

Phân lớp thư rác với giải thuật học Arcx4-rMNB

Đỗ Thanh Nghị

BM. Khoa học máy tính
Khoa Công nghệ thông tin
Số 1 Lý Tự Trọng, Ninh Kiều, Cần Thơ
dtnghi@cit.ctu.edu.vn

Ngày 7 tháng 10 năm 2011



Nội dung

Giới thiệu

Giải thuật Arcx4 of random MNB

Kết quả thực nghiệm

Kết luận, hướng phát triển

Dịch vụ liên lạc phổ biến

Dịch vụ thư điện tử

- ▶ ưu điểm: đơn giản, nhanh chóng, chi phí thấp
- ▶ nhiều người sử dụng

Vấn nạn thư rác

- ▶ quảng cáo, khiêu dâm, phản động, thậm chí là những đoạn mã độc hại đính kèm
- ▶ gây lãng phí, phiền toái



Ví dụ minh họa

5-12 Free Bonus Viagra100mg pills with every order. All Major Credit Cards Accepted. Cialis20mg \$2.97, Viagra100mg \$1.87 37w [Spam](#)

 Tawana Celestine <tawana_celestinebj@charleskendall.com> [to show details](#) Aug 17 [Reply](#) | [▼](#)

Foreign Pharmacy | Discount & Cheap Drugstore

An Online Pharmacy provides Cheap drugs at your doorstep

CialisPlusViagra Powerpack special price
ViagraAs low as \$1.85
ViagraPlus as low as \$2.95
ViagraProfessional as low as \$3.95
ViagraSuperForce as low as \$4.99
CialisAs low as \$2.40
CialisSuperActive+ as low as \$3.33
LevitraAs low as \$2.50
SomaAs low as \$1.06

Order with us and save your medical bills up to 80-90%.

We have worldwide customers ...



Thiệt hại

Thống kê về thiệt hại do thư rác

- ▶ nghiên cứu cho thấy tổn thất năm 2002 ước tính khoảng **2 tỷ đô la** (Sung-jin, 2003)
- ▶ một nghiên cứu khác về tổn thất tháng 10 năm 2003 ước tính khoảng **10,4 tỷ đô la** (Mi2g, 2003)
- ▶ theo (Doug, 2003): nếu một người phát tán thư rác thu được **lợi 10 ngàn đô la** trong 1 tháng thì tổn thất do họ gây ra là **100 ngàn đô la**.

Lọc thư rác

- ▶ giảm bớt tổn thất



Lọc thư rác dựa trên địa chỉ người gửi

Dễ bị qua mặt: sinh ngẫu nhiên địa chỉ người gửi

S: 220 crepes.fr
C: HELO cit.ctu.edu.vn
S: 250 Hello cit.ctu.edu.vn, pleased to met you
C: MAIL FROM: <alice@cit.ctu.edu.vn>
S: 250 alice@cit.ctu.edu.vn ... Sender ok
C: RCPT TO: <peter@crepes.fr>
S: 250 peter@crepes.fr ... Recipient ok
C: DATA
S: 354 Enter mail, end with “.” On a line by itself
C: How are you?
C: Do you like internet?
C: .
S: 250 Message accepted for delivery
C: QUIT



