

ITEC 301 - Machine Learning

Machine Problem 1

Spam Filtering using Naive Bayes Classifier

1. Introduction

Spam detection in email and messages like SMS or chat in general may be considered as a two-category classification problem. The first category represents legitimate messages or ham and the other unsolicited messages or spam. In this assignment you will be given the opportunity to make your own Naive Bayes classifier for filtering spam. For a successful implementation you need to train your classifier on some messages labeled as ham or spam and evaluate its performance on unseen test messages.

Recall that the Naive Bayes formula is given by

$$P(\omega|x_1, x_2, \dots, x_d) = \frac{\prod_{i=1}^d P(x_i|\omega)P(\omega)}{\sum_{\omega} \prod_{i=1}^d P(x_i|\omega)P(\omega)}$$

where d is the vocabulary size $|V|$. For our spam filter, the class conditional likelihood for word x_i is as follows:

$$P(x_i|\omega) = \sum_{D \in D_{\omega}} \frac{I(x_i \in D)}{|D_{\omega}|}$$

where $I(\cdot)$ is the indicator function, D_{ω} is the set of documents belonging to class ω in the training set and $|D_{\omega}|$ is its cardinality.

Note that when a word does not occur in a document class in the training data, it does not mean that it will never appear in any document of that class. We circumvent this problem by using Lambda Smoothing which uses a modified formula for the class conditional likelihood:

$$P(x_i|\omega) = \sum_{D \in D_{\omega}} \frac{I(x_i \in D) + \lambda}{|D_{\omega}| + \lambda|V|}$$

where $|V|$ is the vocabulary size. When $\lambda = 1$, this formula is called Laplace Smoothing. In order to avoid loss of precision during multiplication of the likelihoods, I suggest that you instead add the logarithm of the likelihoods and exponentiate the result.

2. Naive Bayes for the SMS Spam Collection Dataset

You will design a classifier for a subset of the SMS Spam Collection Dataset which is a dataset for benchmarking spam algorithms. This dataset is provided to you in the attached file. After downloading the dataset, you should see CSV files containing a TrainingData.csv (containing 3900 data) with values label and message, and a TestData.csv (containing 1672 data) with message only, without label. Use the TrainingData.csv to build your Naive Bayes Classifier and use the TestData.csv to test the classifier as demonstrated in the class, and discussed below.

3. Classifier Construction and Evaluation

1. Parse the documents in the training set (TrainingData.csv). For simplicity, first remove all the special characters, then define a word as any sequence of alphabetic characters [a-zA-Z] delimited by a white space in front and a white space, comma, or period at the end. Each data should be delimited by a new line. Form the vocabulary V of unique words in the training data, count their statistics and report the prior probabilities for spam and ham.

2. Construct and train a Naive Bayesian Classifier from the count statistics above. Use the Laplace Smoothing $\lambda = 1$.

3. Implement the code for classifying an unknown message and try it on the test set (TestData.csv).

4. Output the file ResultData.csv which contains the TestData.csv but with labels like the TrainingData.csv as shown below.

The TrainingData.csv in Excel should appear as:

