



UNIVERSITÀ POLITECNICA DELLE
MARCHE

DIGITAL ADAPTIVE CIRCUITS AND LEARNING
SYSTEMS

Sound Event Detection con la tecnica del “*few-shot learning*”

Matteo Orlandini e Jacopo Pagliuca

Prof. Stefano SQUARTINI
Dott.ssa Michela CANTARINI

15 luglio 2021

Indice

1	Introduzione	1
2	Convolutional Neural Network (CNN)	2
3	Area under precision-recall curve (AUPRC)	2
4	Few-shot learning	3
4.1	Introduzione	3
4.2	The meta learning framework	4
4.3	Approcci al meta-apprendimento	4
5	Prototypical Network	5
6	Relation Network	7
7	Dataset Spoken Wikipedia Corpora	9
7.1	Annotazioni	9
7.2	Lettori	11
7.3	Parole	13
8	Protonet codice	17
8.1	Protonet training	17
8.2	Protonet validation	18
8.3	Protonet test	18

1 Introduzione

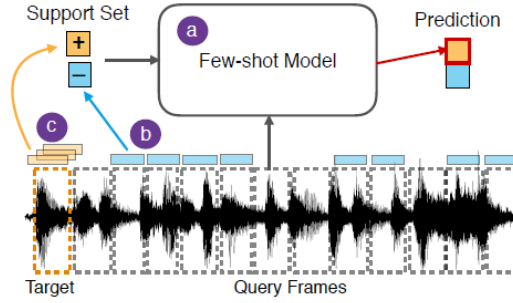


Figura 1: Metodo proposto per il few-shot sound event detection. (a) Applicazione del modello few-shot, (b) costruzione del set di esempi negativi, in blu, e (c) data augmentation per la generazione di più esempi positivi, in arancione. [4]

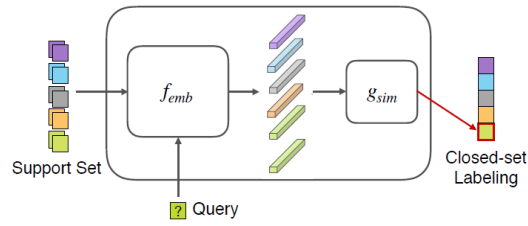


Figura 2: Modello few-shot learning nel caso 5-way 2-shot. [4]

2 Convolutional Neural Network (CNN)

Una rete neurale convoluzionale è una rete feedforward pensata per applicazioni che richiedono l'elaborazione di immagini o di grandi moli di dati. L'ingresso della rete è costituito da una o più matrici che attraversano delle fasi di convoluzione che ne riducono le dimensioni. I filtri di convoluzione sono detti kernel o di pooling. Infine, il risultato viene passato a una rete fully connected che ha il compito di classificazione.

La prima parte del filtro consiste nella convoluzione dell'immagine in ingresso con una serie di kernel i cui parametri vengono allenati. Il filtro viene traslato in entrambe le dimensioni della matrice producendo una serie di feature bidimensionali. Il parametro *stride* definisce il passo con cui il kernel viene traslato sulla matrice di ingresso. Solitamente, viene effettuato uno zero-padding sui contorni del volume in modo da poter caratterizzare anche i valori che si trovano ai bordi delle matrici. In seguito alla convoluzione viene applicata una funzione non-lineare come ad esempio una sigmoide o una ReLu con lo scopo di aumentare la proprietà di non linearità.

Il pooling layer è uno strato della rete che ha lo scopo di ridurre le dimensioni della matrice prodotta dal convolutional layer. Combinando gruppi di elementi della matrice (solitamente 2x2) restituisce il valore massimo tra essi, che andrà a sostituire il blocco stesso.

3 Area under precision-recall curve (AUPRC)

Nella fase di test, la rete elabora una predizione dell'output a partire da un ingresso noto che poi viene confrontata con il valore effettivo. Nel nostro caso è presente un positive set e un negative set, si tratta quindi di classificazione binaria. Il confronto tra predizione e label può produrre quattro risultati:

- Vero Negativo (TN): il valore reale è negativo e il valore predetto è negativo;
- Vero Positivo (TP): il valore reale è positivo e il valore predetto è positivo;
- Falso Negativo (FN): il valore reale è positivo e il valore predetto è negativo;
- Falso Positivo (FP): il valore reale è negativo e il valore predetto è positivo;

Questi valori vanno a comporre la confusion matrix. Nel nostro caso siamo interessati a Precision e Recall.

Il parametro Precision rappresenta quanti tra i casi predetti come positivi sono realmente positivi.

Il parametro Recall indica quanti tra i casi realmente positivi è stato predetto in modo corretto.

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

4 Few-shot learning

(<https://www.borealisai.com/en/blog/tutorial-2-few-shot-learning-and-meta-learning-i/>)

4.1 Introduzione

Gli esseri umani possono riconoscere nuove classi di oggetti partendo da pochissimi esempi. Tuttavia, la maggior parte delle tecniche di machine learning richiedono migliaia di esempi per ottenere prestazioni simili a quelle umane. L'obiettivo del *few-shot learning* è classificare i nuovi dati dopo aver visto solo pochi esempi di training. Nel caso estremo, potrebbe esserci solo un singolo esempio per ogni classe (*one shot learning*). In pratica, il few-shot learning è utile quando è difficile trovare esempi di training (ad es. casi di una malattia rara) o quando il costo dell'etichettatura dei dati è elevato.

