarità del coseno è una misura della somiglianza tra due punti dati in un piano. La somiglianza del coseno viene utilizzata per determinare la distanza tra i vicini, ad esempio nei sistemi di raccomandazione viene utilizzata per consigliare film.

I risultati ottenuti sono accettabili ma molte volte consiglia delle ricette per niente simili.

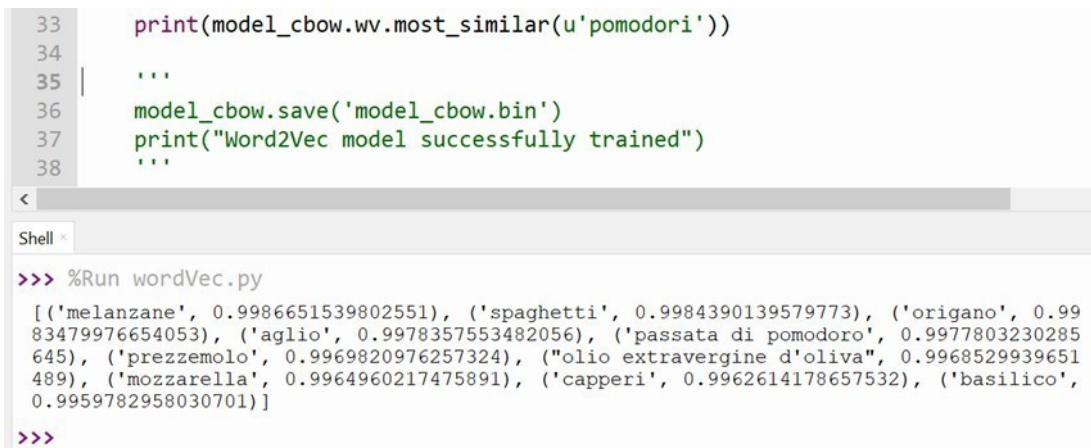
Per migliorarlo avevamo bisogno di tenere traccia degli ingredienti comunemente usati insieme.

6.3 Word2vec

IL problema era: Come facciamo a far comprendere all'agente intelligente che Alcuni ingredienti stanno bene insieme dato che sono presenti nelle ricette nel dataset in cui compaiono gruppi di ingredienti insieme? Per fare questo ho addestrato un modello chiamato Word2vec in quanto è altamente efficiente nel raggruppare insieme parole simili.

Per esempio i pomodori sono simili ad altri ingredienti usati spesso insieme.

```
33 print(model_cbow.wv.most_similar(u'pomodori'))
34
35 ...
36 model_cbow.save('model_cbow.bin')
37 print("Word2Vec model successfully trained")
38 ...
```



```
<
Shell
>>> %Run wordVec.py
[('melanzane', 0.9986651539802551), ('spaghetti', 0.9984390139579773), ('origano', 0.9983479976654053), ('aglio', 0.9978357553482056), ('passata di pomodoro', 0.9977803230285645), ('prezzemolo', 0.9969820976257324), ('olio extravergine d'oliva', 0.9968529939651489), ('mozzarella', 0.9964960217475891), ('capperi', 0.9962614178657532), ('basilico', 0.9959782958030701)]
>>>
```

Word2vec è una semplice rete neurale artificiale a due strati progettata per elaborare il linguaggio naturale. L'algoritmo richiede in ingresso un corpus e restituisce un insieme di vettori che rappresentano la distribuzione semantica delle parole nel testo.

Per addestrare il Modello di input include il numero di neuroni alle parole del vocabolario. La dimensione dello strato di uscita e dello strato di ingresso rimane la stessa. Tuttavia, la dimensione dello strato nascosto è impostata in base ai Vettori delle dimensioni delle parole risultanti. È possibile eseguire l'incorporazione di parole con Word2Vec attraverso due metodi. In entrambi i metodi sono necessarie reti neurali artificiali. Questi metodi sono:

- **CBOW** : è denominato metodo "Continuous Bag of Words" (CBOW), per predire una parola target basata su una sequenza contigua di n-grammi.
- **Skim Gram** : usa un vettore di parole per predire un contesto target

Nel mio caso ho usato CBOW

Per realizzarlo in python ho usato la libreria gensim e seguito alcune guide online come questa: guida.

Dovevamo rappresentare ogni documento corpus (lista di ingredienti per ricetta) come un singolo incorporamento in modo tale da calcolare le somiglianze.