Lọc thư rác dựa trên từ khóa

Dễ bị qua mặt: biến thể

				1 - 50 of 94 Older >
<input type="checkbox"/>	Mitchell Cindie	BetterEjaculation control, Experience Rock-HardErections on yourPenis fmitg - You Have	6:18 pm	
<input type="checkbox"/>	Liza Yasmin	Explosive, intenseOrgasms, Increase Volume ofEjaculate, Doctor designed and endorsed	Oct 17	
<input type="checkbox"/>	me	Dear hptoan, 70% off and more. units - If you have any difficulty viewing this email click here	Oct 17	
<input type="checkbox"/>	Best Pharma Online	Mr. hptoan, exclusive 60% off. a Juliana - This message contains graphics. If you do not see the	Oct 17	
<input type="checkbox"/>	Dreama Jeanette	ViagraDiscounts, CheapCialis & Much More. Discreet Packaging and Fast Shipping tjws -	Oct 17	
<input type="checkbox"/>	Al Starr	Cheapest Cializ+Viagre = \$78, this is 80% lower than Retail price!! We are here to save... -	Oct 16	
<input type="checkbox"/>	Iyriuzak5055	Hi hptoan, it's your Sale notifier! - http://drugstorema.ru	Oct 16	
<input type="checkbox"/>	Pfizer Online	Dear hptoan, It's your personal discount. th - If you have any difficulty viewing this email click	Oct 16	
<input type="checkbox"/>	Talitha Billi	High QualityMedications + Discount On All Reorders = Best Deal Ever! Viagra50/100mg - ...	Oct 16	
<input type="checkbox"/>	Genuine Pfizer	Mr. hptoan, 80% OFF for you. protecting used - This message contains graphics. If you do not	Oct 16	
<input type="checkbox"/>	Neva Shan	Discount CialisViagra from \$1.30, Express delivery, 90000+ Satisfied US, UK, CANADIAN C...	Oct 15	
<input type="checkbox"/>	Pfizer Online	Dear hptoan, we start Sale. countries - If you have any difficulty viewing this email click here	Oct 14	
<input type="checkbox"/>	Best-quality Pfizer	Mr. hptoan, exclusive deal for you. northern Paris heavy - If you are unable to see the message	Oct 14	
<input type="checkbox"/>	Pfizer	Hi hptoan, our Sale starts. the - Viewing difficulties? Check out the online version of this email.	Oct 14	
<input type="checkbox"/>	Lory Graciela	Codeine/Phentermine/Hydrocodone/Vicodin 7.5/750mg \$3.90/pill, NoPrescription, Shipping	Oct 14	
<input type="checkbox"/>	Rosenda Sumiko	High QualityMedications + Discount On All Reorders = Best Deal Ever! Viagra50/100mg - ...	Oct 14	
<input type="checkbox"/>	Pfizer	Hi hptoan, our Sale starts. save Alaska Great had the - Viewing difficulties? Check out the online	Oct 13	
<input type="checkbox"/>	me	Hi hptoan, our Sale starts. the is m battle - Viewing difficulties? Check out the online version of	Oct 13	
<input type="checkbox"/>	Chante Tamika	Need affordable Drugs?? Purchase Online here: GenericViagr \$2.23, GenericCialis \$2.80 ...	Oct 13	
<input type="checkbox"/>	Best-quality Pfizer	Hi hptoan, gets discount today. the follows Wounds the ancient - Viewing difficulties? Check	Oct 13	
<input type="checkbox"/>	Ria Melinda	V.i.a.g.r.a.	Oct 13	

Lọc thư rác dựa trên nội dung

Mô hình túi từ và học máy để phát hiện thư rác

- ▶ nội dung thư (không có cấu trúc): biểu diễn về cấu trúc bảng
- ▶ mô hình túi từ: thư điện tử biểu diễn dạng véctơ có giá trị thành phần thứ i là tần số xuất hiện từ thứ i trong thư
- ▶ tập thư điện tử: bảng (ma trận), mỗi dòng là một thư, mỗi cột tương ứng với một từ trong tự điển
- ▶ xây dựng mô hình phân lớp thư rác: **số chiều rất lớn** đến vài chục ngàn, **mỗi chiều chứa ít thông tin** cho phân lớp



Mô hình phân lớp dữ liệu có số chiều rất lớn

Đề xuất giải thuật mới (Arcx4-rMNB)

- ▶ Arcx4 of random Multinomial Naive Bayes (Arcx4-rMNB)
- ▶ đơn giản, nhanh, cho độ chính xác cao
- ▶ kết quả thử nghiệm với 1921 thư điện tử
- ▶ đạt được độ chính xác đến **95.81%**
- ▶ cải thiện hơn **45%** độ chính xác của MNB (Lewis et al., 1994)
- ▶ phân lớp chính xác như SVM (Vapnik, 1995)



Multinomial Naive Bayes (MNB)

Giải thuật

- ▶ tập học gồm T văn bản, N từ vựng và C lớp
- ▶ gán văn bản t_i tới lớp có ước lượng xác suất cao nhất $P(c|t_i)$.

$$P(c|t_i) = \frac{P(c)P(t_i|c)}{P(t_i)} \quad c \in C \quad (1)$$

$$P(t_i|c) = (\sum_n f_{ni})! \prod_n \frac{P(w_n|c)^{f_{ni}}}{f_{ni}!} \quad (2)$$

f_{ni} là tần số xuất hiện từ w_n trong văn bản t_i



Arcx4 of random Multinomial Naive Bayes (Arcx4-rMNB)

Đánh giá giải thuật MNB

- ▶ đơn giản, dễ cài đặt, tốc độ nhanh
- ▶ độ chính xác chưa cao khi **xử lý dữ liệu có số chiều rất lớn, mỗi chiều chứa ít thông tin**

Cải tiến MNB

- ▶ sử dụng **tập con các chiều ngẫu nhiên** từ tập ban đầu
- ▶ rMNB: xử lý tốt dữ liệu số chiều lớn (Breiman, 2001)
- ▶ sử dụng **kỹ thuật Arcx4** (Breiman, 1998)
- ▶ Arcx4-rMNB: cải thiện mô hình phân lớp dữ liệu số chiều lớn



Arcx4 of random Multinomial Naive Bayes (Arcx4-rMNB)