L'apprendimento few-shot viene solitamente studiato utilizzando la classificazione *N-way-K-shot*. L'obiettivo è quello di discriminare le N classi composte da K esempi ciascuna. Una tipica dimensione del problema potrebbe essere quella di discriminare tra $N = 10$ classi con solo $K = 5$ campioni ciascuno di training. Non possiamo allenare un classificatore usando metodi convenzionali; qualsiasi algoritmo di classificazione moderno dipenderà da molti più parametri rispetto agli esempi di addestramento e generalizzerà male.

Se i dati non sono sufficienti per ridurre il problema, una possibile soluzione è acquisire esperienza da altri problemi simili. A tal fine, la maggior parte degli approcci caratterizza l'apprendimento a breve termine con un problema di meta-apprendimento.

4.2 The meta learning framework

Nel framework di apprendimento classico, impariamo come classificare dai dati di training e valutiamo i risultati utilizzando i dati di test. Nel quadro del meta-apprendimento, *impariamo come imparare* a classificare in base a una serie di *training task* e valutiamo utilizzando una serie di test task; In altre parole, usiamo un insieme di problemi di classificazione per aiutare a risolvere altri insiemi non correlati.

Qui, ogni attività imita lo scenario few-shot, quindi per la classificazione N-way-K-shot, ogni attività include N classi con K esempi ciascuna. Questi sono noti come support set per il task e vengono utilizzati per apprendere come risolvere il task. Inoltre, esistono ulteriori esempi delle stesse classi, note come set di query, utilizzate per valutare le prestazioni del task corrente. Ogni task può essere completamente unico; potremmo non vedere mai le classi di un task in nessuno degli altri. L'idea è che il sistema veda ripetutamente istanze (task) durante l'addestramento che corrispondono alla struttura dell'attività finale di few-shot, ma che contengono classi diverse.

Ad ogni fase del meta-apprendimento, aggiorniamo i parametri del modello in base ad un training task selezionato casualmente. La funzione di loss è determinata dalle prestazioni di classificazione sul set di query del task, in base alla conoscenza acquisita dal relativo set di supporto. Poiché la rete viene sottoposta ad un compito diverso in ogni fase temporale, deve imparare a discriminare le classi di dati in generale, piuttosto che un particolare sottoinsieme di classi.

Per valutare le prestazioni in pochi colpi, utilizziamo una serie di attività di test. Ciascuno contiene solo classi invisibili che non erano in nessuna delle attività di formazione. Per ciascuno, misuriamo le prestazioni sul set di query in base alla conoscenza del loro set di supporto.

4.3 Approcci al meta-apprendimento

Gli approcci al meta-apprendimento sono diversi e non c'è consenso sull'approccio migliore. Tuttavia, esistono tre famiglie distinte, ognuna delle quali sfrutta un diverso tipo di conoscenza a priori:

- Conoscenze a priori sulla somiglianza: apprendiamo degli embedding nei training task che tendono a separare classi diverse anche quando non sono visibili.
- Conoscenze a priori sull'apprendimento: utilizziamo le conoscenze a priori per vincolare l'algoritmo di apprendimento a scegliere parametri che generalizzino bene da pochi esempi.

- Conoscenza a priori dei dati: sfruttiamo le conoscenze a priori sulla struttura e la variabilità dei dati e questo ci consente di apprendere modelli praticabili da pochi esempi.

5 Prototypical Network

L'approccio si basa sull'idea che esiste un embedding in cui i punti delle istanze di una classe si raggruppano attorno a una singola rappresentazione prototipo per ogni classe. Nella classificazione few-shot per ogni episodio viene fornito un support set di N esempi etichettati $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ dove $\mathbf{x}_i \in \mathbb{R}^D$ rappresenta il vettore D -dimensionale della feature (nel nostro caso spettrogrammi di dimensione 128×51) e $y_i \in \{1, \dots, K\}$ la rispettiva label. S_k denota il set di esempi etichettati con la classe k .

La rete Prototypical calcola una rappresentazione M -dimensionale \mathbf{c}_k , o *prototipo*, di ogni classe tramite una funzione di embedding $f_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$ con parametri da allenare ϕ . La funzione di embedding è rappresentata da una rete convoluzionale. Ogni prototipo è il vettore media tra gli embedding delle istanze della stessa classe.

$$\mathbf{c}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}_i, y_i) \in S_k} f_\phi(\mathbf{x}_i) \quad (3)$$

Data una funzione distanza $d : \mathbb{R}^M \times \mathbb{R}^M \rightarrow [0, +\infty)$, la rete Prototypical calcola la relazione di una query \mathbf{x} rispetto ai prototipi tramite la funzione softmax delle distanze prese con segno negativo.

$$p_\phi(y = k|\mathbf{x}) = \frac{\exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k))}{\sum_k' \exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k'))} \quad (4)$$

Il processo di training avviene minimizzando il negativo del logaritmo della probabilità $J(\phi) = -\log(p_\phi(y = k|\mathbf{x}))$ considerando la distanza fra query e il prototipo della sua classe. Gli episodi di training sono formati campionando C classi di parole da un lettore e selezionando per ognuna casualmente K istanze e Q query.

