

RPH. Naive Bayes spam filter

Vadym Ostapovych, Roman Novikov

4. ledna 2022

Úvod

Úkolem dané úlohy je vytvoření spam filtru, který bude umět odfiltrovat spam a korektní emaily. Pro její řešení bylo rozhodnuto vytvořit klasifikátor založený na Bayesove podmíněné pravděpodobnosti.

1. Matematický princip filtru

Naive Bayes je trénovací filtr, který vyžaduje pevnou sadu trénovacích dat pro určení parametrů klasifikátoru. Za nich budeme považovat klasifikované emaily (třídy)

- SPAM - email je spamový.
- OK - email je korektní (ham).

Hlavní idea uvedeného spam filtru je z sady trénovacích emailů zformovat svůj slovník "rysových slov" (feature words), taký k nim mohou patřit nějaké větší fráze. Pod slovníkem budeme rozumět dvě množiny slov, které budou patřit nebo v spam, nebo v ham emailech. Pro konkrétní slovo budeme počítat dva typy slovníků:

- 1) Kolikrát slovo se objevuje v konkrétní třídě (class), např. při učení slovo "slave" mezi třemi spam emaily, objevilo 5-krát., proto do toho slovníku přidáme s klíčem "slave" bude číslo 5.
- 2) Kolikrát slovo se objevuje v konkrétní třídě na počet emailů tohoto classu. Uvažujeme minulý příklad se slovem "slave", ale nechť z těch 3 spam emailu, jen ve dvou z nich se objevily ty 5 slov, proto do toho slovníku přidáme číslo 2.

Z dvou uvedených slovníků musíme vypočítat parametry filtru pro další klasifikaci emailu. V našem případě to budou podmíněné pravděpodobnosti "likelihoods" pro jednotlivá slova, že email je spam, resp. ham, za podmínky, že obsahuje nějaké slovo. Uvedeme vzorec pro výpočet spamovosti emailu

$$P(S|W) = \frac{N_{W|S}}{N_{spam\ words} + N_{uniq}} \cdot \frac{N_{W|per\ spam\ email}}{N_{spams}} \quad (1)$$

kde

- $P(S|W)$ - pravděpodobnost, že email je spam za podmínky, že obsahuje slovo W .
- $N_{W|S}$ - počet slova W ve všech emailech.
- $N_{spam\ words}$ - celkový počet objevování rysových slov v spam emailech.
- N_{uniq} - celkový počet rysových (unikátních) slov, klíčů.
- $N_{W|per\ spam\ email}$ - v kolika spamových emailech se objevuje dané slovo.
- N_{spams} - celkový počet spamových emailů.

Z rovnice lze vidět, že podmíněná pravděpodobnost se skládá z dvou jednotlivých zlomků, které identifikují Bayesove pravděpodobnosti objevování slova vůči celkovému počtu spamových slov a vůči celkovému počtu spamových emailů. Stejný vzorec použijeme pro výpočet pro hamovosti slova, kde budeme počítat už vůči hamovým slovům a hamovým emailem. Po výpočtu dostaneme dva slovníky s pravděpodobnostmi pro hamové a spamové slova F_{hams} , F_{spams} . Z nich už budeme klasifikovat nové emaily.

V novém emailu nalezneme unikátní nebo rysová slova a přidáme do nějaké množiny např. F_{email} , počet slov v emailu nás nezajímá. Pak z té množiny vypočítáme

$$F_{email} \cap F_{spams} - \text{pravděpodobnosti slov, které se objevily v spamu během trenování} \quad (2)$$

$$F_{email} \cap F_{hams} - \text{pravděpodobnosti slov, které se objevily v hamu během trenování} \quad (3)$$

Pro každou z nalezených množin spočítáme pravděpodobnost, že email je spam, resp. ham. Ukážeme výpočet pro spam, pro ham je totiž stejný.

$$P_{S_{email}} = P(S) \cdot P(S|W_1) \cdot P(S|W_2) \cdot \dots \cdot P(S|W_n) \quad (4)$$

kde

- $P(S|W_i)$ - pravděpodobnost, že email je spam za podmínky, že obsahuje slovo W_i .
- P_S - pravděpodobnost spamu v novém emailu, může být vypočten z počtu spam a ham emailů nebo rozumně nastaven.

Z pravděpodobností $P_{S_{email}}$, $P_{H_{email}}$ už dáváme nový email do třídy podle vyšší pravděpodobnosti

$$\text{Email class} = \text{argmax}(P_{S_{email}}, P_{H_{email}}) \quad (5)$$

2. Implementace

Všechna trénovací část je implementována v trénovacím corpusu. V třídě **Filter** jen zavoláme metodu z trénovacího corpusu pro nastavení parametrů filtru.

Vyber rysových slov

Pro výpočet rysových slov potřebujeme dva slovníky, v každém z nich klíčem je jednotlivý class SPAM nebo OK. Hodnotami klíčů jsou Countery, které identifikují počet slov. První slovník, který představuje celkový počet slov v jednotlivých třídách slov se jmenuje **vocabulary**, druhý představuje počet těch slov na každý email, jmenuje **words_per_email** a oba jsou atributy trénovacího corpusu.