Giải thuật

- ▶ xây dựng tuân tự tập các mô hình $\{rMNB_i\}$
- ▶ mỗi $rMNB_i$; chỉ sử dụng ngẫu nhiên $n = \text{sqrt}(N)$ chiều
- ▶ mô hình xây dựng sau tập trung vào khắc phục lỗi từ các mô hình xây dựng trước đó
- ▶ phân lớp: bình chọn số đông từ $\{rMNB_i\}$
- ▶ hiệu quả: phân lớp rất chính xác dữ liệu số chiều lớn

Chuẩn bị dữ liệu

Tạo dữ liệu

- ▶ thu thập 1921 thư (1143 thư rác và 778 không phải thư rác)
- ▶ tiền xử lý với BoW (McCallum, 1998): bỏ qua các từ không chứa nhiều thông tin để nhận dạng thư rác, quy về từ gốc
- ▶ mô hình túi từ: bảng dữ liệu, 1921 phần tử (thư), 28719 thuộc tính (từ) và 2 lớp (thư rác hay không phải thư rác)
- ▶ nghi thức kiểm tra chéo 3-fold

Chuẩn bị chương trình

Chương trình

- ▶ giải thuật Arcx4-rMNB, MNB, Arcx4-MNB: C/C++
- ▶ máy học SVM: LibSVM (Chang & Lin, 2001)

Tiêu chí đánh giá

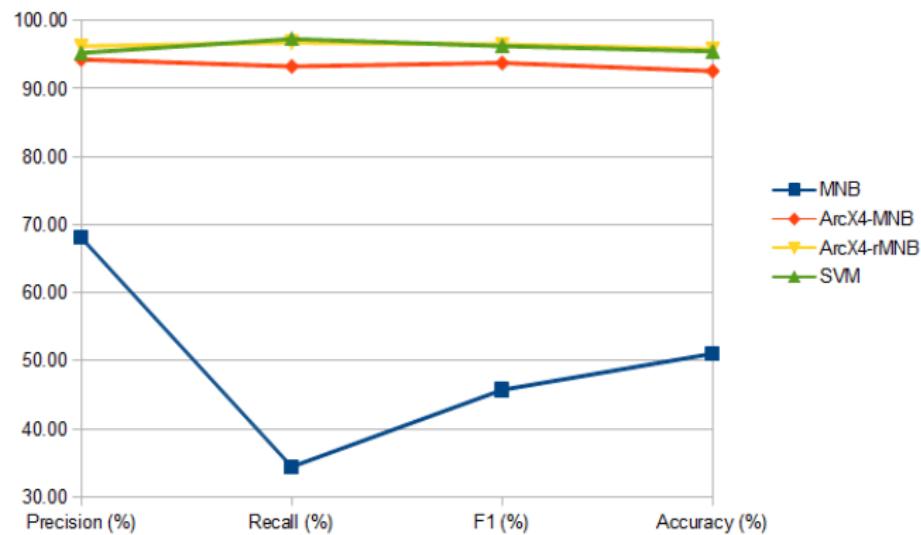
Recall, Precision, F1, Accuracy

- ▶ Recall: số thư rác được phát hiện đúng là thư rác chia cho tổng số thư rác
- ▶ Precision: số thư rác được phát hiện đúng là thư rác chia cho tổng số thư được dự báo là thư rác
- ▶ F1: trung bình điều hòa giữa Precision và Recall
- ▶ Accuracy: số thư được dự báo đúng chia cho tổng số thư



So sánh kết quả Precision, Recall, F1, Accuracy

Kết quả phân lớp thư rác



So sánh kết quả Precision, Recall, F1, Accuracy

Bảng: Kết quả phân lớp thư rác

	MNB	ArcX4-MNB	ArcX4-rMNB	LibSVM
TP	128	347	360	362
FN	244	25	12	10
TN	189	228	235	231
FP	60	21	14	18
Precision (%)	68.09	94.29	96.26	95.26
Recall (%)	34.41	93.28	96.77	97.31
F1 (%)	45.71	93.78	96.51	96.28
Accuracy (%)	51.05	92.59	95.81	95.49

Kết luận

Giải thuật mới Arcx4-rMNB cho phân lớp thư rác

- ▶ tiếp cận: học từ nội dung
- ▶ mô hình túi từ: số chiều dữ liệu lớn
- ▶ giải thuật phân lớp hiệu quả: Arcx4-rMNB
- ▶ đạt được độ chính xác đến **95.81%**
- ▶ cải thiện hơn **45%** độ chính xác của MNB (Lewis et al., 1994)
- ▶ phân lớp chính xác như SVM (Vapnik, 1995)

Phát triển

Tiếp tục nghiên cứu

- ▶ sưu tập thêm dữ liệu
- ▶ tích hợp vào hệ thống thư điện tử
- ▶ cải tiến tốc độ xử lý

Cám ơn & câu hỏi thảo luận

