

0 : Cosa c'è in questo notebook

Scaricamento ed utilizzo modelli di Natural Language Processing da huggingface.com per fare un po di esperimenti ed esempi con dati testuali per il progetto **geografia**.

In particolare si esplorano un po di metodi ed strumenti per la comprensione e annotazione dei testi, sfruttando modelli già allenati su task che ci possono interessare come ad esempio :

- Mask-Filling** (con [huggingface transformer](https://huggingface.com)): predice un token mancante in una frase
- Zero-Shot classification** (con [huggingface transformer](https://huggingface.com)): classifica una frase rispetto a delle classi variabili (utile per topic analysis, sentiment retrieval e menate generiche)
- Pos-Tagging** (con la libreria [spacy](https://spacy.io)) : analisi grammaticale e sintattica del testo

0.1 Transformer Models da hugging face

I modelli pre-allenati sono stati sottoposti ad estensivi training con tantissimi dati, in genere macinati per giorni, i task non-supervisonati tipici utilizzati per l'allenamento sono in genere:

- next sentence prediction (prendi un paragrafo, lo spitti sulla base del carattere '.' ed usi la prima frase come dato e la seconda come label da predirre)
- fill-mask: predire la parola mancante di una frase (nel training si prendono frasi complete e si maschera un termine a random per frase)

Accanto a questi task usati per allenare la struttura della rete neurale (aka "il grosso" dei nodi profondi), i modelli di [huggingface](https://huggingface.com) possono essere già "fine-tunati" rispetto ad ulteriori downstream task (tipo sentiment analysis o zero shot) e quindi utilizzabili direttamente. per alcune coppie di modelli-task invece è necessario un'ulteriore sessione di fine-tuning con un dataset rilevante per il task d'interesse. Qua c'è una lista delle tipologie di task piu comuni: <https://huggingface.com/tasks>

Infine un'altra differenza sostanziale che ci interessa è se il modello sia:

- monolingua
- multilingua

I modelli considerati per questo girovagare sono:

- il modello **dbmdz/bert-base-italian-cased** è più leggero e solo in italiano, ma di task ready-to-go disponibili c'è solo fill-mask e pos-tag
- il modello **facebook/bart-large-mnli** è molto pesante, ma è multilingua ed ha implementato lo zero-shot-classification
- il modello **bert-base-cased** solo inglese ma con praticamente tutti i downstream task già pronti all'utilizzo

0.2 Spacy Library

<https://spacy.io> libreria python piu semplice ed efficiente rispetto ai modelli transformer di [hugging face](https://huggingface.com), ma meno versatile per quantità di task possibili, ottima per pulire testi e tokenization (preprocessing), molto buono il pos-tagging

```
In [1]: from transformers import AutoModel, AutoTokenizer
from transformers import pipeline
multilingual_model_name = "facebook/bart-large-mnli" # 4.5 GB
italian_model_name = "dbmdz/bert-base-italian-cased" # 422 MB
alltask_model_name = "bert-base-cased"

## Download model and configuration from huggingface.co and cache.
multilingual_tokenizer = AutoTokenizer.from_pretrained(multilingual_model_name)
multilingual_model = AutoModel.from_pretrained(multilingual_model_name)

2022-05-03 17:55:08.489149: W tensorflow/stream_executor/platform/default/dso_loader.cc:64] Could not load dynamic library 'libcuda.so.11.0'; dLError: libcuda.so.11.0: cannot open shared object file: No such file or directory
2022-05-03 17:55:08.489192: I tensorflow/stream_executor/cuda/cudart_stub.cc:29] Ignore above cudart dLError if you do not have a GPU set up on your machine.
```

