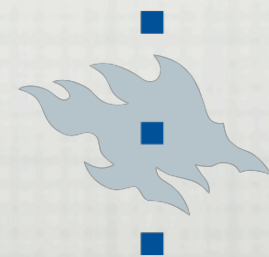


# JOHDATUS TEKÖÄLYYN

TEEMU ROOS

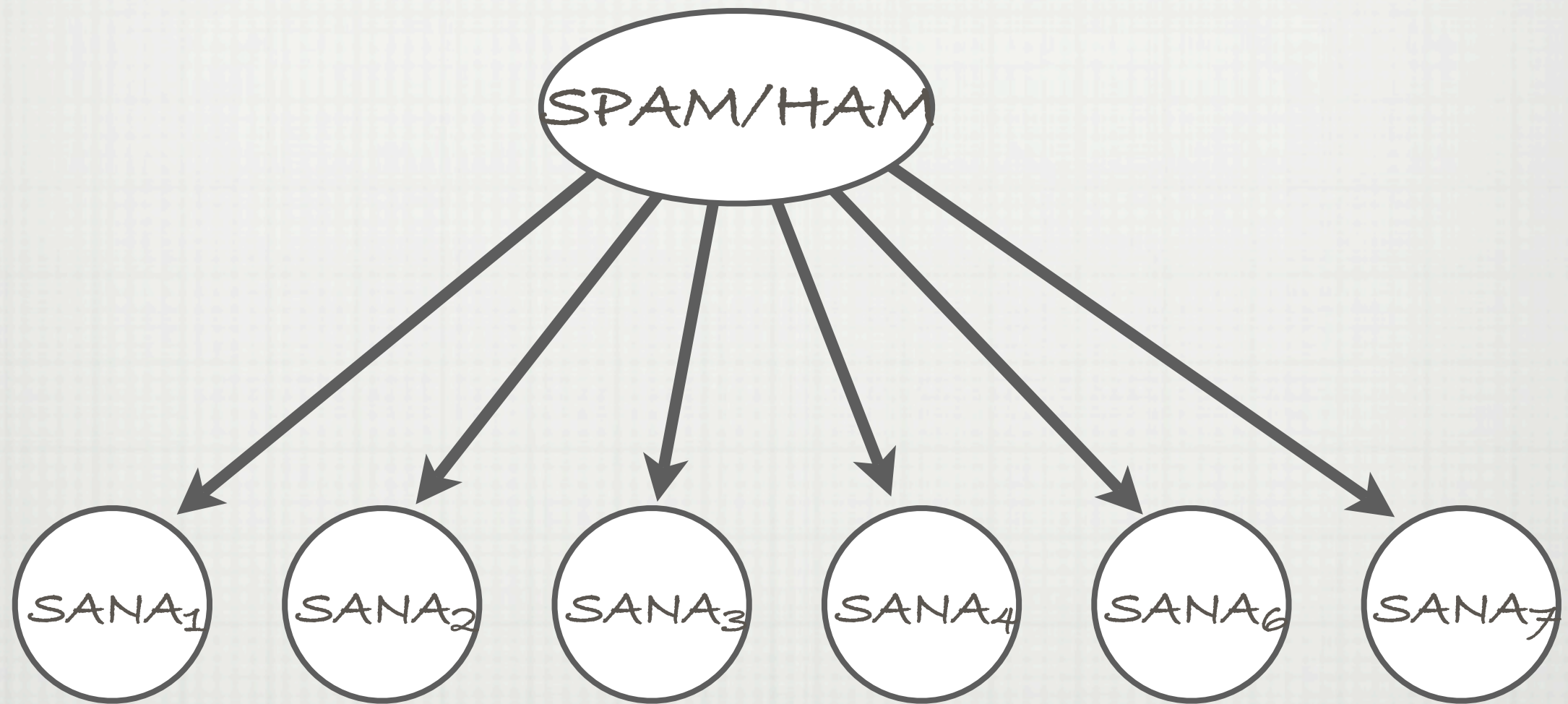


HELSINGIN YLIOPISTO



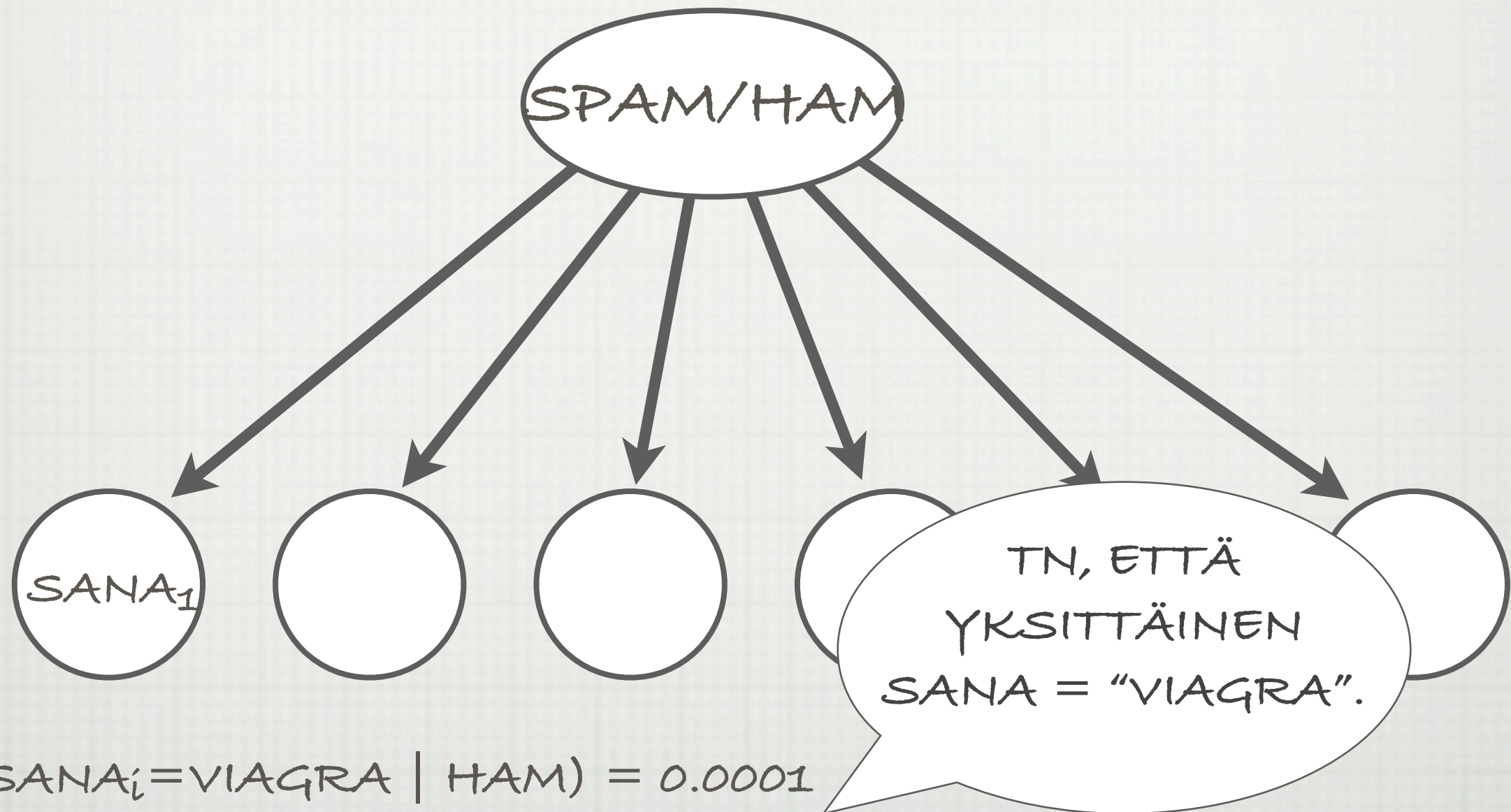
# NAIVI BAYES

---



# NAIVI BAYES

---



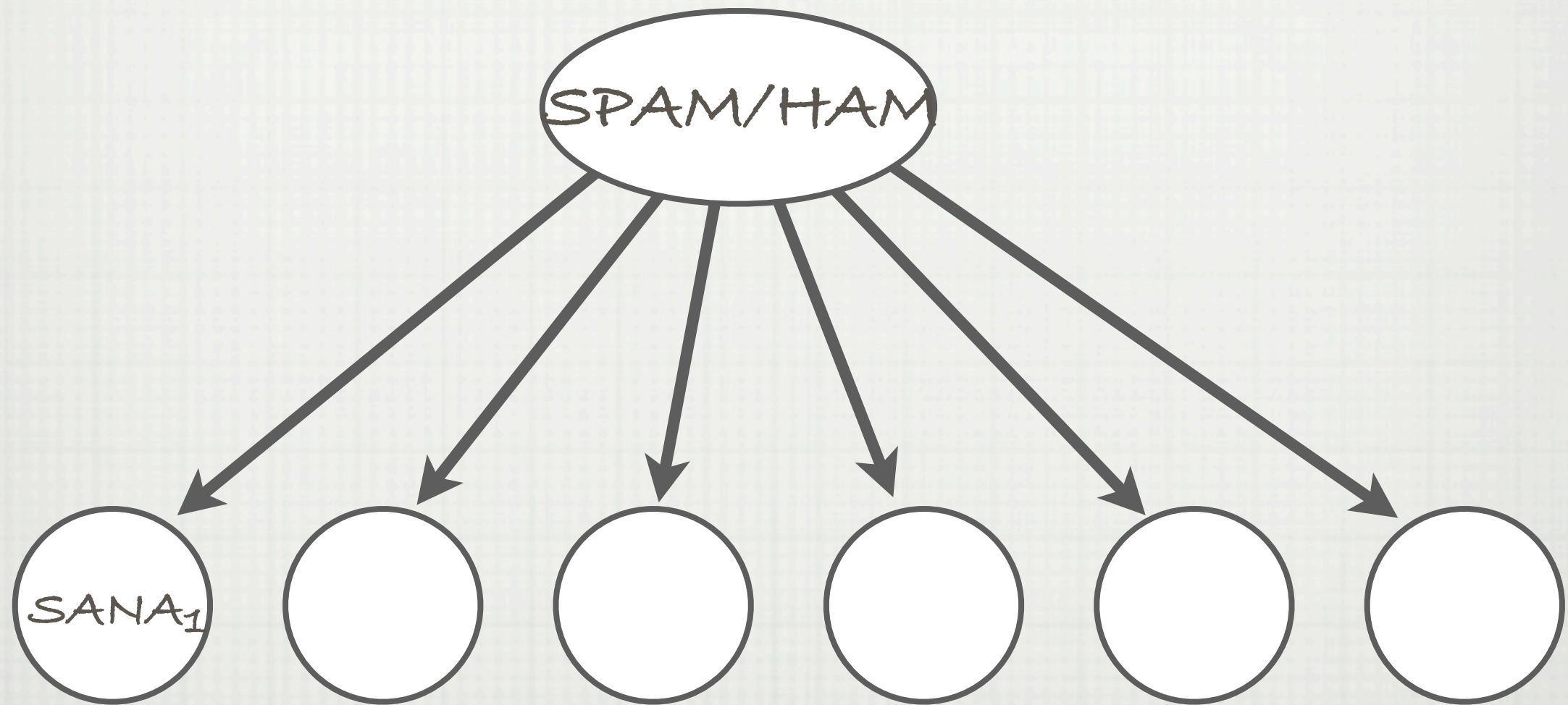
$$P(SANA_i = VIAGRA \mid HAM) = 0.0001$$

$$P(SANA_i = VIAGRA \mid SPAM) = 0.002$$



# NAIVI BAYES

---



$$P(SANA_i = \$ \mid HAM) = 0.0002$$

$$P(SANA_i = \$ \mid SPAM) = 0.005$$

# NAIVI BAYES

---

FROM: "MARGARETTA NITA" <MARGUERITesebrina@wmle.com>  
SUBJECT: **SPECIAL OFFER : VIAGRA ON SALE AT \$1.38 !!!**  
X-BOGOSITY: YES, TESTS=BOGOFILTER, SPAMICITY=0.99993752,  
VERSION=2011-08-29  
DATE: MON, 26 SEP 2011 21:52:26 +0300  
X-CLASSIFICATION: **JUNK - AD HOC SPAM DETECTED (CODE = 73)**

**SPECIAL OFFER : VIAGRA ON SALE AT \$1.38 !!!**

COMPARE THE BEST ONLINE PHARMACIES TO **BUY VIAGRA**. ORDER  
**VIAGRA** ONLINE WITH HUGE **DISCOUNT**.

MULTIPLE BENEFITS INCLUDE **FREE SHIPPING**, REORDER **DISCOUNTS**,  
BONUS PILLS

[HTTP://RXPHARMACYCVS.RU](http://rxpharmacycvs.ru)



# JOHDATUS TOD.NÄK.LASKENTAAN

---

1.  $P(A,B,C) = P(A) P(B|A) P(C|A,B)$  // KETJUSÄÄNTÖ
2.  $P(A) = P(A,B) + P(A, \neg B)$  // "MARGINALISOINTI"
3.  $P(A|B) = P(A,B) / P(B)$  // EHDOLLINEN TN.

$$\underline{P(A)} P(B|A) = \underline{P(B)} P(A|B) ??? \text{ (KS. SÄÄNTÖ 1)}$$

4.  $P(B|A) = P(B) P(A|B) / P(A)$  // BAYESIN KAAVA
5.  $A \perp B \Rightarrow P(A|B) = P(A)$  // RIIPPUMATTOMUUS



# NAIVI BAYES

---

PÄÄTTELY:

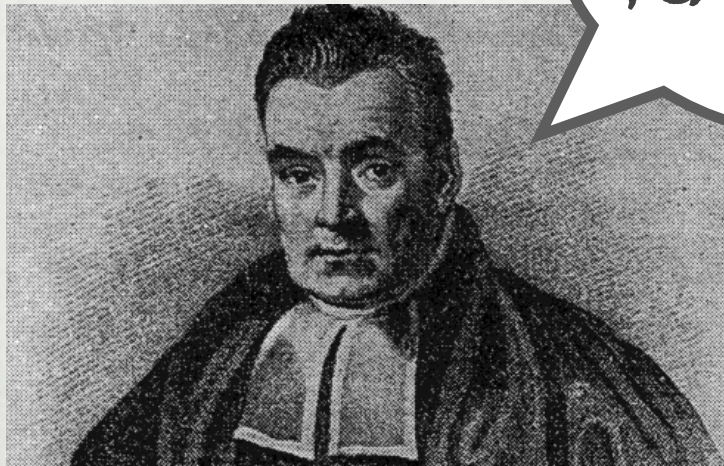
1.  $P(\text{SPAM}) = 0.5$