2 ham message
3 ham Go until jurong point, crazy.. Available only in bugis n great world e buffet... Cine there got amore wat...
3 ham Ok lar... joking w/ u onl...
4 spam Free entry in a 2 wky comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entry confirmation(std txt rate)T&C's apply 08452810075over18's
5 ham U dun say so early hor... U c already then say...
6 ham Nah I don't think he goes to usf, he lives around here though
7 spam FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for it still? Tb ok! XxX std chgs to chg, ££1.50 to rcv
8 ham Even my brother is not able to speak with me. They treat me like aids patient.
9 ham As per your request 'Melle Melle (Oru Minnaminunginte Nuruugu Vettam)' has been set as your calltune for all Callers. Press *9 to copy your friends Calltune only.
10 spam WINNER!! As a valued network customer you have been selected to receive £4900 prize reward! To claim call 09061701461. Claim code KL341. Valid 12 hours only.
11 spam Had your chance to win CASH or a TV? U're entitled to Update to the latest colour mobiles with camera for Free! Call The Mobile Update Co Free on 08002986030
12 ham I'm going to be home soon and don't want to talk about this stuff anymore tonight, k? I've cried enough today.
13 spam SX chances to win CASH! From 100 to 20,000 pounds txt: CSH1 and text to 87575. Cost 150p/day, 6days, 16+ StdTxts apply Reply HL 4 info
14 spam URGENT! You have won a 1 week FREE membership in our £100,000 Prize Jackpot! Text the word: CLAIM to no: 810110 T&C www.dbuk.net/LCCLTD/POBOX4403/DNN/1A7WR/1W
15 ham I've been searching for the right words to use for this award. I promise I won't joke! I want to help for granted and will fulfil my promise. You have been wonderful and a blessing at all times.
16 ham I HAVE A DATE ON SUNDAY WITH WILL!!
17 spam XXXMobileMovieClub: To use your credit, click the WAP link in the next txt message or click here> <http://wap.xxxmobilemovieclub.com?n=QGKIGHJUGCL
18 ham Oh k...i'm watching here)
19 ham Eh u remember how 2 spell his name.... Yes i did. He v naughty make up till i v wet.
20 ham Fine if thata0s the way i feel. Thata0s the way its gots b
21 spam England v Macedonia - dont miss the goals/team news. Txt ur national team to 87077 eg ENGLAND to 87077 Try-WALES, SCOTLAND 4txt/1x1.20 POBOXxc36504/45WQ16+
22 ham Is that seriously how you spell his name?
23 ham PS4om going to try for 2 months ha ha only joking
24 ham So l_ pay first lar... Then when is de stock comin...
25 ham Aft i finish my lunch then i go str down lor. Ard 3 smth lor. U finish ur lunch already.
26 ham Fffffffh. Alright no way I can meet up with you sooner?
27 ham Just forced myself to eat a slice. I'm really not hungry tho. This sucks. Mark is getting worried. He knows I'm sick when I turn down pizza. Lol
28 ham Lol your always so convincing.
29 ham Did you catch the bus ? Are you frying an egg ? Did you make a tea? Are you eating your mom's left over dinner ? Do you feel my Love ?
30 ham I'm back &mp; we're packing the car now, I'll let you know if there's room
31 ham Ahhh. Wow. I vaguely remember that! What does it feel like! Lol
32 ham Wat that's still not all that clear, were you not sure about me being sarcastic or that that's why x doesn't want to live with us
33 ham Yeah he got hit in 2 and now we v apologize. n had fallen out and she was actin like spoilt child and he got caught up in that. Till 2! But we won't go there! Not doing too badly cheers. You?
34 ham K tell me anything about you.
35 ham For fear of fanning out of the of all that housework you just did? Quick have a cuppa
36 spam Thanks for your subscription to Ringtone UK your mobile will be charged ££5/month Please confirm by replying YES or NO. If you reply NO you will not be charged
37 ham Yup... Ok i go home look at the timings then i msg u_ again... Xuhui going to learn on 2nd may too but her lesson is at 8m

```

label,message
nam,"Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat..."
nam,"Ok lar... Joking wif u cond..."
nam,"Free entry in 2 a wily comp to win FA Cup final tixs 21st May 2005. Text FA to 87121 to receive entry question(std tx rate)&C; say
86452310050ver18"
nam,"U dun say so early hor... U c already then say..."
nam,"ah i think he will go to us, he lives around here though"
nam,"FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for it still? Tb ok! XxX std chgs to send, 3ei.58 to rcv"
nam,"Even my brother is not like to speak with me. Treat me like a sick patient."
nam,"Press your request 'Melle Melle (Oru Minaminimunte Nuruugu Vettam)' has been set as your caller tune for all Callers. Press *9 to copy your friends
Callertune"
nam,"Congrats! As a valued network customer you have been selected to receive a £900 prize reward! To claim call 08061701661. Claim code KL341. Valid 12
months only."
nam,"Have your mobile 11 months or more? U R entitled to Update to the latest colour mobiles with camera for Free! Call The Mobile Update Co FREE on
08002956830"
nam,"I'm gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough today."
nam,"SIX thanks for winn CASH! From 100 to 20,000 pounds Txt CASH1 and send to 87575. Cost 150p/day, 6days, 16p. Tansds apply Reply HL 4 info"
nam,"URGENT! You have won a 1 WEEK FREE membership in our £100,000 Prize Jackpot! Txt the word: CLAIM to No: 81418 Txt: www.dub.net LCLTD POBOX
440109"
nam,"I've been searching for the right words to thank you for this birthday. I promise i will take you time help for granted and will fulfil my promise. You
have been wonderful and a blessing at all times."
nam,"I HAVE A DATE ON SUNDAY WITH MIL!! I'm gonna have a blast! Txt: XXXX000b11eMovieClub! to use your credit, click the MAP link in the next txt message or click here>> http://wap. xxx000b11eMovieClub.com"
nam,"Q=UKG(H)GJGCL"
nam,"Oh K...i'm watching here!"
nam,"Eh i remember how 2 spell his name... Yes i did. He v naughty make until i v wet."
nam,"Fine if that's the way u feel. That's the way ita gots b"
nam,"England v Macedonia - dont miss the goals/team news. Txt ur national team to 87077 eg ENGLAND to 87077 Try:WALES, SCOTLAND 4txt/1ht.10
pobxxbxxkxkxkxkx10"
nam,"Is that seriously how you spell his name?"
nam,"Illian going to try for 2 months ha he only joking"
nam,"So i pay first lar... Then when is da stock coming?"
nam,"Aft i finish my lunch then i go str down lor. Ard 3 seth lor. u finish ur lunch already?"
nam,"Ffffffffrhh... Alright no way i can meet up with you sooner!"
nam,"Just forced myself to eat a slice. I'm really not hungry tho. This sucks. Mark is getting worried. He knows I'm sick when i turn down pizza. Lol"
nam,"lol you always so convincing."

