



CMSC 170: Introduction to Artificial intelligence

Week 05: Augmenting the Spam Filter with Laplace Smoothing

John O-Neil V. Geronimo
Institute of Computer Science
University of the Philippines Los Baños



Content

- I. Review on Naïve Bayes Spam Filter
- II. Augmenting the Spam Filter with Laplace Smoothing



Content

- I. Review on Naïve Bayes Spam Filter
- II. Augmenting the Spam Filter with Laplace Smoothing

Naïve Bayes Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Spam Filtering Using Naïve Bayes

$$P(\text{Spam}|\text{message}) = \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})}$$

Formula:

$$\begin{aligned} & P(\text{Spam}|\text{message}) \\ &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \end{aligned}$$

- $P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$
- $P(\text{Ham}) = 1 - P(\text{Spam})$
- $P(\text{message}|\text{Spam})$
 $= P(w_0|\text{Spam}) \dots P(w_n|\text{Spam})$
- $P(w_n|\text{Spam})$
 $= \frac{\text{count}(w \text{ in Spam})}{\text{count}(\text{total no of words in Spam})}$
- $P(\text{message}|\text{Ham})$
 $= P(w_0|\text{Ham}) \dots P(w_n|\text{Ham})$
- $P(w_n|\text{Ham})$
 $= \frac{\text{count}(w \text{ in Ham})}{\text{count}(\text{total no of words in Ham})}$
- $P(\text{message})$
 $= P(\text{message}|\text{Spam})P(\text{Spam})$
 $+ P(\text{message}|\text{Ham})P(\text{Ham})$



Putting it all together,

$$\begin{aligned} & P(S|m) \\ &= \frac{P(w_0|S)P(w_1|S) \dots P(w_n|S)P(S)}{P(w_0|S)P(w_1|S) \dots P(w_n|S)P(S) + P(w_0|H)P(w_1|H) \dots P(w_n|H)P(H)} \end{aligned}$$

where $m = \text{message}$, $S = \text{Spam}$, and $H = \text{Ham}$

Spam Filter



Spam Dataset



You will love this product!



I used this cream and it removed my acne!!!



Buy now!

Ham Dataset



Lucky I am in love with my bestfriend



You are still the one that I love



Could I love you anymore?

Spam Bag-of-Words

Index	Word	Frequency
0	you	1
1	will	1
2	love	1
3	this	2
4	product	1
5	i	1
6	used	1
7	cream	1

Index	Word	Frequency
8	and	1
9	it	1
10	removed	1
11	my	1
12	acne	1
13	buy	1
14	now	1

TNOW:

16

DS:

15

Ham Bag-of-Words

Index	Word	Frequency
0	lucky	1
1	i	3
2	am	1
3	in	1
4	love	3
5	with	1
6	my	1
7	bestfriend	1

Index	Word	Frequency
8	you	2
9	are	1
10	still	1
11	the	1
12	one	1
13	that	1
14	could	1
15	anymore	1

TNOW:

21

DS:

16

Spam Filter



Is this message ham or spam?

I love you.

Compute for $P(\text{Spam}|\text{message})$:

$$\begin{aligned}P(\text{Spam}|\text{message}) &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \\&= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message}|\text{Spam})P(\text{Spam}) + P(\text{message}|\text{Ham})P(\text{Ham})} \\&= \frac{\frac{1}{4096} (0.5)}{\frac{1}{4096} (0.5) + \frac{2}{1029} (0.5)} = 0.1115931027 = 11.16\%\end{aligned}$$

Spam Filter



Since $P(\text{Spam}|\text{message}) = 11.16\%$

The message “I love you.”
is classified as a **ham** message!

Spam Filter



Is this message ham or spam?

I hate you.

Spam Bag-of-Words

Index	Word	Frequency
0	you	1
1	will	1
2	love	1
3	this	2
4	product	1
5	i	1
6	used	1
7	cream	1

Index	Word	Frequency
8	and	1
9	it	1
10	removed	1
11	my	1
12	acne	1
13	buy	1
14	now	1

TNOW:

16

DS:

15

Spam Filter



Compute for $P(\text{message}|\text{Spam})$:

$$P(\text{"i"}|\text{Spam}) = \frac{\text{count}(\text{"i" in Spam})}{\text{count}(\text{total no of words in Spam})} = \frac{1}{16}$$

$$P(\text{"hate"}|\text{Spam}) = \frac{\text{count}(\text{"hate" in Spam})}{\text{count}(\text{total no of words in Spam})} = \frac{0}{16}$$

$$P(\text{"you"}|\text{Spam}) = \frac{\text{count}(\text{"you" in Spam})}{\text{count}(\text{total no of words in Spam})} = \frac{1}{16}$$