$P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM})$

2.  $P(\text{SPAM} | \text{SANA}_1 = \text{VIAGRA}) = \frac{\text{-----}}{\text{-----}}$

$P(\text{SANA}_1 = \text{VIAGRA})$

BAYESIN  
KAAVA!





# NAIVI BAYES

---

PÄÄTTELY:

1.  $P(\text{SPAM}) = 0.5$

$$P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM})$$

$$2. P(\text{SPAM} | \text{SANA}_1 = \text{VIAGRA}) = \frac{P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM})}{P(\text{SANA}_1 = \text{VIAGRA})}$$

3.  $P(\text{SPAM} | \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS})$

$$P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS} | \text{SPAM})$$

$$= \frac{P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS} | \text{SPAM})}{P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS})}$$

4.  $P(\text{SPAM} | \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM})$

$$P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM} | \text{SPAM})$$

$$= \frac{P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM} | \text{SPAM})}{P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM})}$$



# NAIVI BAYES

PÄÄTTELY:

1.  $P(\text{SPAM}) = 0.5$

$P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM})$

2.  $P(\text{SPAM} | \text{SANA}_1 = \text{VIAGRA}) =$  -----

$P(\text{SANA}_1 = \text{VIAGRA})$

3.  $P(\text{SPAM} | \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}) =$

$P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM}) P(\text{SANA}_2 = \text{IS} | \text{SPAM})$

=

-----

$P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS})$

PARI TÄRKEÄÄ  
JUTTUA TÄSTÄ...