1 : Fill-Mask task example

```
In [2]: fill = pipeline('fill-mask', model=italian_model_name)
masked_sentence = 'Umberto Eco è [MASK] un grande scrittore'
fill(masked_sentence)

Some weights of the model checkpoint at dbmdz/bert-base-italian-cased were not used when initializing BertForMaskedLM: ['cls.seq_relationship.bias', 'cls.seq_relationship.weight'] - This is expected if you are initializing BertForMaskedLM from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPretraining model). - This IS NOT expected if you are initializing BertForMaskedLM from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

Out[2]: [{'score': 0.91440952267674585, 'token': 402, 'token_str': 'stato', 'sequence': 'Umberto Eco è stato un grande scrittore'}, {'score': 0.02569987252354622, 'token': 5801, 'token_str': 'diventato', 'sequence': 'Umberto Eco è diventato un grande scrittore'}, {'score': 0.022715188562878026, 'token': 409, 'token_str': 'anche', 'sequence': 'Umberto Eco è anche un grande scrittore'}, {'score': 0.006274290382862091, 'token': 1402, 'token_str': 'oggi', 'sequence': 'Umberto Eco è oggi un grande scrittore'}, {'score': 0.004773843102157116, 'token': 14743, 'token_str': 'divenuto', 'sequence': 'Umberto Eco è divenuto un grande scrittore'}]]

In [3]: masked_sentence = 'a che ora ci [MASK]?'
fill(masked_sentence)

Out[3]: [{'score': 0.1162104457616806, 'token': 7607, 'token_str': 'troviamo', 'sequence': 'a che ora ci troviamo?'}, {'score': 0.06763437390327454, 'token': 17567, 'token_str': 'aspettiamo', 'sequence': 'a che ora ci aspettiamo?'}, {'score': 0.06530724465847015, 'token': 1383, 'token_str': 'sarà', 'sequence': 'a che ora ci sarà?'}, {'score': 0.05703260377049446, 'token': 4238, 'token_str': 'vediamo', 'sequence': 'a che ora ci vediamo?'}, {'score': 0.053022969514131546, 'token': 17307, 'token_str': 'vedete', 'sequence': 'a che ora ci vedete?'}]]

In [4]: masked_sentence = 'a che [MASK] ci troviamo?'
fill(masked_sentence)

Out[4]: [{'score': 0.3658718466758728, 'token': 510, 'token_str': 'cosa', 'sequence': 'a che cosa ci troviamo?'}, {'score': 0.04696609797080885, 'token': 739, 'token_str': 'modo', 'sequence': 'a che modo ci troviamo?'}, {'score': 0.043758004903793335, 'token': 212, 'token_str': 'non', 'sequence': 'a che non ci troviamo?'}, {'score': 0.023330725729465485, 'token': 1302, 'token_str': 'livello', 'sequence': 'a che livello ci troviamo?'}, {'score': 0.0231430269777748, 'token': 1711, 'token_str': 'condizioni', 'sequence': 'a che condizioni ci troviamo?'}]]
```

2 : Zero-Shot classification task example

link modello multilingua (unico disponibile su [huggingface](https://huggingface.com) in grado di performare l'inferenza zero-shot in italiano)