Di seguito il codice per allenare il modello Word2Vec

```
df = pd.read_csv('ricette6k.csv', encoding='iso-8859-1')

df.dropna(subset=['ingredients'], inplace=True)
# ottieni il corpus
corpus = get_and_sort_corpus(df)

# allena e salva il modello Word2Vec CBOW
model_cbow = Word2Vec(corpus, sg=0, workers=8, window=get_window(corpus), min_count=1, vector_size=100 )
```

- Vector size= n, permette di decidere la dimensione dei vettori che l'algoritmo andrà a creare.
- window= n, permette di decidere la massima distanza tra la parola corrente e quella predetta all'interno di una frase. Nel mio caso è dato dalla media di ingredienti in una ricetta (9)
- min count=n, utile per capire se una determinata ricetta piace agli utenti meno (non verrà usata nella soluzione finale)
- workers= n, permette di decidere quanti thread del processore usare per addestrare il modello.
- sg= 0, questo crea un modello CBOW

Usiamo TD-IDF per aggregare gli accorpamenti.

Calcolo la media ponderata di tutti gli incorporamenti di parole di ogni documento (lista ingredienti) e ridimensioniamo i vettori usando la sua frequenza inversa del documento (IDF).

IDF appesantirà i termini molto comuni in un corpus.

In questo modo si ha una capacità di separare le ricette migliore. Word2Vec cerca di prevedere le parole in base all'ambiente circostante, Quindi era fondamentale ordinare gli ingredienti in ordine alfabetico

```
#ottiene il corpus con gli ingredienti in ordine alfabetico
def get_and_sort_corpus(data):
    corpus_sorted = []
    for doc in data.ingredients.values:
        doc=doc.split(",")
        doc.sort()
        corpus_sorted.append(doc)
    return corpus_sorted

#calcola la lunghezza media di ogni doc (lista ingredienti ricetta)
def get_mean(corpus):
    lengths = [len(doc) for doc in corpus]
    avg_len = float(sum(lengths)) / len(lengths)
    return round(avg_len)
```

Abbiamo poi utilizzato sia la similarità sia knn per trovare le ricette simili.

6.4 Similarità del coseno

La similarità del coseno :

```
# usa TF-IDF come pesi per ogni incorporamento di parole
tfidf_vec_tr = TfidfEmbeddingVectorizer(model)
tfidf_vec_tr.fit(corpus)
doc_vec = tfidf_vec_tr.transform(corpus)
doc_vec = [doc.reshape(1, -1) for doc in doc_vec]
assert len(doc_vec) == len(corpus)

print('Shape of word-mean doc2vec...')
print(tfidf_vec_tr.transform(corpus).shape)

input = ingredients
input = input.split(",")

input_embedding = tfidf_vec_tr.transform([input])[0].reshape(1, -1)

# ottieni la similarità del coseno e tutti gli embeddings nel documento
cos_sim = map(lambda x: cosine_similarity(input_embedding, x)[0][0], doc_vec)
scores = list(cos_sim)
# prendi le prime N recommendations
recommendations = get_recommendations(N, scores)
return recommendations
```

6.5 KNN

Ho testato il funzionamento anche con KNN.

Il k-nearest neighbors (traducibile come primi k-vicini), abbreviato in K-NN, `e un algoritmo utilizzato nel riconoscimento di pattern per la classificazione di oggetti basandosi sulle caratteristiche degli oggetti vicini a quello considerato.

In questo caso abbiamo usato KNN per trovare le ricette più vicine rispetto agli ingredienti che l'utente ha fornito. Per fare questo, prima ho convertito i nomi degli ingredienti in vettori utilizzando un modello di embedding. Questi vettori hanno una rappresentazione in uno spazio a più dimensioni che rappresenta la somiglianza semantica degli ingredienti.

Successivamente, ho calcolato la distanza tra il vettore di embedding della ricetta dell'utente e il vettore di embedding di ogni ricetta nella mia base di dati utilizzando la distanza Euclidea. Infine, ho selezionato le K ricette più vicine (dove K `e un parametro che possiamo impostare) e ho restituito i nomi di queste ricette come raccomandazioni, usando la funzione get recommendations.

In sintesi, KNN `e stato utilizzato per trovare le ricette più simili a quella fornita dall'utente utilizzando la distanza Euclidea tra i vettori di embedding degli ingredienti.