4.  $P(\text{SPAM} | \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{KATU}) =$

$P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM}) P(\text{SANA}_2 = \text{IS} | \text{SPAM}) P(\text{SANA}_3 = \text{KATU} | \text{SPAM})$

=

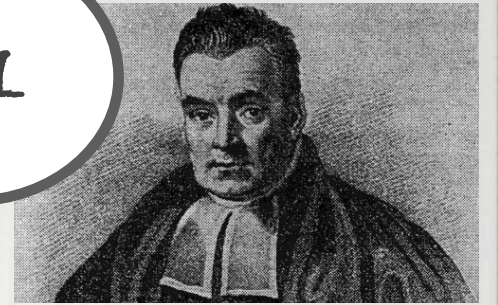
-----  
 $P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{KATU})$





# NAIVI BAYES

#1



PÄÄTTELY:

$$4. P(\text{SPAM} \mid \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM}) \\ P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM} \mid \text{SPAM}) \\ =$$

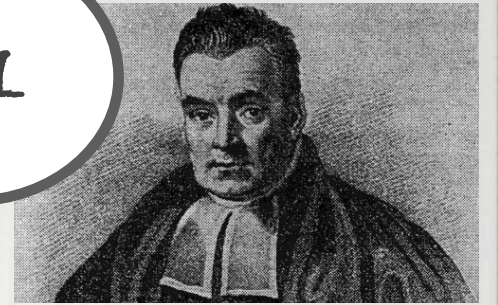
$$P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM})$$

$$P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM}) \\ = P(\text{SPAM}, \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM}) \\ + P(\neg \text{SPAM}, \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM})$$



# NAIVI BAYES

#1



PÄÄTTELY:

$$4. P(\text{SPAM} \mid \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM}) \\ P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM} \mid \text{SPAM}) \\ =$$

$$P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM})$$

$P(\text{EVIDENSSI})$

$= P(\text{SPAM}, \text{EVIDENSSI})$

// MARGINALISOINTI

$+ P(\neg \text{SPAM}, \text{EVIDENSSI})$

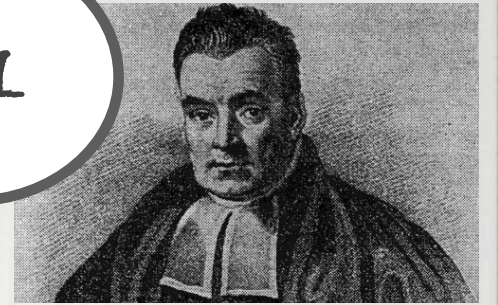
$P(\text{SPAM})P(\text{EVIDENSSI} \mid \text{SPAM})$

$$P(\text{SPAM} \mid \text{EVIDENSSI}) = \frac{P(\text{SPAM})P(\text{EVIDENSSI} \mid \text{SPAM})}{P(\text{SPAM}, \text{EVIDENSSI}) + P(\neg \text{SPAM}, \text{EVIDENSSI})}$$



# NAIVI BAYES

#1



PÄÄTTELY:

$$4. P(\text{SPAM} \mid \text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM}) \\ P(\text{SPAM}) P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM} \mid \text{SPAM}) \\ =$$

$$P(\text{SANA}_1 = \text{VIAGRA}, \text{SANA}_2 = \text{IS}, \text{SANA}_3 = \text{ALGORITHM})$$

$P(\text{EVIDENSSI})$

$$= P(\text{SPAM}, \text{EVIDENSSI}) \\ + P(\neg \text{SPAM}, \text{EVIDENSSI})$$

$$P(\text{SPAM})P(\text{EVIDENSSI} \mid \text{SPAM})$$

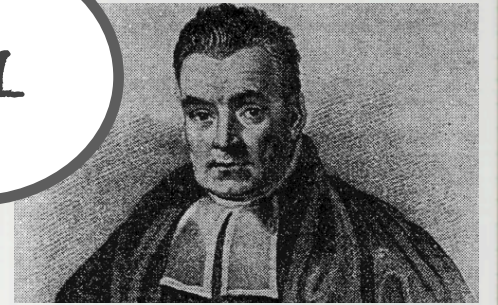
$$P(\text{SPAM} \mid \text{EVIDENSSI}) = \frac{P(\text{SPAM}, \text{EVIDENSSI})}{P(\text{SPAM}, \text{EVIDENSSI}) + P(\neg \text{SPAM}, \text{EVIDENSSI})} \\ \frac{P(\text{SPAM})P(\text{EVIDENSSI} \mid \text{SPAM})}{P(\text{SPAM})P(\text{EVIDENSSI} \mid \text{SPAM}) + P(\neg \text{SPAM})P(\text{EVIDENSSI} \mid \neg \text{SPAM})}$$

$$P(\neg \text{SPAM} \mid \text{EVIDENSSI}) = \frac{P(\neg \text{SPAM}, \text{EVIDENSSI})}{P(\text{SPAM}, \text{EVIDENSSI}) + P(\neg \text{SPAM}, \text{EVIDENSSI})} \\ \frac{P(\neg \text{SPAM})P(\text{EVIDENSSI} \mid \neg \text{SPAM})}{P(\text{SPAM})P(\text{EVIDENSSI} \mid \text{SPAM}) + P(\neg \text{SPAM})P(\text{EVIDENSSI} \mid \neg \text{SPAM})}$$



# NAIVI BAYES

#1



PÄÄTTELY:

$$4. P(\text{SPAM} \mid \text{SANA}_1=\text{VIAGRA}, \text{SANA}_2=\text{IS}, \text{SANA}_3=\text{ALGORITHM}) \\ P(\text{SPAM}) P(\text{SANA}_1=\text{VIAGRA}, \text{SANA}_2=\text{IS}, \text{SANA}_3=\text{ALGORITHM} \mid \text{SPAM}) \\ =$$

$$P(\text{SANA}_1=\text{VIAGRA}, \text{SANA}_2=\text{IS}, \text{SANA}_3=\text{ALGORITHM})$$

$$P(\text{EVIDENSSI}) \\ = P(\text{SPAM}, \text{EVIDENSSI}) \\ + P(\neg\text{SPAM}, \text{EVIDENSSI})$$

$$P(\text{SPAM} \mid \text{EVIDENSSI})$$

$$P(\neg\text{SPAM} \mid \text{EVIDENSSI})$$

$$P(S \mid E) \quad P(S)P(E \mid S)$$

$$P(\neg S \mid E) \quad P(\neg S)P(E \mid \neg S)$$

$$P(\text{SPAM})P(\text{EVIDENSSI} \mid \text{SPAM})$$

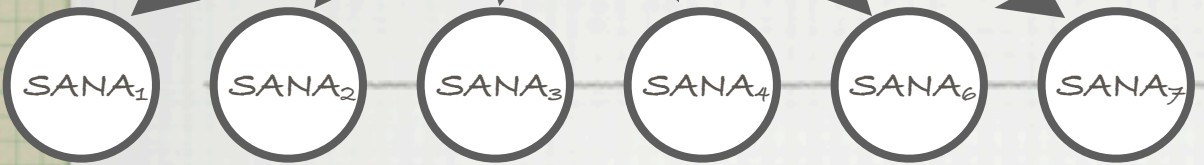
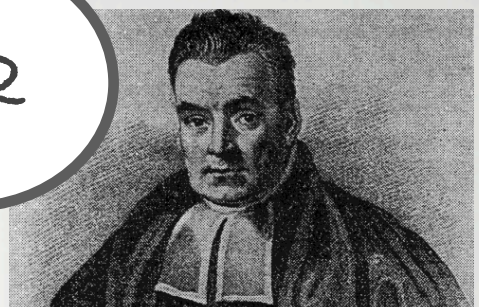
$$P(\neg\text{SPAM})P(\text{EVIDENSSI} \mid \neg\text{SPAM})$$