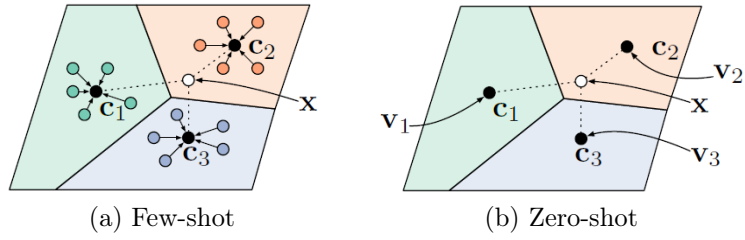


Figura 3: Reti Prototypical negli scenari few-shot e zero-shot. (a): i prototipi few-shot \mathbf{c}_k sono calcolati come la media degli embedding del support set per ogni classe. (b): i prototipi zero-shot \mathbf{c}_k sono prodotti facendo l'embedding dei metadati \mathbf{v}_k . In entrambi i casi i punti degli embedding delle query sono classificati facendo il softmax sulle distanze del prototipo delle classi: $p_\phi(y = k|\mathbf{x}) \propto \exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k))$. [2]

6 Relation Network

La rete Relation ha lo scopo di associare due istanze alla volta per determinare la loro similarità. Questo viene effettuato concatenando gli embedding di più istanze in un unico elemento che sarà dato in ingresso a una rete decisionale i cui parametri saranno aggiornati in modo che una concatenazione di elementi simili restituisca un risultato vicino a 1. La Relation Network è costituita da due moduli: un modulo di *embedding* f_ϕ (equivalente a quello nella Prototypical) e un modulo di *relation* g_ϕ , come illustrato in figura 4. Le istanze x_i del query set \mathcal{Q} e quelle x_j del support set \mathcal{S} vengono date in ingresso al modulo di embedding producendo dei vettori (feature maps) $f_\phi(x_i)$ e $f_\phi(x_j)$. Questi ultimi vengono poi dati all'operatore $\mathcal{C}(\cdot, \cdot)$ che ne fa la concatenazione: $\mathcal{C}(f_\phi(x_i), f_\phi(x_j))$. Le feature map concatenate passano poi attraverso il modulo di decisione che restituisce uno scalare da 0 a 1, il quale rappresenta la somiglianza tra x_i e x_j . Per il caso C -way one-shot, viene concatenata la query con le istanze delle C classi producendo C punteggi di somiglianza.

$$r_{i,j} = g_\phi(\mathcal{C}(f_\phi(x_i), f_\phi(x_j))), \quad i = 1, 2, \dots, C \quad (5)$$

Nel caso C -way K-shot invece, la query viene concatenata con la somma elemento per elemento degli embedding di ogni istanza delle classi. Quindi, in entrambi i casi i confronti $r_{i,j}$ sono C per ogni query.

Per allenare il modello viene usato l'errore quadratico medio (MSE) in modo che l'uscita del modulo di decisione produca 1 se i vettori concatenati sono della stessa classe e 0 altrimenti.

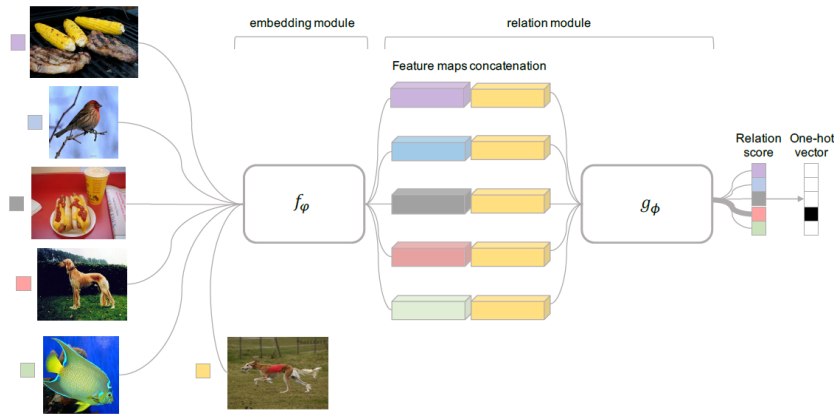


Figura 4: Architettura della Relation Network nel caso 5-way 1-shot con un esempio di query. [3]