```

1 message

2 That depends. How would you like to be treated? :)

3 Right on brah, see you later

4 Waiting in car at c my mum lor. U leh? Reach home already?

5 Your 2004 account for 07000000000 shows 786 unredeemed points. To claim call 08719181259 Identifier code: XXXXXX Expires 26.03.05

6 Do you want a new video handset? 750 anytime any network mins? Half Price Line Rental? Camcorder? Reply or call 08000930705 for delivery tomorrow

7 Went fast asleep take care.

8 No that just means you have a fat head

9 Sounds like a plan! Cardiff is still here and still cold! I'm sitting on the radiator!

10 Serious? What like proper tongue'd her

11 She's good. She was wondering if you want say hi but she's smiling now. So how are you coping with the distance

12 How i noe... She's in da car now... Later then car lar... I'm wearing shorts...

13 You have an important customer service announcement. Call FREEPHONE 0800 542 0825 now!

14 Yeah whatever lol

15 Today is ACCEPT DAY. U Accept me as? Brother Sister Lover Dear1 Best1 Clos1 Lubefnd1 Jstfrnd Ufpartr Belovd Suthwart Bstfrnd No rply means enemy

16 Ard 530 lor. I ok then message u. lor.

17 Ok. C u then.

18 Eh ur laptop got no stock iell... He say mon muz come again to take a look got a not...

19 No need to qk... I too bored iztz y suddenly tk of this...

20 I wish I don't think its gonna snow that much. But it will be more than those flurries we usually get that melt before they hit the ground. Ekl! We haven't had snow since before I was even born!