SPAM/HAM

# NAIVI BAYES

#2



PÄÄTTELY:

$$P(SANA_1=VIAGRA, SANA_2=IS, SANA_3=ALGORITHM|SPAM)$$

$$= P(SANA_1=VIAGRA|SPAM) \quad // \text{KETJUSÄÄNTÖ}$$

$$P(SANA_2=IS | \cancel{SANA_1=VIAGRA}, SPAM)$$

$$P(SANA_3=ALGORITHM | \cancel{SANA_1=VIAGRA}, \cancel{SANA_2=IS}, SPAM)$$

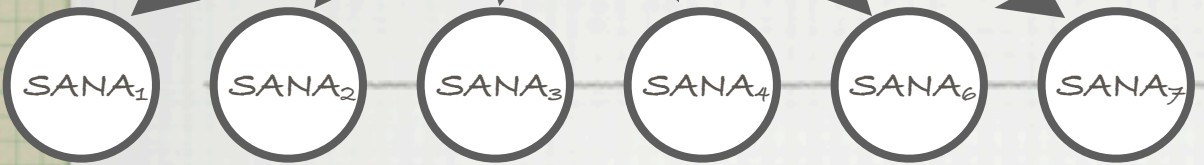
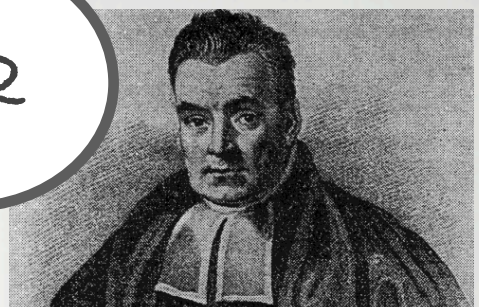
// RIIPPUMATTOMUUS



SPAM/HAM

# NAIVI BAYES

#2



PÄÄTTELY:

$$P(SANA_1=VIAGRA, SANA_2=IS, SANA_3=ALGORITHM|SPAM)$$

$$= P(SANA_1=VIAGRA|SPAM) \\ P(SANA_2=IS|SPAM) \\ P(SANA_3=ALGORITHM|SPAM)$$

$$\frac{P(S|E) \quad P(S)P(E|S)}{P(\neg S|E) \quad P(\neg S)P(E|\neg S)}$$

$$P(SANA_1=VIAGRA, SANA_2=IS, SANA_3=ALGORITHM|\neg SPAM)$$

$$= P(SANA_1=VIAGRA|\neg SPAM) \\ P(SANA_2=IS|\neg SPAM) \\ P(SANA_3=ALGORITHM|\neg SPAM)$$

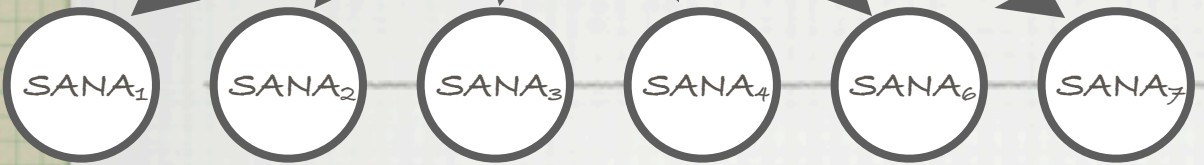
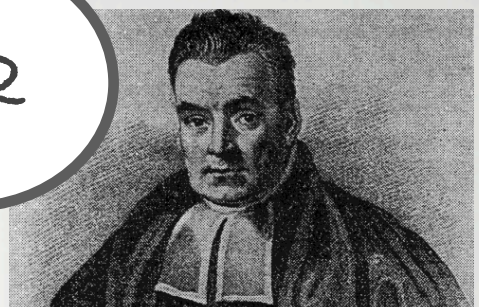
// KETJUSÄÄNTÖ



SPAM/HAM

# NAIVI BAYES

#2



PÄÄTTELY:

$$P(\text{SPAM} | \text{EVIDENSSI}) / P(\neg \text{SPAM} | \text{EVIDENSSI})$$

$$= P(\text{SPAM}) / P(\neg \text{SPAM})$$

$$P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM}) / P(\text{SANA}_1 = \text{VIAGRA} | \neg \text{SPAM}) > 1$$

$$P(\text{SANA}_2 = \text{IS} | \text{SPAM}) / P(\text{SANA}_2 = \text{IS} | \neg \text{SPAM}) = 1$$

$$P(\text{SANA}_3 = \text{ALGORITHM} | \text{SPAM}) / P(\text{SANA}_3 = \text{ALGORITHM} | \neg \text{SPAM}) < 1$$

...

...

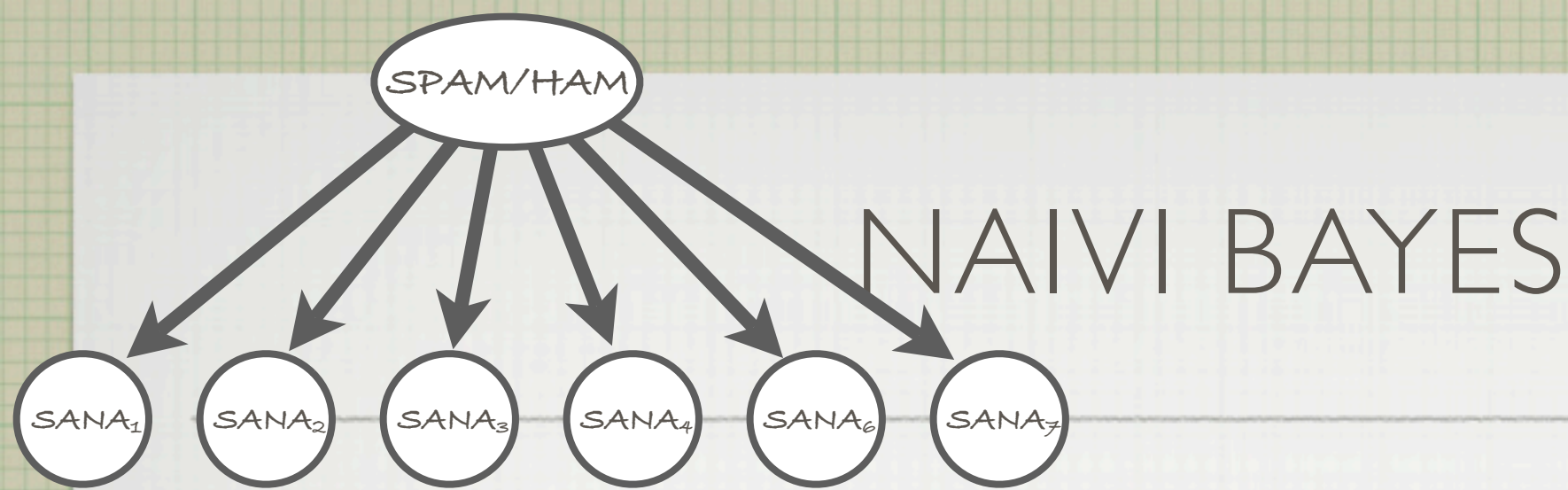
...

