

Temă practică – Învățare automată 2023 Fall

Deadline: end-of-day 12 ianuarie 2024

Studiați, din punct de vedere teoretic și experimental, gradul de adaptare/adekvare al algoritmilor studiați până acum în raport cu problema de clasificare a email-urilor spam pentru setul de date Ling-Spam:

http://www.aueb.gr/users/ion/data/lingspam_public.tar.gz.

Cerințe:

- Înțelegeți setul de date (care sunt atributele, etichetele, cum le extrageți din reprezentarea textuală) — este necesar ca această cerință să fie, pe scurt, documentată sub forma procesării datelor înainte de introducerea în algoritm. Mesajele de tip spam conțin acest indiciu în titlul fișierului (sub forma prefixului “*spm*”). Utilizați, în cele ce urmează, 9 foldere (de la part1 la part9) pentru antrenare și păstrați câte unul pentru testare (cel intitulat part10), din fiecare dintre categoriile *lemon*, *bare*, *stop*, *lemon_stop*.
- Selectați și implementați un algoritm, dintre cei învățați, pe care îl considerați adecvat rezolvării acestei probleme. (0.75 puncte)
- Justificați într-un raport LaTeX alegerea făcută, din punct de vedere teoretic și experimental, atât în mod individual, cât și prin comparație cu ceilalți algoritmi candidați la tipul de problemă studiată. (0.75 puncte)
- Implementați strategia de cross-validare Leave-One-Out și atașați raportului un grafic care să ilustreze statistic rezultatele. (0.25 puncte)
- Adăugați la raport un grafic care să dovedească performanța algoritmului vostru pe setul de date de testare, din punct de vedere al acurateții obținute. Acuratețea obținută trebuie să fie relevantă, moderat mai bună decât orice strategie trivială (dat cu banul sau ales mereu aceeași clasa). Dacă ați testat mai mulți algoritmi, includeți grafice comparative, atât prin raportare la această cerință, cât și la cea precedentă. (0.25 puncte)
- Explicați și alte detalii ale experimentului care vi se par relevante, în cuvinte sau în mod grafic. Puteți cerceta și variante îmbunătățite față de varianta clasică a algoritmului, studiată la seminar, spre a le implementa și a spori acuratețea.

Bonus: Presupunând că unele date nu ar fi etichetate cu titlul de spam/non-spam, ci ar fi disponibile doar reprezentările textuale fără a fi clasificate în prealabil, implementați și descrieți un mod prin care ați putea preprocesa datele, tot prin strategii cunoscute, astfel încât să le introduceți apoi în cadrul aceluiași algoritm implementat, nemodificat. Scopul acestui exercițiu este să găsiți un mod de a folosi eșantioane (“sample”-uri) fără etichetă pentru a îmbunătăți performanța algoritmului. Pe lângă descrierea efectivă a pasului de preprocesare, trebuie să includeți în raport grafice pentru eroarea/acuratețea la CVLOO și acuratețea la testare și pentru această metodă. Considerați ca fiind neetichetate datele din folderele part1 și part2 ale fiecărei categorii (ignorați dacă fișierul are sau nu prefixul “*spm*”). (0.5 puncte)

Lucru în echipe de maxim 2 persoane.

Punctaj total: 2 puncte (temă) + 0.5 puncte (bonus)