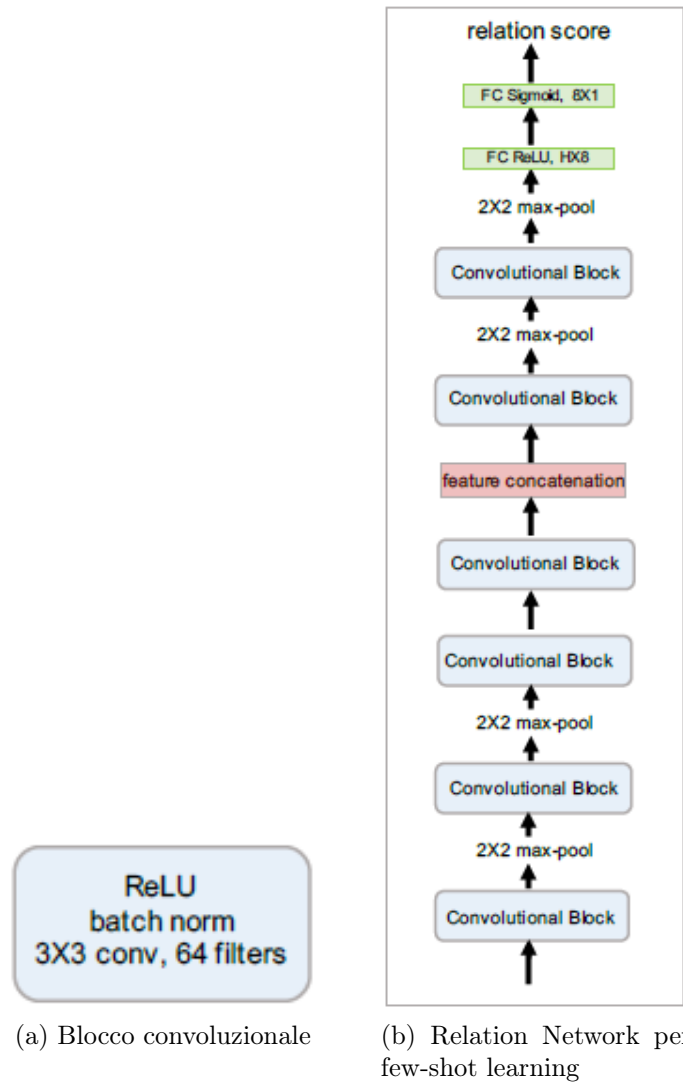


Figura 5: Architettura della Relation Network per few-shot learning (b) composta dagli elementi inclusi nel blocco convoluzionale (a). [3]

7 Dataset Spoken Wikipedia Corpora

Il progetto Spoken Wikipedia unisce lettori volontari di articoli di Wikipedia. Sono disponibili centinaia di articoli in inglese, tedesco e olandese per gli utenti che non sono in grado o non vogliono leggere la versione scritta dell'articolo. Il dataset trasforma i file audio in un corpus, cioè una raccolta ordinata e completa di opere o di autori, allineato nel tempo, rendendolo accessibile per la ricerca.[1]

Questo corpus ha diverse caratteristiche importanti:

- centinaia di ore di audio allineato
- lettori eterogenei
- diversi argomenti
- genere testuale ben studiato
- le annotazioni delle parole possono essere mappate all'html originale
- allineamenti a livello di fonema

Ogni articolo è suddiviso in sezioni, frasi e token. Ogni token è normalizzato e la normalizzazione è allineata all'audio.

7.1 Annotazioni

Ogni annotazione è racchiusa in un tag chiamato “**article**”, che contiene una sezione “**meta**” per i metadati, come si può vedere in 1 e una sezione “**d**”, come in 2, contenente l'articolo e le relative annotazioni.

```
<article>
<meta>
  <link key="DC.conformsto" value="http://nats.gitlab.io/swc/schema/swc-1.0.rnc"/>
  <prop key="DC.creator" value="Spoken Wikipedia Corpus Collection Software"/>
  <prop key="DC.publisher" value="Universität Hamburg"/>
  <link key="DC.reference" value="http://nbn-resolving.de/urn:nbn:de:gbv:18-228-7-2209"/>
  <prop key="DC.type" value="dataset"/>
  <prop key="DC.license" value="CC-BY-SA"/>
  <prop key="DC.title" value="Limerence"/>
  <prop key="DC.language" value="en"/>
  <prop key="DC.identifier" value="Limerence"/>
  <prop key="DC.date.read" value="2005-04-29 00:00:00"/>
  <link key="DC.source" value="https://en.wikipedia.org/wiki/Limerence"/>
  <prop key="DC.source.wikiID" value="154147"/>
  <prop key="DC.source.revision" value="791627969"/>
  <link key="DC.source.text" value="https://en.wikipedia.org/w/index.php?title=Limerence&oldid=13811989"/>
  <link key="DC.source.audio" value="https://upload.wikimedia.org/wikipedia/commons/aa/Limerence.ogg" group="audiol"/>
  <prop key="DC.source.audio.offset" value="0.0" group="audiol"/>
  <link key="DC.source.audio.page" value="https://en.wikipedia.org/w/index.php?title=File%3aLimerence.ogg" group="audiol"/>
  <link key="DC.source.audio.date" value="2007-09-14 19:23:05" group="audiol"/>
  <prop key="reader.name" value="the Epot"/>
  <prop key="processing.step" value="tokenize" group="tokenize"/>
</article>
```

```

<prop key="processing.step.date" value="2017-08-10T07:44:44.095+02:00[Europe/Berlin]"
  group="tokenize"/>
<prop key="processing.step.options" value="Namespace(output=articles/Limerence/
  tokenized.swc, all_sections=false, null_normalize=false, raw_output=null,
  subparser_name=tokenize, lang=en, no_introduction=false, article_dir=articles/
  Limerence)" group="tokenize"/>
<prop key="processing.step.git.commit.id" value="7431
  dbf93f212ad828208abaf8f518fb8de11ff3" group="tokenize"/>
<prop key="processing.step.git.commit.time" value="09.08.2017 @ 15:21:55 CEST" group
  ="tokenize"/>
<prop key="processing.step" value="align" group="align"/>
<prop key="processing.step.date" value="2017-08-11T16:12:18.423+02:00[Europe/Berlin]"
  group="align"/>
<prop key="processing.step.options" value="Namespace(output=articles/Limerence/
  aligned.swc, transcript=articles/Limerence/tokenized.swc, g2p=../model_en/model.
  fst.ser, phone=false, subparser_name=align, dict=../model_en/empty.dic,
  acoustic_model=../model_en/, audio=articles/Limerence/audio.wav)" group="align"
  />
<prop key="processing.step.git.commit.id" value="7431
  dbf93f212ad828208abaf8f518fb8de11ff3" group="align"/>
<prop key="processing.step.git.commit.time" value="09.08.2017 @ 15:21:55 CEST" group
  ="align"/>
</meta>