> 1 => "SPAM"

< 1 => "HAM"

**OSAMÄÄRÄ MÄÄRÄÄ!**





YHTEENVETO TOISTAISEKSI:

\* TARVITAAN:

- "PRIORIJAKAUMA"  $P(\text{SPAM}) = 0.\_\_\_$

- "LUOKKAEHDOLLISET" JAKAUMAT

$$P(\text{SANA}_i = \text{VIAGRA} | \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{VIAGRA} | \neg \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{IS} | \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{IS} | \neg \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{ALGORITHM} | \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{ALG.} | \neg \text{SPAM}) = 0.\_\_\_$$

\* OLETETAAN ETTÄ  $P(\text{SANA}_i | \text{SANA}_j, \text{SPAM}) = P(\text{SANA}_i | \text{SPAM})$   
(EHDOLLINEN RIIPPUMATTOMUUS)

\* OLENNAINEN ON OSAMÄÄRÄ

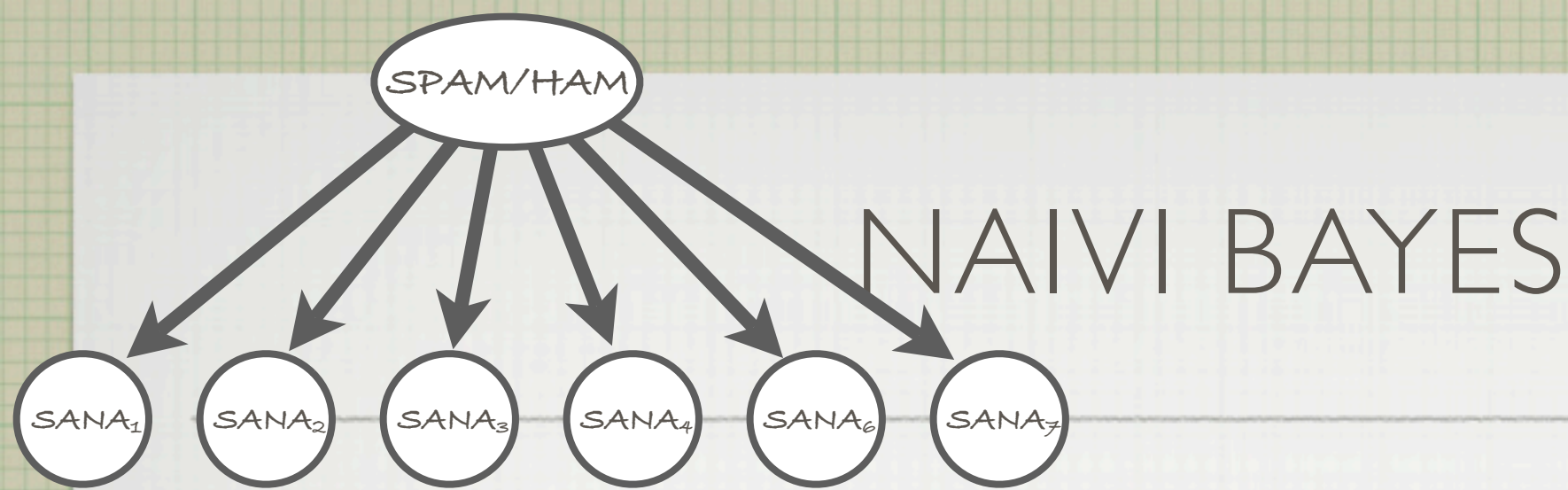
$$P(\text{SANA}_i = \text{VIAGRA} | \text{SPAM})$$

(OTETAAN NÄIDEN TULO)

---


$$P(\text{SANA}_i = \text{VIAGRA} | \neg \text{SPAM})$$





PSEUDOKOODINA:

SPAMICITY(Viesti, P):

$$\text{Odds} = \frac{P(\text{Spam})}{P(\text{noSpam})} \quad // \quad \frac{P(\text{Spam})}{P(\text{Spam}) + P(\text{noSpam})} = 1$$

for each Sana in Viesti

$$\text{Odds} = \text{Odds} * \frac{P(\text{Sana\_Spam}(\text{Sana}))}{P(\text{Sana\_noSpam}(\text{Sana}))}$$

return(Odds)

PÄÄTTELY:

$$P(\text{SPAM} | \text{EVIDENSSI}) / P(\neg \text{SPAM} | \text{EVIDENSSI})$$

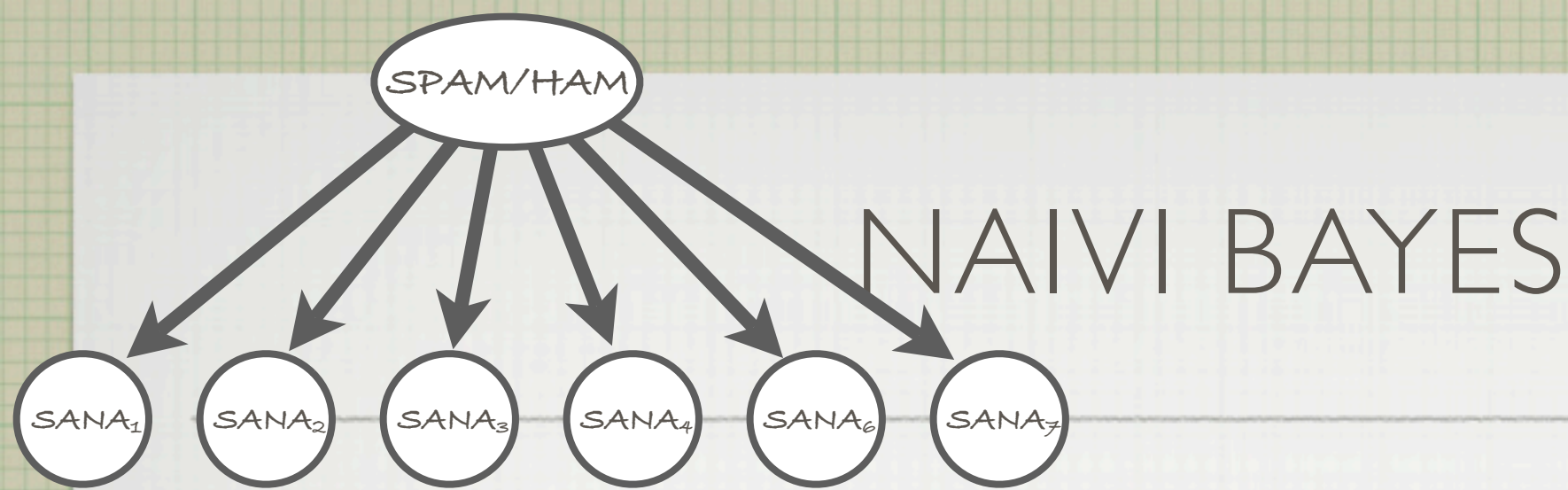
$$= \frac{P(\text{SPAM})}{P(\neg \text{SPAM})}$$

$$\frac{P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM})}{P(\text{SANA}_1 = \text{VIAGRA} | \neg \text{SPAM})}$$

$$\frac{P(\text{SANA}_2 = \text{IS} | \text{SPAM})}{P(\text{SANA}_2 = \text{IS} | \neg \text{SPAM})}$$

$$\frac{P(\text{SANA}_3 = \text{ALGORITHM} | \text{SPAM})}{P(\text{SANA}_3 = \text{ALG} | \neg \text{SPAM})}$$





PSEUDOKOODINA:

SPAMICITY(Viesti, P):

$$\text{Odds} = \frac{P(\text{Spam})}{P(\text{noSpam})} \quad // \quad \frac{P(\text{Spam})}{P(\text{Spam}) + P(\text{noSpam})} = 1$$

for each Sana in Viesti

$$\text{Odds} = \text{Odds} * \frac{P(\text{Sana\_Spam}(\text{Sana}))}{P(\text{Sana\_noSpam}(\text{Sana}))}$$

return(Odds)

PÄÄTTELY:

$$\frac{P(\text{SPAM} | \text{EVIDENSSI})}{P(\neg \text{SPAM} | \text{EVIDENSSI})}$$