Codice:

```
tfidf_vec_tr = TfidfEmbeddingVectorizer(model)
tfidf_vec_tr.fit(corpus)
doc_vec = tfidf_vec_tr.transform(corpus)

input = ingredients
input = input.split(",")

input_vec = tfidf_vec_tr.transform([input])

# ottieni KNN tra input embedding e tutti gli embeddings nel documento
knn = NearestNeighbors(n_neighbors=N, algorithm='ball_tree')
knn.fit(doc_vec)
distances, indices = knn.kneighbors(input_vec)
recommendations = get_recommendations(N, indices[0])

return recommendations
```

La riga `knn = NearestNeighbors(n_neighbors=N, algorithm='ball tree')` crea un oggetto di tipo `NearestNeighbors` da `scikit-learn`.

La classe `NearestNeighbors` ha due parametri importanti:

- `n_neighbors`: questo parametro specifica il numero di vicini più vicini che vogliamo trovare per ogni punto del dataset. In questo caso, `n_neighbors=N` significa che stiamo cercando i N vicini più vicini.
- `algorithm`: questo parametro specifica l'algoritmo di ricerca delle vicinanze che si desidera utilizzare. Nel caso specifico, l'algoritmo specificato `e 'ball tree'.

Il modello viene addestrato sulla matrice dei vettori dei documenti con il metodo fit. Infine il metodo neighbors viene utilizzato per ottenere le K (in questo caso, N) vicinanze più vicine tra il vettore di input e tutti i vettori dei documenti.

Metriche valutazione KNN:

	precision	recall	f1-score	support
0	0.43	0.58	0.50	162
1	0.57	0.40	0.47	10
2	0.36	0.35	0.35	40
3	0.83	0.96	0.89	368
4	0.25	0.08	0.12	12
5	0.57	0.51	0.54	68
6	1.00	0.11	0.20	18
7	0.29	0.16	0.21	43
8	0.75	0.72	0.73	291
9	0.21	0.16	0.18	19
10	0.66	0.56	0.60	158
11	0.00	0.00	0.00	19
accuracy			0.67	1208
macro avg	0.49	0.38	0.40	1208
weighted avg	0.66	0.67	0.65	1208

KNN: input: "farina,uova,grana padano,pangrattato"

```

---Gnocchi alla parigina
latte,uova, farina,parmigiano reggiano,burro,noce moscata,sale,groviera,pepe

---Bombe di patate veloci
patate, farina, lievito,uova,sale,prosciutto cotto,scamorza,origano,semi

---Crepe alla crema di asparagi
latte, farina,uova,sale,burro,noce moscata,asparagi,groviera,pepe,parmigiano reggiano

---Sformati di broccoli con salsa al caprino
broccoli,pepe,sale,uova,grana padano,burro,latte, farina,noce moscata, caprino

---Funghi fritti con maionese di barbabietola
funghi,pangrattato, farina,timo,uova,latte,sale,barbabietole,semi,limone,zenzero

```

Per cui non è possibile valutare tramite precision, recall, ecc. Mi sono basato su Score

COS-SIM: input "farina,uova,grana padano,pangrattato"

```
--Ravioli ricotta e spinaci
farina,uova,semola,spinaci,ricotta,parmigiano reggiano,noce moscata,sale,pepe
Score: 1.000
--Patate duchessa
patate,burro,grana padano,uova,sale,pepe
Score: 1.000
--Biscotti salati senza glutine
farina,burro,grana padano,uova,olive,rosmarino,sale
Score: 0.999
--Frittelle di zucchine croccanti
zucchine,farina,birra,acqua,grana padano,sale,pepe,semi
Score: 0.999
--Uova gratinate
uova,pancetta,burro,farina,latte,parmigiano reggiano,noce moscata,sale,pepe,prezzemolo
Score: 0.999
```


Chapter 7

Quarta Soluzione : Clustering

Ho ripreso la soluzione che usava sono TD-IDF e cosine similarity e abbiamo usato la cosine similarity come input per un algoritmo di clustering .Ed il risultato `e stato ottimo.

Per il clustering ho scelto l'algoritmo K-Means uno dei più diffusi. Attuando due strategie:

Nella prima strategia ho calcolato la cosine similarity di tra tutte le ricette del dataset a priori e l'abbiamo passata in input al clustering.

Nel seconda strategia ho creato il clustering una volta calcolato la cosine similarity tra la ricetta scelta e tutte le altre ricette del dataset.