```

Codice 1: Metadati delle annotazioni delle parole in un audio

Il documento, contrassegnato dal tag “d”, può contenere parti diverse, ciascuna spiegata di seguito. La sezione “extra” contiene il testo che abbiamo incluso ma non fa parte dell’articolo, “ignored” contiene ciò che fa parte del testo ma viene ignorato per l’allineamento, “section” contiene un titolo e un contenuto, “p” contiene un paragrafo e “s” una frase che a sua volta contiene dei token “t”. In quest’ultimo è contenuta la singola parola originale e le normalizzazioni. Per esempio, la punteggiatura non ha annotazioni di normalizzazione in quanto non è pronunciata, ma il numero 500 ne ha due - “cinque” e “cento”. un token stesso non ha allineamento, solo la sua normalizzazione “n” è allineata. La normalizzazione ha una “pronunciation” e può avere un tempo di “start” ed “end”, se è allineata. La normalizzazione, a sua volta, può contenere dei fonemi “ph”.

```

<d>
  <extra text="LimerenceFrom wikipedia, the free encyclopedia at e n dot wikipedia dot
    org.">
    <s text="LimerenceFrom wikipedia, the free encyclopedia at e n dot wikipedia dot org
      .">
      <t text="LimerenceFrom">
        <n pronunciation="LimerenceFrom" start="140" end="1190"/>
      </t>
      <t text="wikipedia">
        <n pronunciation="wikipedia" start="1250" end="1950"/>
      </t>
      <t text=","/>
      <t text="the">
        <n pronunciation="the" start="1950" end="2070"/>
      </t>
      <t text="free">
        <n pronunciation="free" start="2070" end="2300"/>
      </t>
      <t text="encyclopedia">
        <n pronunciation="encyclopedia" start="2300" end="3220"/>
      </t>
      <t text="at">
        <n pronunciation="at"/>
      </t>
      <t text="e">
        <n pronunciation="e" start="3490" end="3710"/>
      </t>
    </s>
  </extra>
</d>

```

Codice 2: Annotazioni delle parole in un audio

7.2 Lettori

Il dataset Spoken Wikipedia Corpora contiene un totale di 1340 audio di diversi lettori, ma nel progetto vengono presi solo gli audio che contengono annotazioni a livello di parola. In questo modo vengono presi solo 208 lettori e partizionati come in [4] con un rapporto 138 : 15 : 30 tra lettori di training, validation e test. I lettori e le parole sono state estratte dai file “aligned.swc” contenuti in ogni audio e, successivamente, salvati in diversi file json. Per la gestione dei json sono state create due funzioni utili per la lettura e scrittura come mostrato nel codice 3.

```
1 import json
2
3 def write_json_file(filename, list_of_dict, indent = 0):
4     f = open(filename, "w")
5     f.write(json.dumps(list_of_dict, indent = indent))
6     f.close
7
8 def read_json_file(filename):
9     f = open(filename, "r")
10    list_of_dict = json.load(f)
11    f.close
12    return list_of_dict
```

Codice 3: json_manager.py

Il codice 4 salva un file json chiamato “readers_paths.json” formato dalle chiavi “reader_name” e “folder” i cui valori sono rispettivamente il nome del lettore e le cartelle in cui sono salvati i file audio registrati dal relativo lettore, come si può vedere nel json 5. Questo codice, usando la libreria `xml.etree.ElementTree` e la funzione `ET.parse`, rappresenta l'intero documento XML come un albero. La funzione `getroot` ne trova la radice e, successivamente, si scorre l'albero iterando finché non si trova il tag “prop” in cui è contenuta la chiave “reader.name”. Inoltre, si effettua un controllo sul nome del lettore perché può capitare che viene salvato nel file “aligned.swc” con nomi diversi. Ad esempio, in alcuni file si può trovare nella chiave “reader.name” il valore “:en:user:alexkillby|alexkillby”, mentre in altri solo “alexkillby” oppure si può trovare “user:popularoutcast”, mentre in altri solo “popularoutcast”. Una volta noto il nome del lettore si vede se è già presente nella lista di dizionari e se è un lettore nuovo lo si aggiunge con il relativo nome del file audio, mentre se il lettore era già presente si aggiunge solamente il nome del file audio.

```
1 import xml.etree.ElementTree as ET
2 from tqdm import tqdm
3 import os
4 import json
5 from json_manager import *
6
7 # Initialize list with empty dictionaries
```

```

8 readers = []
9
10 source_path = "./Dataset/English spoken wikipedia/english/"
11 filename = "aligned.swc"
12
13 # search for readers for each folder
14 for audio_path in os.scandir(source_path):
15     # save only folder name from entire path
16     folder = os.path.basename(audio_path)
17     if (os.path.exists(source_path + "/" + folder + "/" + filename)):
18         # parse the xml file "aligned.swc"
19         tree = ET.parse(source_path + "/" + folder + "/" + filename)
20         # getroot returns the root element for this tree
21         root = tree.getroot()
22         #root.iter creates a tree iterator with the current element as the root.
23         # The iterator iterates over this element and all elements below it, in
24         # document (depth first) order.
25         for property in root.iter(tag = 'prop'):
26             # if the key "reader.name" exists
27             if (property.attrib['key'] == 'reader.name'):
28                 # save the reader name taking the value of the attribute
29                 reader_name = property.attrib['value'].lower()
30                 # fix readers names that contain "user:"
31                 if ("user:" in reader_name):
32                     # fix readers names that contain "|"
33                     if ("|" in reader_name):
34                         # example reader_name = [[:en:user:alexkillby|alexkillby]] ->
35                         # reader_name = alexkillby
36                         reader_name = reader_name[reader_name.find("user:") + 5:
37                                                 reader_name.find("|")]
38                     # fix readers names that contain "|"
39                     elif ("]]" in reader_name):
40                         # example reader_name = [[user:popularoutcast]] ->
41                         # reader_name = popularoutcast
42                         reader_name = reader_name[reader_name.find("user:") + 5:
43                                                 reader_name.find("]]")]
44                 # if the reader is not yet on the list create a dict and append to
45                 # the readers list
46                 if not any(reader['reader_name'] == reader_name for reader in
47                             readers):
48                     dictionary = {'reader_name': reader_name, 'folder': [folder]}
49                     readers.append(dictionary)
50                 else:
51                     # if the reader is already on the list add the folder name
52                     for reader in readers:
53                         if (reader['reader_name'] == reader_name):
54                             reader['folder'].append(folder)
55 # print the number of the readers
56 print("The readers are:", str(len(readers)))
57
58 # save a "readers_paths.json" with the name of the readers and the relative
59 # file audio folders
60 write_json_file("readers_paths.json", readers, indent = 4)