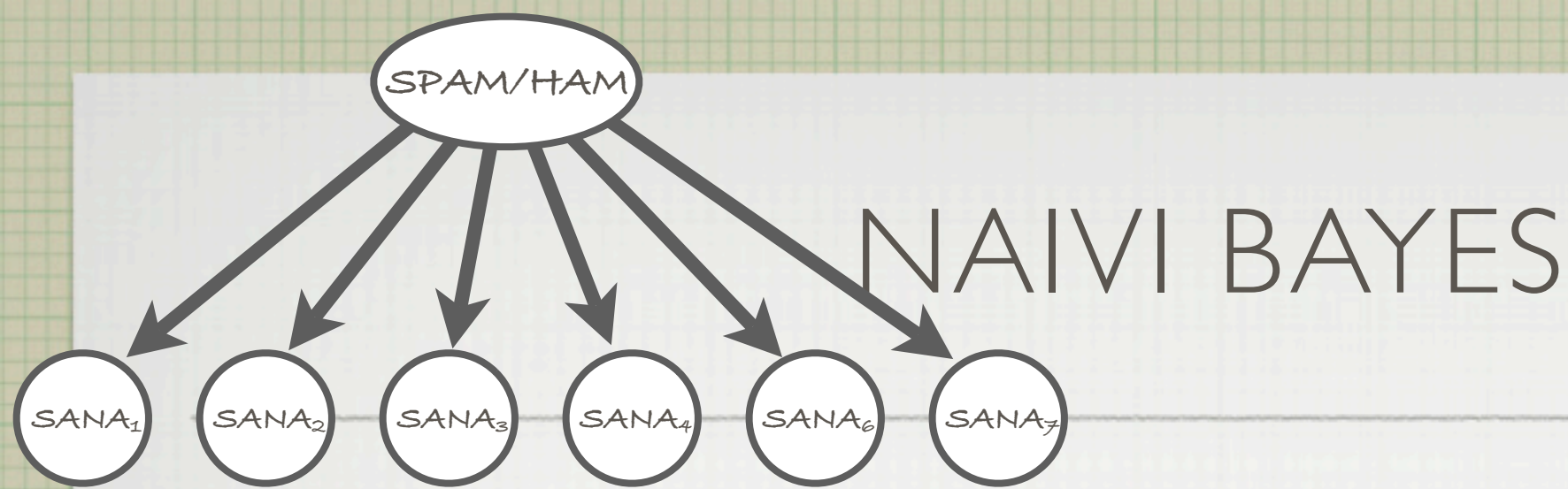
$$= \frac{P(\text{SPAM})}{P(\neg \text{SPAM})}$$

$$\frac{P(\text{SANA}_1 = \text{VIAGRA} | \text{SPAM})}{P(\text{SANA}_1 = \text{VIAGRA} | \neg \text{SPAM})}$$

$$\frac{P(\text{SANA}_2 = \text{IS} | \text{SPAM})}{P(\text{SANA}_2 = \text{IS} | \neg \text{SPAM})}$$

$$\frac{P(\text{SANA}_3 = \text{ALGORITHM} | \text{SPAM})}{P(\text{SANA}_3 = \text{ALG} | \neg \text{SPAM})}$$





PSEUDOKOODINA:

SPAMICITY(Viesti, P):

$$\text{Odds} = \text{P.Spam} / \text{P.noSpam} \quad // \quad \text{P.Spam} + \text{P.noSpam} = 1$$

for each Sana in Viesti

$$\text{Odds} = \text{Odds} * \text{P.Sana\_Spam}(\text{Sana}) / \text{P.Sana\_noSpam}(\text{Sana})$$

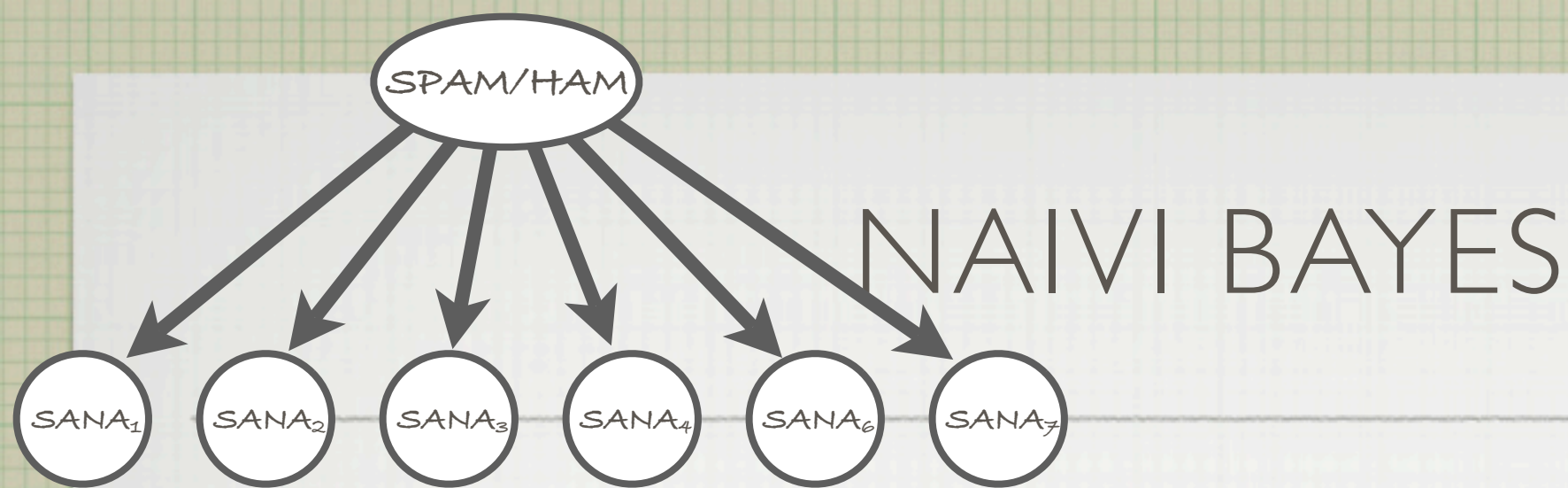
return(Odds)

JOS SPAMICITY(Viesti, P) >1, LUOKITTELE VIESTI SPAMIKSI

JOS SPAMICITY(Viesti, P) <1, LUOKITTELE VIESTI HAMIKSI

JOS SPAMICITY(Viesti, P) =1, EN TIEDÄ





PSEUDOKOODINA:

SPAMICITY(Viesti, P):

$$\text{Odds} = \frac{P.\text{Spam}}{P.\text{noSpam}} \quad // \quad \frac{P.\text{Spam}}{P.\text{Spam} + P.\text{noSpam}} = 1$$

for each Sana in Viesti

$$\text{Odds} = \text{Odds} * \frac{P.\text{Sana\_Spam}(\text{Sana})}{P.\text{Sana\_noSpam}(\text{Sana})}$$

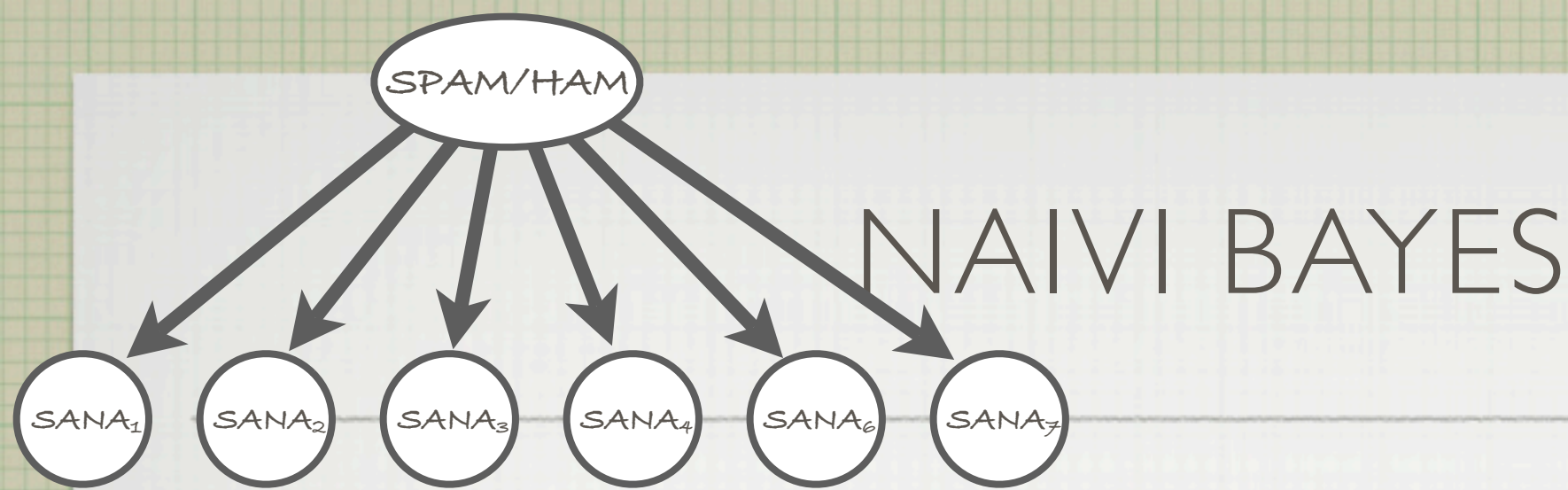
return(Odds)