usare generico in italiano sullo **zero-shot** learning e inference: <https://zephyrnet.com/tutorials/Zero-shot-learning-puoi-classificare-un-oggetto-senza-vederlo-prima/>

```
In [5]: zs = pipeline("zero-shot-classification", model=multilingual_model_name)

In [6]: zs('che bella giornata di sole', candidate_labels=['positivo', 'negativo']) #multiclass

Out[6]: {'sequence': 'che bella giornata di sole', 'labels': ['positivo', 'negativo'], 'scores': [0.9892858266830444, 0.010714183561503087]}

In [7]: zs(sequences='dove vai stasera?', hypothesis_template='questa frase è una {}', candidate_labels=['affermazione', 'domanda'])

Out[7]: {'sequence': 'dove vai stasera?', 'labels': ['domanda', 'affermazione'], 'scores': [0.006775164604107, 0.9932248651981354]}

In [8]: zs(sequences='where are you going?', hypothesis_template='this phrase is a {}', candidate_labels=['question', 'affirmation'])

Out[8]: {'sequence': 'where are you going?', 'labels': ['question', 'affirmation'], 'scores': [0.895594909324646, 0.10440445691347122]}

In [9]: zs(sequences='Voglio uscire con voi stasera, ma devo pulire casa prima', candidate_labels=['dubbio', 'invito', 'carro', 'positivo', 'negativo'], multiclass=True)

Out[9]: {'sequence': 'Voglio uscire con voi stasera, ma devo pulire casa prima', 'labels': ['invito', 'carro', 'negativo', 'dubbio', 'positivo'], 'scores': [0.7079709093203008, 0.00898021140655315, 0.00931985914707184, 0.053561460226774216, 0.02815048614131927]}

In [10]: #https://www.google.it/search?q=thayer+wood+plane+italia&biw=1528&bih=788&tbm=isch&source=iu&ictx=1&vet=1&fir=1dr_Cwra3dTA3M%252C3vr2PPAMlr9Tdm%252C%253BEHLhCfWfKHMM%252CE6xJH
moods=['fastidio', 'rabbia', 'sconforto', 'tristezza', 'noia', 'pace', 'relax', 'calma', 'felicità', 'gioia', 'indifferenza', 'interesse', 'dubbio', 'eccitazione']
zs(sequences='Sono felice di uscire con voi stasera', candidate_labels=moods, multiclass=True)

Out[10]: {'sequence': 'Sono felice di uscire con voi stasera', 'labels': ['felicità', 'relax', 'calma', 'interesse', 'eccitazione', 'fastidio', 'pace', 'indifferenza', 'gioia', 'noia', 'dubbio', 'rabbia', 'sconforto', 'tristezza'], 'scores': [0.4919441342353821, 0.1050914758082251, 0.11903546450130092, 0.05410005897283554, 0.03366483747959137, 0.03167872875928879, 0.022022755161573065, 0.0174749000787735, 0.015620590630428314, 0.013213590306741230, 0.00887034647166729, 0.00653752079235272, 0.005970512183010578, 0.00316450695340600014]}

In [11]: sentenza = 'non capisco perchè te la prendi con me, ti ho già regalato tutte le mi capre, non posso darti anche il bue'
zs(sequences=sentenza, candidate_labels=moods, multiclass=True)

Out[11]: {'sequence': 'non capisco perchè te la prendi con me, ti ho già regalato tutte le mi capre, non posso darti anche il bue', 'labels': ['sconforto', 'dubbio', 'indifferenza', 'calma', 'relax', 'tristezza', 'interesse', 'pace', 'eccitazione', 'gioia', 'noia', 'fastidio', 'rabbia', 'felicità'], 'scores': [0.30105090081993103, 0.1431449055671692, 0.12648096806041174, 0.09243334829087281, 0.00247055490091848, 0.07852306216955185, 0.05074809119105339, 0.032442253082990046, 0.027599514051570313, 0.023189708590507507, 0.015240885317325592, 0.011493487283507456, 0.008943037129938002, 0.006272169205990404]}

In [12]: zs(sequences='le tue orecchie sono bruttissime, non uscirò mai con te', candidate_labels=['estetica', 'logica', 'attualità', 'cotillons'], multiclass=True)

Out[12]: {'sequence': 'le tue orecchie sono bruttissime, non uscirò mai con te', 'labels': ['estetica', 'logica', 'attualità', 'cotillons'], 'scores': [0.3290303647510158, 0.27184581750591797, 0.206147315502016075, 0.1329764872789383]}

*** Pos-Tagging
WORD POS-TAG SYNTACTIC DEP
-----
San PRON nsubj
Francisco PRON flat:name
prevede VERB ROOT
di ADP mark
bandire VERB xcomp
i DET det
robot NOUN obj
di ADP case
consegna NOUN nmod
porta VERB advcl
a ADP case
porta NOUN obl

*** Analisi Entità specifiche nel testo:
ENTITY NAME LABEL
-----
San Francisco LOC
```

```
In [20]: doc = nlp(sentences[3])
print('testo originale --> ', doc.text)
print('\n*** Pos-Tagging')
print('WORD', ' '*10-len('word'), 'POS-TAG', ' '*10-len('POS-TAG'), 'SYNTACTIC DEP')
print('-----')
for token in doc:
    print(token.text, ' '*10-len(token.text), token.pos_, ' '*10-len(token.pos_), token.dep_)

print('\n*** Analisi Entità specifiche nel testo:')
print('ENTITY NAME', ' '*20-len('ENTITY NAME'), 'LABEL')
print('-----')
for ent in doc.ents:
    print(ent.text, ' '*20-len(ent.text), ent.label_)

testo originale --> London è una grande città del Regno Unito.

*** Pos-Tagging
WORD POS-TAG SYNTACTIC DEP
-----
London PRON nsubj
è AUX cop
una DET det
grande ADJ amod
città NOUN obj
del ADP case
Regno PRON nmod
Unito PRON flat:name
. PUNCT punct

*** Analisi Entità specifiche nel testo:
ENTITY NAME LABEL
-----
London LOC
Regno Unito LOC
```

```
In [21]: doc = nlp(sentenza.replace('me', 'Giacomo').replace('mi', 'sue').replace('ho', 'ha').replace('posso', 'può'))
print('testo originale --> ', doc.text)
print('\n*** Pos-Tagging')
print('WORD', ' '*10-len('word'), 'POS-TAG', ' '*10-len('POS-TAG'), 'SYNTACTIC DEP')
print('-----')
for token in doc:
    print(token.text, ' '*10-len(token.text), token.pos_, ' '*10-len(token.pos_), token.dep_)

print('\n*** Analisi Entità specifiche nel testo:')
print('ENTITY NAME', ' '*20-len('ENTITY NAME'), 'LABEL')
print('-----')
for ent in doc.ents:
    print(ent.text, ' '*20-len(ent.text), ent.label_)

testo originale --> non capisco perchè te la prendi con Giacomo, ti ha già regalato tutte le sue capre, non può darti anche il bue

*** Pos-Tagging
WORD POS-TAG SYNTACTIC DEP
-----
non ADV advmod
capisco VERB ROOT
perchè SCONJ mark
te PRON expl
lì PRON obj
prendi VERB ccomp
con ADP case
Giacomo PRON obl
di PUNCT punct
i PRON nsubj
ha AUX aux
già ADV advmod
regalato VERB conj
tutte DET det:predet
le DET det
sue DET det:poss
capre NOUN obj
non ADV advmod
può AUX aux
darti VERB conj
anche ADV advmod
il DET det
bue NOUN obj

*** Analisi Entità specifiche nel testo:
ENTITY NAME LABEL
-----
Giacomo PER
```

```
In [ ]: 
```