```

Codice 4: xml_parser_readers.py

```

{
  "reader_name": "the epopt",
  "folder": [
    "(I_Can%27t_Get_No)_Satisfaction",
    "Ceremonial_ship_launching",

```

```

    "Limerence",
    "Revolt_of_the_Admirals",
    "Ship_commissioning"
  ]
},
{
  "reader_name": "wodup",
  "folder": [
    "0.999..%2e",
    "Execution_by_elephant",
    "Hell_Is_Other_Robots",
    "Tom_Bosley",
    "Truthiness"
  ]
},

```

Codice 5: Formato del file readers_paths.json

7.3 Parole

Il codice 6.....

```

1  import xml.etree.ElementTree as ET
2  from tqdm import tqdm
3  import os
4  import collections
5  from json_manager import *
6  import random
7
8  folder = []
9  source_path = "./Dataset/English spoken wikipedia/english/"
10 filename = "aligned.swc"
11 file_audio_name = "audio.ogg"
12
13 # iterate in each folder of the dataset
14 for folder in tqdm(os.scandir(source_path), desc = "Folder number"):
15     # initialize the list of dict "target_words"
16     target_words = []
17     if (os.path.exists(folder.path + "/" + filename) \
18         and os.path.exists(folder.path + "/" + file_audio_name) \
19         and os.path.isdir(folder.path)):
20         # parse the xml file aligned.swc
21         tree = ET.parse(folder.path + "/" + filename)
22         # getroot returns the root element for this tree
23         root = tree.getroot()
24         # initialize the list "words"
25         words = []
26         # root.iter creates a tree iterator with the current element as the
27         # root. The iterator iterates over this element and # all elements
28         # below it, in document (depth first) order.
29         for token_normalization in root.iter(tag = 'n'):
30             # we take only the words with the timestamp, so only if there is
31             # 'start' (or 'end') tag
32             if 'start' in token_normalization.keys():
33                 # add every word with "start" key
34                 words.append(token_normalization.attrib['pronunciation'].
35                             lower())
36
37     # collections.Counter stores elements as dictionary keys, and their
38     # counts are stored as dictionary values.

```

```

37 unique_words = collections.Counter(words)
38 # for each key (word) in "unique_words" append a new target_word if
39 # the number of occurrence is at least 10
40 for key in unique_words.keys():
41     # we only consider words that occur at least 10 times in the
42     # recording. Note that unique_words[key] is the word frequency
43     if (unique_words[key] >= 10):
44         # add a new target word
45         target_words.append({'word' : key, \
46                             'frequency' : unique_words[key], \
47                             'start' : [], \
48                             'end' : []})
49 # for each "target_words" append the relative "start" and "end"
50 # timestamp
51 for token_normalization in root.iter(tag = 'n'):
52     # we take only the words with the timestamp, so only if there is
53     # 'start' (or 'end') tag
54     if 'start' in token_normalization.keys():
55         # iterate over the "target_words"
56         for target_word in target_words:
57             # add start and end timestamp only to the relative
58             # "target_word"
59             if target_word['word'] == token_normalization.attrib['
60                 pronunciation'].lower():
61                 target_word['start'].append(int(token_normalization.
62                 attrib['start']))
63                 target_word['end'].append(int(token_normalization.
64                 attrib['end']))
65
66 write_json_file(folder.path+"/word_count.json", target_words, indent
67                = 4)
68
69 # If there are more than 10 words that occur at least 10 times in
70 # the recording, we sort the words by their number of occurrences,
71 # divide the sorted list into 10 equally sized bins, and sample
72 # one keyword per bin.
73 if (len(target_words) >= 10):
74     target_words = random.sample(target_words, 10)
75 # save the "target_words.json"
76 write_json_file(folder.path+"/target_words.json", target_words,
77                indent = 4)

```

Codice 6: find_target_words.py

```

{
    "word": "i",
    "frequency": 12,
    "start": [
        660,
        8800,
        115050,
        ...
    ],
    "end": [
        870,
        8940,
        115240,
        ...
    ]
},
{
    "word": "the",

