

## *Název práce:*

Spam filter

## *Autory:*

Nazar Ponomarev, Nikita Kisel.

## *Stručný úvod:*

Cílem bylo napsat program spam filter, který bude hodnotit emaily jako spam nebo ham.

## *Popis principu/algoritmu použitého spam filtru:*

Spam filter se učí na datasetu a pak hodnotí emaily na základě získané informace. Spam filter na konci učení má dva slovníky. První je báze spam a ham odesílatelů. Druhý obsahuje slova a jejich výskyt ve spam nebo ham emailech. Hodnocení emailu probíhá nejprve porovnáním adresy odesílatelů emailů, jestli ona se nachází v bázi spam a nenachází v hamu, filter hodnotí email jako spam, a naopak jestli adresa se nachází v bázi ham a nenachází ve spamu, filter hodnotí email jako ham. Když žádná ze dvou možností nenastala filter pokračuje práce se druhým slovníkem. Pro každé unikátní slovo z emailu filter počítá jeho výskyt ve spam a ham položkách slovníku. Jako výsledek na konci filter má dva čísla - výskyt slov ve spamu a hamu. Podle těchto čísel filter spočítá procento spamu v emailu. Když procento spamu v emailu větší než 60, filter hodnotí email jako spam. Číslo 60 (ne 50) není náhodné, protože v případě chyby je lepší ohodnotit email jako ham. Filter se umí adaptovat k datasetu, proto existuje koeficient násobení. Když spam emailu více než ham emailů, koeficient je  $> 1$ , a naopak, když spam emailu méně než ham emailů, koeficient je  $< 1$ . Přesný koeficient spočítá poměrem čísla spam a ham emailů v datasetu. Také filter má jednu další možnost ve případě výjimky, když součet čísla výskytu slov ve spamu a hamu nulový, filter hodnotí email jako ham (protože v případě chyby je lepší ohodnotit email jako ham).

## *Popis způsobu trénování filtru:*

Rozhodli jsme se, že budeme trénovat a testovat filter na doporučených datech v coursewaru. Bylo to lepší pro nás, protože ze začátku věděli, jak budou vypadat emaily a co můžeme z nich použít pro trénování. Pro trénování filter používá způsob tokenizaci emailů. Všimli jsme si, že filter může vždy vzít adresu odesílatele z emailu, protože první slovo v úhlových (<>) závorkách je vždycky adresa emailu odesílatele. Pak jsme si všimli, že zpráva v emailu vždy jde po prvním "\n\n". Filter na základě našich znalostí postupně zaplňuje z emailu slovník:

- Jestli email je ve bázi 64, převádí ho do normálního textu.
- Dělá všechny písmena ve zprávě malými.
- Nahrzuje interpunkční znaky mezerami.
- Dělí zprávu podle mezer a nových řádků (\n).
- Odstraní neúčinná slova ze zprávy.
- Plnuje slovník všimly zbývajícími slovy s délkou větší než 3 s výjimkou několika slov.

Filter přiřadí druhé položce slova ve slovníku jeho výskyt ve spamu nebo hamu.

## *Výsledky dosažené naším spam filtrem:*

Používali jsme jenom jeden základní typ filtru. Experimentovali s různými funkcemi a koeficienty uvnitř nich (používali jsme data z coursewaru):

1. Empiricky jsme se rozhodli o hodnocení emailu jako spam jenom z 60%, protože 50% je špatný případ pro chybové rozhodnutí. Kvalita filtru se zvýšila o asi 17%.
2. Odstranili jsme z porovnání slova s délkou menší než 3, kvalita filtru se zvýšila o asi 22%.
3. Udělali jsme báze spam a ham odesílatelů. Která zvýšila kvalitu filtru o asi 5%.

## *Stručný popis rozdělení práce v týmu:*

	utils.py	corpus.py	tokenization.py	training_corpus.py	prediction_corpus.py	filter.py
Nazar Ponomarev	+	+	+	+		
Nikita Kisel	+	+			+	+

## *Stručný popis organizace práce v týmu:*

Bydlíme na jedné koleji, proto dělali jsme filtr spolu ve studovně. Pro přenos modulu jsme používali Telegram (<https://telegram.org>).

## *Zhodnocení, závěr:*

Práce byla velmi zajímavá, naučili jsme se pracovat s datami a práci spolu.

## *Seznam použité literatury, online zdrojů:*

- <https://www.youtube.com/watch?v=VDq8fCW8LdM>
- [https://www.w3schools.com/python/python\\_reference.asp](https://www.w3schools.com/python/python_reference.asp)
- [https://en.wikipedia.org/wiki/Lexical\\_analysis#Tokenization](https://en.wikipedia.org/wiki/Lexical_analysis#Tokenization)
- <https://docs.python.org/release/3.1.3/library/stdtypes.html#string-methods>
- <https://docs.python.org/release/3.1.3/library/re.html?highlight=re#module-re>