JOS SPAMICITY(Viesti, P) >  $1 + \alpha$ , LUOKITTELE VIESTI SPAMIKSI

JOS SPAMICITY(Viesti, P) <  $1 - \beta$ , LUOKITTELE VIESTI HAMIKSI

MUUTEN, EN TIEDÄ





PSEUDOKOODINA:

$$\text{LOG}(A*B) = \text{LOG}(A) + \text{LOG}(B)$$

SPAMICITY(Viesti, P):

$$\text{Odds} = \text{P.Spam} / \text{P.noSpam} \quad // \quad \text{P.Spam} + \text{P.noSpam} = 1$$

for each Sana in Viesti

$$\text{Odds} = \text{Odds} * \text{P.Sana\_Spam}(Sana) / \text{P.Sana\_noSpam}(Sana)$$

return(Odds)

KÄYTÄNNÖN ONGELMA: ALI-JA YLIVUODOT

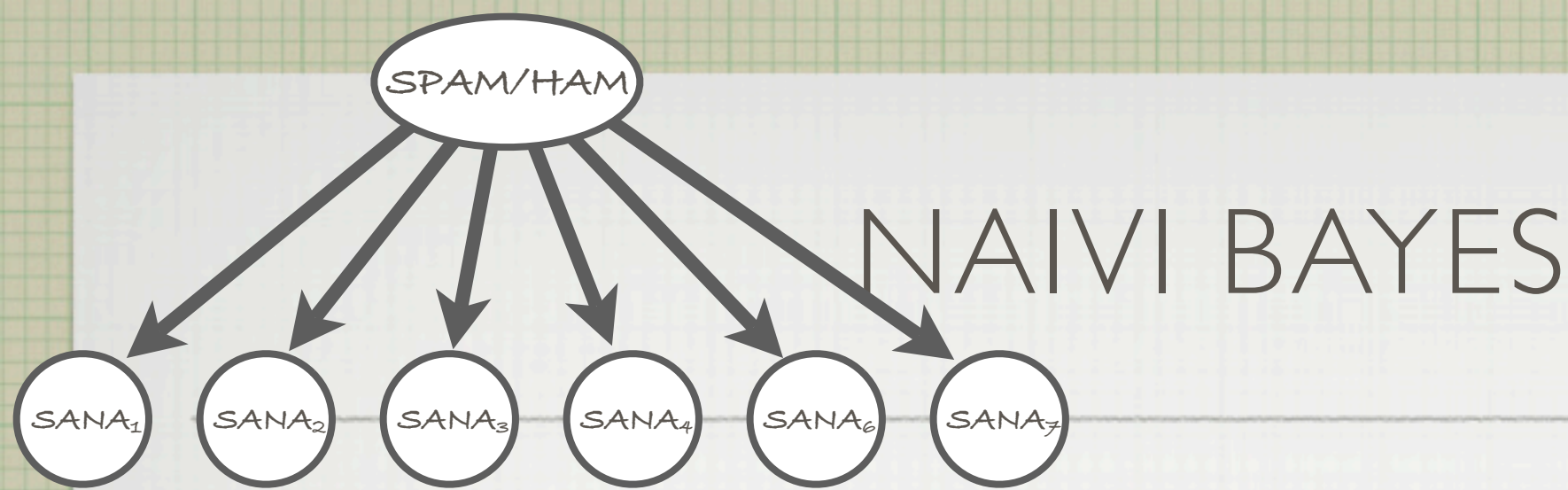
Odds ARVOSTA TULEE HELPOSTI LIIAN PIENI (LÄHELLÄ NOLLAA)

TAI LIIAN SUURI.

RATKAISU:

KÄYTÄ  $\log(\text{Odds})$





PSEUDOKOODINA:

$$\text{LOG}(A * B) = \text{LOG}(A) + \text{LOG}(B)$$

SPAMICITY(Viesti, P):

$$\text{logOdds} = \log(\text{P.Spam} / \text{P.noSpam}) \quad // \quad \text{P.Spam} + \text{P.noSpam} = 1$$

for each Sana in Viesti

$$\text{logOdds} = \text{logOdds} + \log(\text{P.Sana\_Spam}(\text{Sana}) / \text{P.Sana\_noSpam}(\text{Sana}))$$

return( $\exp(\text{logOdds})$ )

KÄYTÄNNÖN ONGELMA: ALI-JA YLIVUODOT

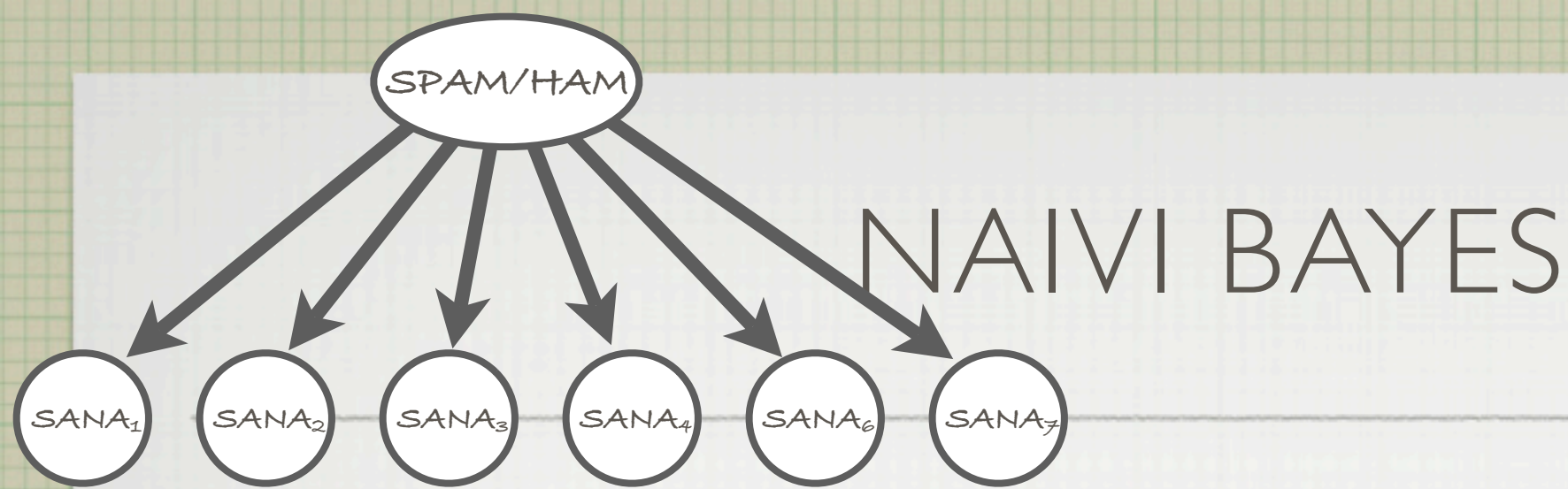
Odds ARVOSTA TULEE HELPOSTI LIIAN PIENI (LÄHELLÄ NOLLAA)

TAI LIIAN SUURI.

RATKAISU:

KÄYTÄ  $\log(\text{Odds})$





## PARAMETRIEN OPPIMISESTA

- \* VAIKEA KEKSIÄ PÄÄSTÄ EHDOLLISIA TN:IA.
- \* HUONOT ARVOT HUONONTAVAT FILTERIN TOIMINTAA
- \* PAREMPI RATKAISU:
  1. KERÄÄ ISO KASA SPAM-VIESTEJÄ
  2. KERÄÄ ISO KASA HAM-VIESTEJÄ
  3. ARVIOI  $P(SANA_i = \underline{\quad}) = 0.$  LASKEMALLA DATASTA  
(VRT. LASKUHARJ. 2.3)
- \* VAROTTAVA NOLLATODENNÄKÖISYYKSIÄ!  
(JOS KASA EI TARPEEKSI ISO, JOTKUT SANAT EIVÄT VAIN SATU ESIINTYMÄÄN SIINÄ.)



