

# Önálló projekt 2

Formális és program nyelvek elemzése gépi tanulási modellekkel

Sisák László Sándor

Témavezető: Lukács András

**Feladat:** Python kódok bináris klasszifikációja gépi tanulás segítségével.

**Nyers adat:** *Python Questions from Stack Overflow*, 607000 kérdés, 987000 válasz. (forrás: Kaggle)

## Python and MySQL

Asked 13 years, 9 months ago Modified 2 years, 7 months ago Viewed 4k times

I can get Python to work with PostgreSQL but I cannot get it to work with MySQL. The main problem is that on the shared hosting account or MySQL, I generally fail when installing the host.

I found [bpgsql](#) really good because I read and then call the functions of it.

python mysql postgresql bpgsql

I don't have any experience with <http://www.SiteGround.com> as a web host personally.

This is just a guess, but it's common for a shared host to support Python and MySQL with the MySQLdb module (e.g., GoDaddy does this). Try the following CGI script to see if MySQLdb is installed.

```
#!/usr/bin/python

module_name = 'MySQLdb'
head = '''Content-Type: text/html

%s is ''' % module_name

try:
    __import__(module_name)
    print head + 'installed'
except ImportError:
    print head + 'not installed'
```

# Az adat megtisztítása

## Python and MySQL

Asked 13 years, 9 months ago Modified 2 years, 7 months ago Viewed 4k times

45 I can get Python to work with PostgreSQL but I cannot get it to work with MySQL. The problem is that on the shared hosting account or PySQL, I generally fail when installing the host.

I found [bpgsql](#) really good because I read and then call the functions of it.

python mysql postgresql bpgsql

7 I don't have any experience with [http://www.mysql.com/doc/connector-python/en/mysql-connector-python.html](#).

This is just a guess, but it's common for a MySQLdb module (e.g., GoDaddy does the same) to not be installed.

```
#!/usr/bin/python
module_name = 'MySQLdb'
head = '''Content-Type: text/html

%s is ''' % module_name

try:
    __import__(module_name)
    print head + 'installed'
except ImportError:
    print head + 'not installed'
```

Tags	Code_snippets
[sql, database]	[fetchall(), curs.execute('select first_name f...
[arrays, iteration]	[#/bin/python bar in dict(Foo) , has_ke...
[arrays, iteration]	[foo = 12 foo in iter_attr('bar', 1d'), foo ...
[mysql, postgresql]	[#/usr/bin/python  module_name = 'MySQL...
[python]	[groups = [ununiquekeys = [n] for n in gro...
[php]	[f = open('logfile', 'rb')in]
[oop, methods]	[setattr]
[urllib]	[info(), urllibobject.info()['Content-Length']]
[windows, image, pdf]	[ps.alpha]
[iteration]	[groups = [ununiquekeys = [n] for n, g in gro...

szűrés

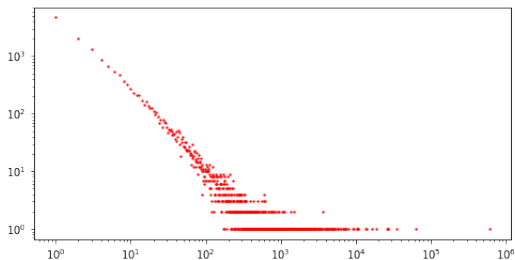
szűrés

A címkéket gyakoriság, a kódokat szintaktikai helyesség alapján szűrtem.

# Az adat megtisztítása

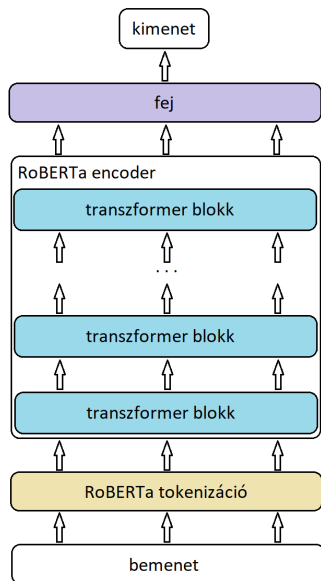
A közel 17000 címkéből csak 207 volt, ami legalább 1000 különböző kérdésnél szerepelt.