```



```

    "frequency": 75,
    "start": [
        4160,
        49930,
        53680,
        ...
    ],
    "end": [
        4320,
        50030,
        53710,
        ...
    ]
},

```

Codice 7: Formato del file word_count.json

Il file “readers_paths.json” viene successivamente usato per salvare le parole di ogni lettore, il nome della relativa cartella in cui vengono pronunciate e i timestamp di inizio e fine della parola per ogni cartella come mostrato nel json 9. Il codice 8 legge il nome di ogni lettore dal json precedentemente salvato, successivamente per ognuna delle cartelle del lettore si.....

```

1  from tqdm import tqdm
2  import os
3  from json_manager import *
4
5  source_path = "./Dataset/English spoken wikipedia/english/"
6  # read the json that contains the readers name and their audio folders
7  readers = read_json_file("readers_paths.json")
8  #initialize a list of dict
9  readers_words = []
10 # for each reader search the word spoken by the reader
11 for reader in tqdm(readers):
12     # create a dict with 'reader_name' and 'words' keys
13     new_readers_words = {'reader_name' : reader['reader_name'], \
14                          'words' : [] }
15
16     # for each reader create a new dict for words
17     words_per_reader = []
18     # flag to signal if "word_count.json" exists
19     json_file_exist = False
20     # search for each folder the file "word_count.json"
21     for reader_folder in reader['folder']:
22         if (os.path.exists(source_path + "/" + reader_folder + "/word_count.
23                               json")):
24             # "word_count.json" exists
25             json_file_exist = True
26             # read "word_count.json"
27             recording_words = read_json_file(source_path + "/" +
28                                               reader_folder + "/word_count.json")
29             # for each audio folder create a new list of dict
30             folder_per_word = []
31             # for each word in the audio save the folder, start and end
32             # timestamps
33             for word in recording_words:
34                 # create a dict with 'folder', 'start' and 'end' keys
35                 folder_per_word = {'folder' : reader_folder, \

```

```

34         'start' : word['start'], \
35         'end' : word['end']]
36     # if the word is not yet in the list add the word
37     if not any (word['word'] == word_per_reader['word'] for
38                 word_per_reader in words_per_reader):
39         words_per_reader.append({ 'word' : word['word'], \
40                                   'folders' : [folder_per_word]})
41
42     # otherwise add 'start' and 'end' if the "reader_folder" is
43     # in word_per_reader['folders'] or add the "folder_per_word"
44     else:
45         for word_per_reader in words_per_reader:
46             if word['word'] == word_per_reader['word']:
47                 if reader_folder in word_per_reader['folders']:
48                     for folder in word_per_reader['folders']:
49                         if reader_folder == folder['folder']:
50                             folder['start'] += word['start']
51                             folder['end'] += word['end']
52                 else:
53                     word_per_reader['folders'].append(
54                         folder_per_word)
55
56     # if "word_count.json" exists add the new reader to "training_readers"
57     # list
58     if (json_file_exist):
59         new_readers_words['words'] = words_per_reader
60         readers_words.append(new_readers_words)
61
62     # save a "readers_words.json" with the name of the readers and the relative
63     # words spoken
64     write_json_file("readers_words.json", readers_words, indent = 0)

```

Codice 8: words_per_reader.py

```

"reader_name": "wodup",
  "words": [
    {
      "word": "the",
      "folders": [
        {
          "folder": "0.999..%2e",
          "start": [
            6950,
            1029740,
            1032520,
            ...
          ],
          "end": [
            7190,
            1029880,
            1032620,
            ...
          ]
        },
        {
          "folder": "Execution_by_elephant",
          "start": [
            3600,
            10680,
            ...

```

Codice 9: Formato del file readers_words.json

8 Protonet codice

8.1 Protonet training

La funzione `get_training_validation_readers`, a partire dalla lista di lettori di training/validation, seleziona quelli che hanno almeno un numero di parole diverse pari a `C` e li suddivide tra training e validation seguendo il rapporto usato in [ARTICOLO RIFERIMENTO]: 138/153 e 15/153. Per fare ciò scorre l'elenco delle directory dei lettori e salva su un array le parole a loro associate. Al termine di questa operazione, se il numero di parole è almeno `C` il lettore è ritenuto valido e il suo nome viene aggiunto a una lista. Infine la lista viene scissa come spiegato sopra.

Il primo for rappresenta l'inizio degli episodi. Il numero totale di iterazioni è 60000 come in [ARTICOLO]; per ogni episodio viene richiamata la funzione `batch_sample`. Essa campiona a caso uno dei lettori di training e di questo seleziona casualmente `C` parole. Per ognuna di queste parole si ricava il numero di istanze presenti del dataset. Viene inizializzato un vettore con gli indici di queste e ne vengono campionati $K+Q$ in `index`. Poi, iterando, gli spettrogrammi il cui indice è presente in `index` vengono salvati e divisi in K istanze del support set e Q istanze del query set.

Una volta ottenute, le feature del support set e query set vengono date in ingresso alla funzione `loss`. Queste ultime vengono messe in fila e concatenate per poter essere passate al modello che calcolerà gli embedding per ogni istanza. Gli elementi del support set che fanno parte della stessa classe andranno a costituire i prototipi tramite una media dei loro valori: `prototypes = embeddings[:n_class*n_support].view(n_class, n_support, embeddings_dim).mean(1)`.