Budeme postupně procházet jednotlivé emaily každé třídy, pro parsing emailu použijeme modul **email**, který vyčleňuje hlavní content. Z toho hlavního contentu pomocí regulárních výrazů nalezneme adresy urls, názvy emailu, které se zmiňují a typy contentu, které jsou součástí modulu **email**. To můžeme považovat za jednotlivé fráze. Dále spočítáme unikátní slova v tomto emailu, pomocí regulárního výrazu odstraníme všechna písmena, která nejsou z anglické abecedy, a nahradíme je prázdným symbolem, uděláme všechna písmena malými a pomocí funkce **split()** zformujeme list slov. Také dáváme pozor na minimální a maximální délku slov a jestli to slovo není v stop words, slova, od kterých v anglickém jazyce veta nemění moc smysl. Délky byly nastavené na [2, 100] a stop slova byly použita stejně jako v modulu **nlTK**. Uvedeme příklad vyber unikátních slov na větu **Van is a Dungeon Master since 1990**

- is, a jsou stop wordy.
- 1990 není patří do anglické abecedy, je to číslo. áme ostatní slova do malého písma a jelikož respektují délky, uděláme list ["van", "dungeon", "master"]

Counter pomáhá při výpočtu objevování slov. Výsledkem by bylo Counter("van":1, "dungeon":1, "master":1). Na základě toho uděláme dva slovníky.

Výpočet podmíněných pravděpodobností z slovníků

Pro výpočet podmíněných pravděpodobností byl použit princip Laplace smoothing pro lepší predikci. Např. my představujeme, že slovo, které se objevilo pouze v hamu se taky objevilo v k-krát spamových emailu, stejně pro ham. Pokud A - množina spamových slov, B - množina hamových slov, lze to zapsat následně pro $k = 1$, což je v podstatě defaultní konstanta

$$A_{new} = A \cup (A \setminus B) \quad (6)$$

$$B_{new} = B \cup (B \setminus A) \quad (7)$$

Z toho zápisu se objevuje, že délka klíče nebo slov pro oba nové slovníky musí být stejná. A proto ve vzorci pro výpočet unikátních slov N_u můžeme použít libovolný z Counteru buď pro ham nebo spam.

Predikce třídy

Pro predikci classu budeme považovat

$$P(S) = P(H) = 0.5 \quad (8)$$

jelikož nelze předpokládat z trénovacích dat, že pravděpodobnost, že nový email patří do nějaké třídy je stejná jako v trénovacích datech. Také pro výpočet predikční pravděpodobnosti byly použité logaritmy, abychom se vyhnuli floating-point underflow. Proto prakticky se koná srovnání logaritmu

$$\log(P_{S_{email}}) = \log(P(S)) + \sum_{i=1}^n \log(P(S|W_i)) \quad (9)$$

$$\log(P_{H_{email}}) = \log(P(H)) + \sum_{i=1}^n \log(P(H|W_i)) \quad (10)$$

A proto rozhodnutí se koná jako

$$\text{Email class} = \operatorname{argmax}(\log(P_{S_{email}}), \log(P_{H_{email}})) \quad (11)$$

3. Trenovani klasifikatoru

Evaluační klasifikátor se konala na datasetech 1, 2 uvedených v zadání. Trenování se konalo také na uvedených datasetech. Výsledky jsou v tabulce

Trénovací dataset	Evaluační dataset	Kvalita filtru q
1	1	0.953
1	2	0.931
2	1	0.836
2	2	0.97

Tabulka 1: Porovnání kvality filtru

Lze vidět, že kvalifikátor má dostatečně vysokou kvalitu, možný problém neúplné klasifikace je nesprávná detekce klíčových slov nebo malý počet klíčových fraz. Různá možná vylepšení jako počet unikátních slov v novém emailu nebo výpočet $P(S)$, $P(H)$ z trénovací dávalo horší výsledky, proto uvedená implementace je považována za nejlepší dosahnutou. Na Brutu byly dosaženy následující výsledky Z výsledku turnaje se objevilo,

Trénovací dataset	Kvalita filtru q	Místo v turnaji
1	0.932	6
2	0.855	8
3	0.853	13

Tabulka 2: Výsledky turnaje

že filtr má velký počet False Negative a malý počet False Positives, což ukazuje, že nesprávně detekuje ham emaily a dost často je považuje za spam, což také potvrzuje možnou nesprávnou detekci unikátních slov.

Výsledek práce

Přínos práce každého z účastníků byl rovnocenný 50 na 50. Každý udělal něco pro jednotlivé části práce. Návrh filtru je docela správný, a detekce je vysoká, ačkoliv není ideální. Možné zlepšení je přidat detekci správných slov, což by šlo použít knihovnu nltk. Jinak pro tuto úlohu byl dosažen maximální počet bodů.

Použitá literatura

- Wikimedia Foundation. (2021, December 3). Naive Bayes spam filtering. Wikipedia. Retrieved January 4, 2022, from https://en.wikipedia.org/wiki/Naive_Bayes_spam_filtering.
- Spam filter in Python: Naive bayes from scratch. KDnuggets. (n.d.). Retrieved January 4, 2022, from <https://www.kdnuggets.com/2020/07/spam-filter-python-naive-bayes-scratch.html>

- Motwani, S. (2020, August 8). Email spam filtering using naive Bayes classifier. Springboard Blog Email Spam Filtering Using Naive Bayes Classifier Comments. Retrieved January 4, 2022, from <https://in.springboard.com/blog/email-spam-filtering-using-naive-bayes-classifier/>