A kapott adathalmaz 492000 adatpontot tartalmazott.



A címkék gyakorisága.  $(x, y)$  pont azt jelöli, hogy  $y$  különböző címke van az adathalmazban, amely pontosan  $x$  különböző kérdésnél szerepel.

# A modell felépítése



**Tokenizáció:** A bemeneti stringet egészértékű vektorra alakítja.

**Encoder:** A tokenizált bemenetből vektorrepresentációt állít elő. Az így kapott vektor a bemenet jellemzőit írja le.

A tokenizációhoz és a beágyazás előállításához az előtanított RoBERTa transzformer alapú modellt használtam.

**Fej:** Egy sűrű réteg, ami a beágyazást klasszifikálja.

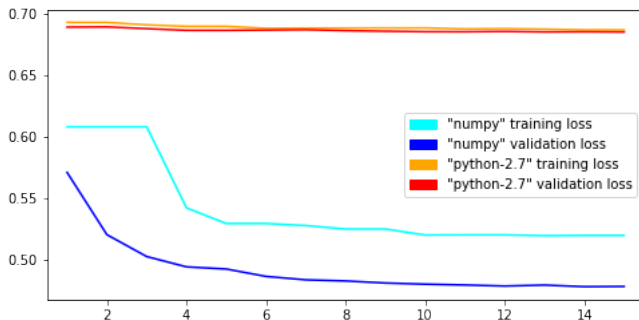
# A modell tanítása

A modell csak az öt leggyakoribb címkén tanult. Az adathalmaz ferdeségét korrigálnom kellett.

Címke	Össz. előfordulás	Előfordulás arány	Tanítóadat mérete
<i>django</i>	55105	11,19%	88168
<i>python-2.7</i>	37313	7,58%	59701
<i>python-3.x</i>	33387	6,78%	53419
<i>numpy</i>	30759	6,25%	49214
<i>pandas</i>	28227	5,73%	45163

# A modell értékelése

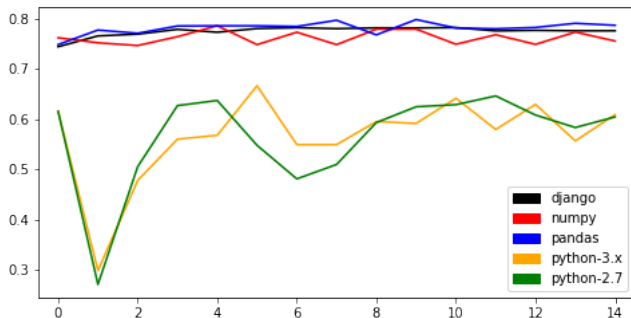
A modell az egyes címkéken 15 epochon keresztül volt tanítva.



Példa a veszteségfüggvény alakulására egy nehezen és egy könnyen tanulható címke esetén.

# A modell értékelése

A teljesítményt  $F_1$ -scoreban mértem. Az ábra a validációs halmazon mért  $F_1$ -score alakulását mutatja a különböző címkékre.





# Kitekintés: adathalmaz

Érdeemes több adatot gyűjteni.

A nagyon rövid kódrészletekből nem tud tanulni a modell, de könnyen szűrhetőek hossz alapján.

A modell figyelembe vehetné a kódrészletek listájának összes elemét.  
Az adat szűrhető pontszám alapján is.

▲ What you're looking for is `setattr` I believe. Use this to set an attribute on an object.

11

```
>>> def printme(s): print repr(s)
>>> class A: pass
>>> setattr(A,'printme',printme)
>>> a = A()
>>> a.printme() # s becomes the implicit 'self' variable
< __main__ . A instance at 0xABCDEF >
```

Share Follow

answered Aug 7, 2008 at 11:30



HS.

14.7k ● 8 ● 38 ● 47

Célom kiterjeszteni a feladatot többcímkes klasszifikációra. Ekkor az egyes címkékre nagyon ferde lesz az adathalmaz.

Érdemes lehet saját tokenizációt tervezni.

A RoBERTán kívül más (előtanított) modellt is érdemes kipróbálni.

A várható jobb beágyazás mellett érdemes összetettebb fejet tanítani.

Köszönöm a figyelmet!