21 FREE-Registration Reply REAl or PDI eg REAl1. P.ushButton. D.ontHt a3. BabyGoodbye4. GoldGidder4. We'llBurnIn! Use one FREE and 6 more when you join for 4&3/wk

22 I don't change that chance either. (because i want 2 concentrate in my educational career in leaving here...)

23 Oh really? perform, write a paper, go to a movie And don't home by midnight, huh?

24 Okay lor... Will they still let us go a not ah? Cos they will not know until later. We drop our cards into the box right?

25 How? Iztz still raining?

26 As if i wasn't having enough trouble sleeping.

27 I havent add u yet right..

28 Lol... I really need to remember to eat when i'm drinking but i do appreciate you keeping me company that night babe "smiles"

29 Babe 7 i lost you... Will you try rebooting?

30 Yes. Hrn you cent shs.

31 Oh i... gotta go home by myself. Cos i'll b going out shopping 4 my frens present.

32 Noooooo I'm gonna be bored to death all day. Cable and internet outage.

33 So! Any amount i can get pls.

34 Playin space poker, u?

35 How come guying go n tell her? Then u told her?

36 You need to get up. Now.

```

message
That depends. How would you like to be treated? :)
"Right on brah, see you later"
Waiting in a car 4 my mum XXXX. U leh? Reach home already?
Your 2004 account of 0700XXXXXXX shows 786 unredeemed points. To claim call 08719181259 Identifier code: XXXXX Expires 26.03.05
You've ordered a new 2004 account 750 anytime any network mins? Half Price Line Rental? Canceled? Reply or call 08000038705 for delivery tomorrow
Went fast asleep dear,take care.
No that just means you have a fat head
Sounds like a plan! Cardiff is still here and still cold! I'm sitting on the radiator!
Serious? What like proper tongued her
She.s good. She was wondering if you want say hi but she.s smiling now. So how are you coping with the long distance
Wow i noe... She's in da car now... Later then c lar... I'm wearing shorts...
You have an important customer service announcement. Call FREE0800 8880 542 8825 now!
Yeah whatever lol
Today is ACCEPT DAY.. U Accept me @ Brother Sister Lover Dear! Best! Clos! Lb!lefrnd Jstrfrnd Lifpartnr Belovd Swthartst Bstfrnd No rplyn mens
Cute!
Ad 530 10r. I ok then message @_lar.
Ok. C u then.
Es un laptop got no stock la... He say mon moe cos again to take a look c got a not...
No need to be ki... BB too bored zzzit y suddenly thk of this..
I wish! I don't think it's gonna stop that much. But it will be more than those flurries we usually get that melt before they hit the ground. Eek! We haven't had snow since it's said; before I was even born!
F&S:Kingtonel Reply REAL or POLY or REAL 1. Pushbutton 2. DonChica 3. BabyGoodbye 4. GoldDigger 5. WeBeBurnin 1st tone FREE and 6 more when u join for B 3/wk
/0 1 thing) Change that sentence into :Because i want 2 concentrate in my educational career in leaving here..!~~~~~
"on really? perform, write a paper, go to a movie And be home by midnight, huh?"
Okay Now... Will they still let us go a not ah? Coz they will not know until later. We drop our cards into the box right?
Wow? Izit still raining?
As if i wasn't having enough trouble sleeping.
I havent add u yet right.
lol ... I really need to remember to eat when I'm drinking but I do appreciate you keeping me company that night babe "smiles"
Babe ? I lost you ??.. Will you try rebooting ?
Yes. Hign you cant cha.
I thk @_gotta go home by urself. Cos i'll b going out shopping 4 my frrens present.
Nooooooo if i gonna be bored to death all day. Cable and internet outage.
Sori any amount i can get pls.
"Playin space poker, w?"

```

The `ResultData.csv` will contain the `TestData.csv` but with a label field like the `TrainingData.csv`. Focus on creating `ResultData.csv` for now. After this activity, you will test the precision and recall of your generated classifier in the next activity to be posted.