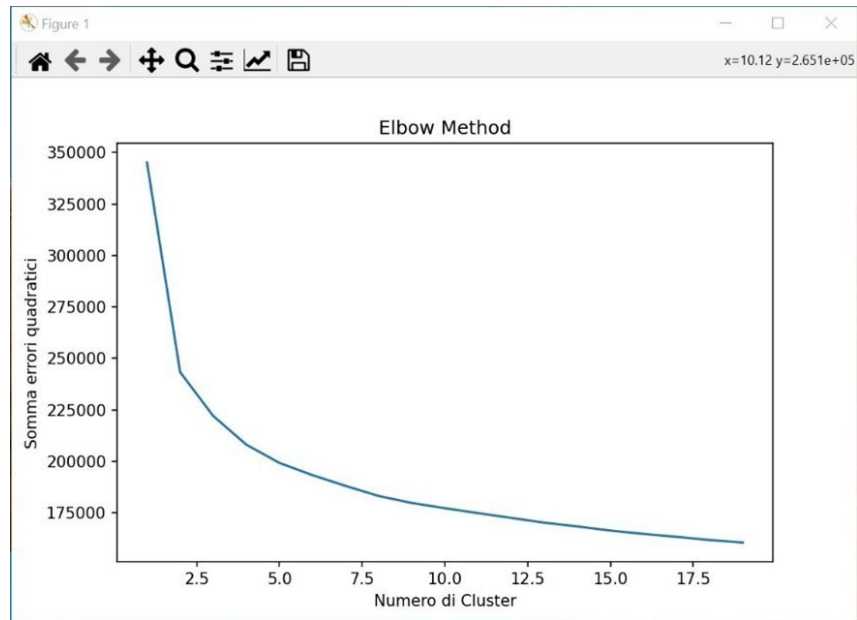
Per la scelta del K abbiamo usato in entrambe le strategie il metodo dell'elbow point per la scelta del k, e Il Silhouette Coefficient per la bontà dei cluster.

L'elbow point `e una metrica nell'analisi del cluster per determinare il numero ottimale di cluster da utilizzare in un modello di clustering. Il punto gomito rappresenta il numero di cluster in cui l'incremento della somma delle distanze quadratiche medie (SSE, Sum of Squared Errors) tra i punti del cluster e il centroide del cluster stesso, inizia a rallentare. Questo punto viene identificato come il punto in cui l'aggiunta di Ulteriori cluster non apporta un significativo miglioramento del modello.In altre parole, l'elbow point rappresenta la trade-off tra il numero di cluster e la qualità del modello di clustering

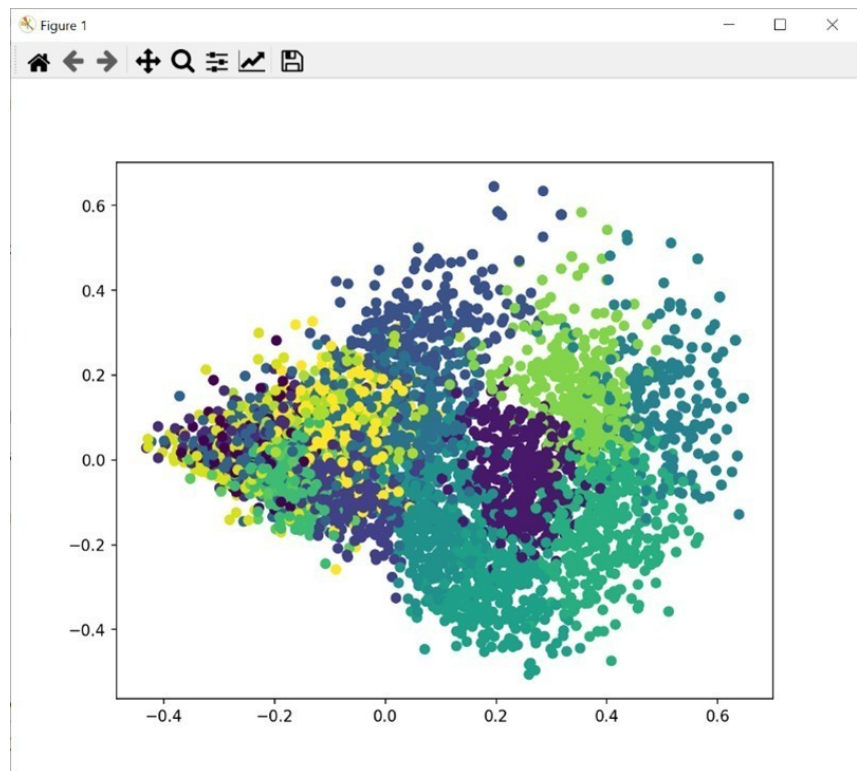
Silhouette Coefficient `e una metrica che valuta la qualità di un clustering rispetto a come le osservazioni sono state assegnate ai vari cluster.Un valore più alto indica un clustering di qualità superiore, con valori compresi tra -1 e 1.