Il parametro `target_inds` viene costruito ripetendo l'indice di ogni classe per il numero di query e sarà utilizzato per poter ricordare a che classe appartiene ogni query.

Per ogni elemento del query set vengono calcolate le distanze rispetto ai prototipi di ogni classe tramite la funzione `euclidean_dist` che rappresenta la distanza euclidea. Per ogni query viene poi usata la funzione di attivazione `log_softmax` fra tutte le sue distanze. La funzione loss da minimizzare è costituita dalla media dei `log_softmax` della classe corrispondente alla query presi con segno negativo. Infatti, nel caso ideale in cui la distanza delle query della stessa classe dal prototipo è 0, la media dei softmax delle distanze dal prototipo assegnato sarà nulla. La accuracy invece è determinata dal numero dei casi in cui il prototipo a distanza minima da una query (o il massimo del softmax) corrisponde con il prototipo della classe associata alla query.

Con `loss_out.backward()` viene utilizzata la loss per la retropropagazione che aggiorna i parametri della rete.

8.2 Protonet validation

Ogni 5000 episodi del training vengono effettuati 1000 episodi di validation. Come nel training vengono ricavati support set e query set ma in questo caso dai lettori di validation. Per questi vengono calcolati con la funzione `loss` allo stesso modo loss e accuracy e vengono salvati. In questo caso però i parametri del modello non vengono aggiornati. Questo processo è necessario per verificare l'effettività del modello su ingressi mai visti in fase di training. Al termine dei 1000 episodi di validation, se la media dell'accuracy è migliore di quella calcolata nel batch precedente il modello viene salvato e sovrascritto al precedente, altrimenti viene ignorato. In questo modo si aggira il problema dell'overfitting fungendo da pseudo early stopping.

8.3 Protonet test

Nella fase di test, seguendo l'articolo di riferimento, l'oggetto di interesse non sono i lettori come nel caso di training e validation. Vengono infatti presi in considerazione i singoli audio. Da questi sarà necessario identificare un gruppo di istanze di una parola (il `positive_set`) e, campionando nell'audio, si potrà verificare se la rete è in grado di riconoscere la parola designata.

Per ogni audio quindi viene ricavato il numero di parole con spettrogramma (in base al contenuto di `Test_features`). Per ognuna di esse viene ricavato un `positive set` e un `negative set` con i quali è possibile verificare la rete. Questo processo viene ripetuto 10 volte come in [ARTICOLO] in modo da generalizzare il più possibile il test. Infatti il `positive set` è ricavato campionando casualmente `p` istanze della parola che si sta analizzando, mentre il `negative set` è ottenuto campionando `n` istanze fra tutte le altre parole dello stesso audio. Per fare questo viene usata la funzione `get_negative_positive_query_set`. Tramite un indice che indica le istanze della parola "positiva" vengono campionati `p` spettrogrammi che vanno a costituire il `positive set`. Le restanti istanze della stessa parola vengono inserite invece nel `query set` come esempio di spettrogramma che la rete dovrebbe riconoscere. Tra l'elenco delle parole viene quindi rimossa la parola positiva con `words.remove(pos_word)`. Usando le restanti vengono ricavate `n` esempi di `negative set` campionando a caso una delle parole e una delle istanze di essa. Allo stesso modo viene formata la seconda metà del `query set` che rappresenta gli esempi negativi. Il numero di `query negative` viene preso uguale a quello delle `positive`.

A questo punto, i tre set e il modello sono parametri della funzione `test_loss`. Allo stesso modo di training e validation vengono calcolati gli embedding per ogni istanza e i prototipi delle due classi: positive set e negative set. Vengono calcolate le distanze che sono passate con segno negativo alla funzione `log_softmax`. Tra i risultati viene selezionato il maggiore che rappresenta la predizione della rete: ovvero la classe a cui viene assegnata la query. Queste predizioni sono contenute nel vettore `y_pred_tmp`, mentre in `y_true_tmp` sono presenti le label reali di ogni query.

Prima di essere utilizzati per calcolare precision e recall, questi due vettori vengono sottoposti a una operazione di complemento a 2 che inverte 0 e 1 in quanto la funzione `precision_recall_curve` richiede che con la label 1 sia identificata la classe positiva. Infine, `sklearn.metrics.auc` calcola la AUPRC. Questo calcolo avviene al termine delle operazioni per ogni audio. Una volta terminati gli audio viene calcolata media e varianza degli AUPRC.

Riferimenti bibliografici

- [1] Arne Köhn, Florian Stegen e Timo Baumann. «Mining the Spoken Wikipedia for Speech Data and Beyond». Inglese. In: *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)* (23–28 mag. 2016). A cura di Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk e Stelios Piperidis. Portorož, Slovenia: European Language Resources Association (ELRA). ISBN: 978-2-9517408-9-1.
- [2] Jake Snell, Kevin Swersky e Richard S. Zemel. *Prototypical Networks for Few-shot Learning*. 2017. arXiv: 1703.05175 [cs.LG].
- [3] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip H. S. Torr e Timothy M. Hospedales. «Learning to Compare: Relation Network for Few-Shot Learning». In: *CoRR* abs/1711.06025 (2017). arXiv: 1711.06025. URL: <http://arxiv.org/abs/1711.06025>.
- [4] Yu Wang, Justin Salamon, Nicholas J. Bryan e Juan Pablo Bello. *Few-Shot Sound Event Detection*. 2020, pp. 81–85. DOI: 10.1109/ICASSP40776.2020.9054708.