# ESIM.

The screenshot shows an email client window titled "Junk\_E-mail". The interface includes a menu bar with "Get Mail", "Write", "Address Book", and "Tag". A search bar at the top right says "Search all messages... <%K>". Below the menu bar, there's a "Quick Filter" section and another search bar "Filter these messages... <%F>".

The left sidebar shows a folder tree under "All Folders" for "cs-mail". The "Junk\_E-mail" folder is selected and highlighted. Other folders include "Inbox (1)", "Drafts", "Templates", "Sent", "Archives", "Trash", "3ci", "ai", "akademia", "anto+peik", "berkeley", "cam", "cosco", "cwi", "eira", "ekahau", "g", "hiit", "icms", "kolmog", "mat", "mdl", "-2006", "helsinki08", "mit", "opetus", "pascal+groups", "review", "AISTATS11", "ECAI08", "ECMLPKDD09", and "ITMB".

The main pane displays a list of emails with columns for "Subject", "From", and "Date". The selected email is:

Subject	From	Date
my new email	hillary	2/28/08 9:17 PM
Penis Enlargement Pills - Enlarge you Penis Naturally Gain Up To ...	Shana	9/24/11 7:33 PM
Replica watches - THE MOST POPULAR MODELS All our replica wa...	Charline Albertine	9/25/11 10:57 PM
Replica watches - THE MOST POPULAR MODELS All our replica wa...	Vanesa Karon	9/26/11 7:06 AM
Replica watches - THE MOST POPULAR MODELS All our replica wa...	SHAYNEKEITHA	4:55 PM
SPECIAL OFFER : VIAGRA on SALE at \$1.38 !!!	Margaretta Nita	9/26/11 9:52 PM
targeted email lists in many different areas	Mindy N Kerr	9/27/11 6:32 PM
The Microsoft Internet E-mail lottery Awards	Microsoft Corporation Sweepstakes Prom...	9/27/11 1:21 PM
Transform HR system to have higher value	Integrated HR Management	9/27/11 10:20 AM
Vahvistus AMSTERDAM BACTH NO: 15/3820/MGL	elizabeth.rice@virgilio.it	9/27/11 4:05 PM
Viagra 100mg x 60 Pills \$125, Free Pills & Reorder Discount, Top...	Creola Astrid	10:41 AM
Which Penis Enlargement Products Work?	Silvana Darcel	9/24/11 11:39 PM
Which Penis Enlargement Products Work?	NEDRA INDIA	11:18 AM

Below the list, the details of the selected email are shown:

from Creola Astrid <alvertaarlinda@onesource.com>☆  
subject Viagra 100mg x 60 Pills \$125, Free Pills & Reorder Discount, Top Selling 100% Quality & Satisfaction guaranteed!  
to sampo.sammalisto@cs.helsinki.fi☆

Buttons for "reply", "reply all", "forward", "archive", and "delete" are visible. The date and time "10:41 AM" and "other actions" are also shown.

The email content is displayed in a "Junk Mail" box with a "Not Junk" button. The text of the email is:

Best Buy Viagra Generic Online

Viagra 100mg x 60 Pills \$125, Free Pills & Reorder Discount, Top Selling 100% Quality & Satisfaction guaranteed!

We accept VISA, Master & E-Check Payments, 90000+ Satisfied Customers!  
<http://tabletpillsapr.ru>

At the bottom right, the status "Unread: 0 Total: 29" is displayed.



ESIM.

---



# ESIM.

## SPAM

1	MONEY	0.04	%
...			
5	VIAGRA	0.21	%
...			
10	IS	0.42	%
...			
19	REPLICA	0.80	%
20	EMAIL	0.84	%
20	YOU	0.84	%
21	DATABASE	0.88	%
25	EMAILS	1.05	%
26	OF	1.09	%
31	TO	1.30	%
43	AND	1.80	%
48	THE	2.01	%
TOTAL	2386		

## HAM

21	ALGORITHM	0.01	%
...			
62	MONEY	0.02	%
...			
2199	FOR	0.78	%
2492	THAT	0.88	%
2990	YOU	1.05	%
3141	IN	1.11	%
3160	I	1.11	%
3218	AND	1.13	%
3283	IS	1.16	%
3472	OF	1.22	%
3874	A	1.37	%
5442	TO	1.92	%
9196	THE	3.24	%
TOTAL	283736		



# ESIM.

## SPAM

## HAM

1 MONEY 0.04 %

...

5 VIAGRA 0.21 %

...

10 I

...

19 R

20 E

20 Y

21 DATABASE 0.88 %

25 EMAILS 1.05 %

26 OF 1.09 %

31 TO 1.30 %

43 AND 1.80 %

48 THE 2.01 %

TOTAL 2386

21 ALGORITHM 0.01 %

...

62 MONEY 0.02 %

...

3218 AND 1.13 %

3283 IS 1.16 %

3472 OF 1.22 %

3874 A 1.37 %

5442 TO 1.92 %

9196 THE 3.24 %

TOTAL 283736

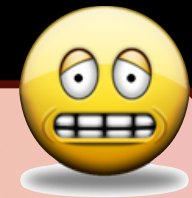
$$P(SANA_i = \text{MONEY} | \text{SPAM})$$

$$0.0004$$

$$\text{-----} = \text{-----} = 1.918 > 1$$

$$P(SANA_i = \text{MONEY} | \neg \text{SPAM})$$

$$0.0002$$





# ESIM.

## SPAM

## HAM

1 MONEY 0.04 %

...

5 VIAGRA 0.21 %

...

10 I

$$P(SANA_i = MONEY | SPAM) \quad 0.0004$$

...

19 R

$$\frac{P(SANA_i = MONEY | SPAM)}{P(SANA_i = MONEY | \neg SPAM)} = \frac{0.0004}{0.0002} = 1.918 > 1$$

20 E

20 Y

$$P(SANA_i = MONEY | \neg SPAM) \quad 0.0002$$

21 DATABASE 0.88 %

3218 AND 1.13 %

25 E

$$P(SANA_i = IS | SPAM) \quad 0.0042$$

26 O

31 T

$$\frac{P(SANA_i = IS | SPAM)}{P(SANA_i = IS | \neg SPAM)} = \frac{0.0042}{0.0116} = 0.3622 < 1$$

43 A

$$P(SANA_i = IS | \neg SPAM) \quad 0.0116$$

48 T

TOTAL 2386

TOTAL 283736





# YHTEENVETO

---

YHTEENVETO NAIVI BAYES-SPAMFILTTERISTÄ:

\* TARVITAAN:

- "PRIORIJAKAUMA"  $P(\text{SPAM}) = 0.\_\_\_$

- "LUOKKA-EHDOLLISET" JAKAUMAT

$$P(\text{SANA}_i = \text{VIAGRA} | \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{VIAGRA} | \neg \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{IS} | \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{IS} | \neg \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{ALGORITHM} | \text{SPAM}) = 0.\_\_\_$$

$$P(\text{SANA}_i = \text{ALG.} | \neg \text{SPAM}) = 0.\_\_\_$$

\* OLETETAAN ETTÄ  $P(\text{SANA}_i | \text{SANA}_j, \text{SPAM}) = P(\text{SANA}_i | \text{SPAM})$   
(EHDOLLINEN RIIPPUMATTOMUUS)

\* OLENNAINEN ON OSAMÄÄRÄ

$$P(\text{SANA}_i = \text{VIAGRA} | \text{SPAM})$$

---

$$P(\text{SANA}_i = \text{VIAGRA} | \neg \text{SPAM})$$



# YHTEENVETO

---

(JATKOA...):

- \* YLI- JA ALIVUOTOJEN VÄLTTÄMISEKSI KANNATTAÄ KÄYTTÄÄ LOGARITMIMUUNNOSTA ( $\text{LOG}(A*B) = \text{LOG}(A) + \text{LOG}(B)$ )
- \* JAKAUMAT PARAS ESTIMOIDA DATASTA
- \* NOLLATODENNÄKÖISYYKSILLE TEHTÄVÄJOTAIN