7.1 Prima Strategia

Elbow Point



Cluster



Output fissato a 10

La ricetta selezionata è: Risotto allo Zafferano
Le ricette simili sono:

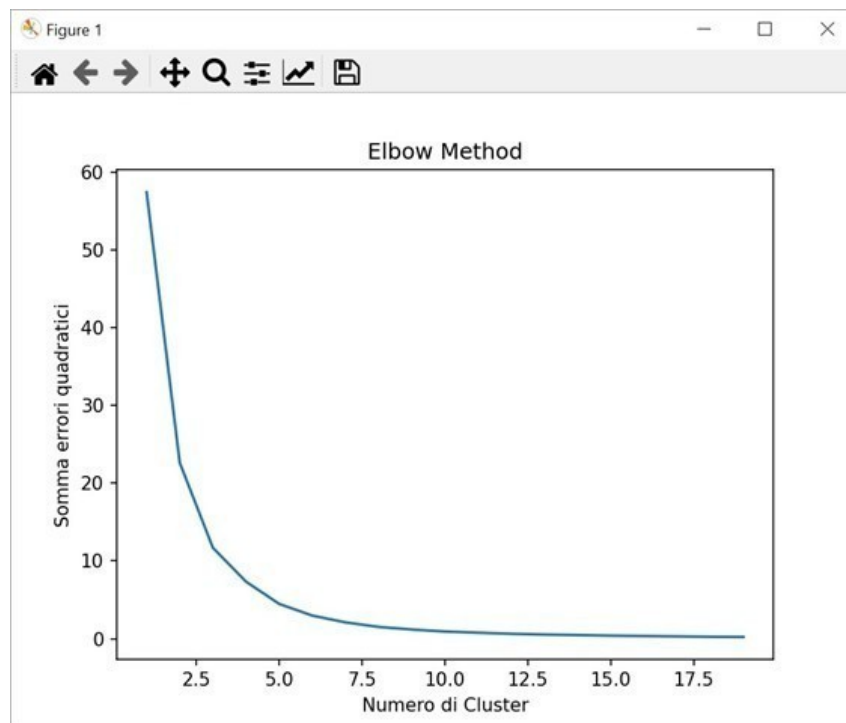
	title	cluster
15	Risotto allo Zafferano	12
30	Risotto alla zucca	12
37	Crema di zucca	12
81	Risotto alla salsiccia	12
82	Vellutata di zucca	12
95	Risotto ai funghi porcini	12
103	Risotto ai funghi	12
108	Pasta con la zucca	12
170	Risotto al pomodoro	12
195	Risotto con zucchini	12

Silhouette Coefficient

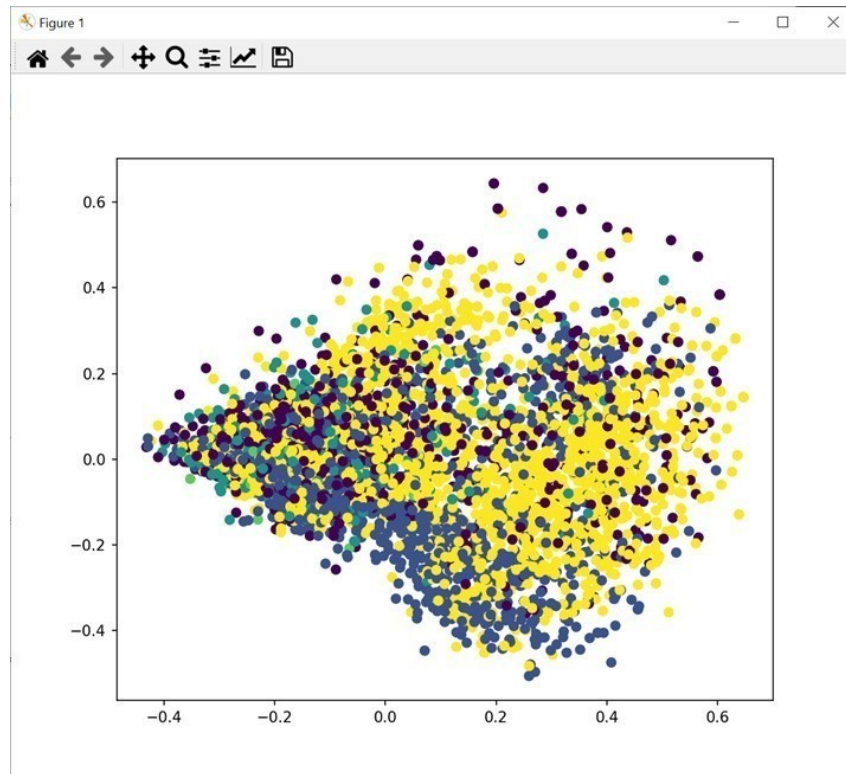
Il Silhouette Coefficient è: 0.06756554372168679

7.2 Seconda Strategia

Elbow Point



Ho visto che già con 5 il risultato era ottimo.
Cluster



Output fissato a 10

```
La ricetta selezionata è: Risotto allo Zafferano
Le ricette simili sono:
```

	title	cluster
15	Risotto allo Zafferano	3
30	Risotto alla zucca	3
31	Riso alla cantonese	3
70	Arancini di riso	3
81	Risotto alla salsiccia	3
95	Risotto ai funghi porcini	3
103	Risotto ai funghi	3
170	Risotto al pomodoro	3
195	Risotto con zucchine	3
224	Risotto con gorgonzola, pere e noci	3

Silhouette Coefficient

Il Silhouette Coefficient è: 0.6863896368535756

Con la prima possiamo effettuare il clustering una sola volta, salvare nel dataset per ogni ricetta il suo cluster di appartenenza, quindi azzerare quasi del tutto il tempo di attesa del risultato. Con la seconda per`o abbiamo raccomandazioni pi`u precise, vedi silhouette Coefficient

Chapter 8

App

Per creare l'interfaccia grafica dell'applicazione ho usato la libreria customtkinter di python che permette di creare una gui semplice e minimale. Inoltre `e possibile non solo inserire una serie di ingredienti che ho , anche selezionare a video delle ricette che ci piacciono per farci consigliare altre ricette simili.

8.1 GUI

Nella pagina iniziale vengono mostrate una serie di ricette che l'utente può selezionare, per ognuna delle quali il programma restituisce 4 ricette consigliate. L'utente pu`o anche scrivere gli ingredienti che ha disposizione nel campo di testo o dirli a voce selezionando il bottone "mic".

Inoltre può selezionare con quale strategia effettuare la ricerca (Word2Vec, Clustering 1, Clustering 2)

Una volta che ha selezionato una ricetta e/o inserito degli ingredienti pu`o procedere con il premere "Avvia Ricerca". Se non vede nessuna ricetta che gli piace può generarne di nuove con il bottone "Nuove ricette", oppure dopo che ha ottenuto dei consigli su una ricetta può tornare alla stessa lista di ricette con il bottone "Lista ricette" per selezionarne una che aveva dimenticato.

8.2 Input Vocale

```
def get_audio_input():  
  
    r = sr.Recognizer()  
  
    with sr.Microphone() as source:  
        print("Dimmi gli ingredienti che hai ")  
        audio = r.listen(source)  
    ingredients = r.recognize_google(audio, language ="it-IT")  
    ingredients_list = ingredients.split()  
    ingredients_string=",".join(ingredients_list)  
    input_text.delete(0,tk.END)  
    input_text.insert(0,ingredients_string)
```

Questo codice `e una funzione chiamata "get audio input" che utilizza la libreria "SpeechRecognition" (sr) per riconoscere la voce dell'utente e tradurla in testo. La funzione utilizza un microfono come fonte di input audio. Quando l'utente inizia a parlare, il codice registra l'audio e lo passa al metodo "recognize google" dell'oggetto "Recognizer" per ottenere la traduzione in testo. Il risultato viene salvato come stringa in "ingredients". Successivamente, la stringa viene divisa in una lista di parole utilizzando il metodo "split" e la lista viene successivamente trasformata in una stringa separata da virgole utilizzando il metodo "join". Infine, il testo viene inserito nella casella di testo visualizzata sull'interfaccia utente.

Chapter 9

Riferimenti

- Word2Vec :[link](#), [link](#) , [link](#).
- NLP: [link](#).
- Web Scraping[link](#).
- Libreria python [link](#).