Compute for $P(\text{message}|\text{Spam})$:

$$\begin{aligned}P(\text{message}|\text{Spam}) &= P(w_0|\text{Spam})P(w_1|\text{Spam}) \dots P(w_n|\text{Spam}) \\&= P("i"|\text{Spam})P("love"|\text{Spam})P("you"|\text{Spam}) \\&= \frac{1}{16} * \frac{0}{16} * \frac{1}{16} \\&= \mathbf{0}\end{aligned}$$

Ham Bag-of-Words

Index	Word	Frequency
0	lucky	1
1	i	3
2	am	1
3	in	1
4	love	3
5	with	1
6	my	1
7	bestfriend	1

Index	Word	Frequency
8	you	2
9	are	1
10	still	1
11	the	1
12	one	1
13	that	1
14	could	1
15	anymore	1

TNOW:

21

DS:

16

Spam Filter



Compute for $P(\text{message}|\text{Spam})$:

$$P(\text{"i"}|\text{Ham}) = \frac{\text{count}(\text{"i"} \text{ in Ham})}{\text{count}(\text{total no of words in Ham})} = \frac{3}{21}$$

$$P(\text{"hate"}|\text{Ham}) = \frac{\text{count}(\text{"hate"} \text{ in Ham})}{\text{count}(\text{total no of words in Ham})} = \frac{0}{21}$$

$$P(\text{"you"}|\text{Ham}) = \frac{\text{count}(\text{"you"} \text{ in Ham})}{\text{count}(\text{total no of words in Ham})} = \frac{2}{21}$$

Compute for $P(\text{message}|\text{Ham})$:

$$\begin{aligned}P(\text{message}|\text{Ham}) &= P(w_0|\text{Ham})P(w_1|\text{Ham}) \dots P(w_n|\text{Ham}) \\&= P("i"|\text{Ham})P("hate"|\text{Ham})P("you"|\text{Ham}) \\&= \frac{3}{21} * \frac{0}{21} * \frac{2}{21} \\&= 0\end{aligned}$$

Spam Filter



Compute for $P(\text{Spam}|\text{message})$:

$$\begin{aligned}P(\text{Spam}|\text{message}) &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message})} \\&= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message}|\text{Spam})P(\text{Spam}) + P(\text{message}|\text{Ham})P(\text{Ham})} \\&= \frac{0 (0.5)}{0(0.5) + 0 (0.5)} = \frac{0}{0} = \boxed{\text{undefined}} \rightarrow \text{overfitting}\end{aligned}$$



Content

- I. Review on Naïve Bayes Spam Filter
- II. Augmenting the Spam Filter with Laplace Smoothing

Updated Formula



Before

$$P(\text{Spam}) = \frac{\text{count}(\text{Spam})}{\text{count}(\text{Spam} \cup \text{Ham})}$$

$$P(\text{Ham}) = \frac{\text{count}(\text{Ham})}{\text{count}(\text{Spam} \cup \text{Ham})}$$

$$= 1 - P(\text{Spam})$$

Now

$$P(\text{Spam}) = \frac{\text{count}(\text{Spam}) + k}{\text{count}(\text{Spam} \cup \text{Ham}) + 2k}$$

$$P(\text{Ham}) = \frac{\text{count}(\text{Ham}) + k}{\text{count}(\text{Spam} \cup \text{Ham}) + 2k}$$

$$= 1 - P(\text{Spam})$$

Updated Formula



Before

$$P(w_n|Spam) = \frac{\text{count}(w \text{ in } Spam)}{\text{count}(\text{total no of words in } Spam)}$$

Now

$$P(w_n|Spam) = \frac{\text{count}(w \text{ in } Spam) + k}{\text{count}(\text{total no of words in } Spam) + (k * (dSize + \text{count}(\text{new words})))}$$

Updated Formula



$$P(w_n|Spam) = \frac{count(w \text{ in Spam}) + k}{count(total \text{ no of words in Spam}) + (k * (dSize + count(new \text{ words})))}$$

dSize (or *dictionary size*) is the number of unique words found on both Spam and Ham dataset

new words refers to the words present in the message to be classified but does not exist in the Ham and Spam dictionary

Spam Filter



Spam	Ham
Secret link offers	Link for ticket offers
Click link win ticket offers	Click link for win
	Reply link for ticket
Secret link offers for secret win	Win ticket for secret place

Is this message ham or spam?

Secret place haven

Consider **$k = 2$**

Spam Bag-of-Words

Index	Word	Frequency
0	secret	3
1	link	3
2	offers	3
3	click	1
4	win	2
5	ticket	1
6	for	1

TNOW:

14

DS:

7

Ham Bag-of-Words

Index	Word	Frequency
0	link	3
1	for	4
2	ticket	3
3	offers	1
4	click	1

Index	Word	Frequency
5	win	2
6	reply	1
7	secret	1
8	place	1

TNOW:

17

DS:

9

Spam Filter



Compute for $P(\text{Spam})$:

Since we have three emails in the spam dataset and a total of seven messages for both spam and ham dataset

$$P(\text{Spam}) = \frac{\text{count}(\text{Spam}) + k}{\text{count}(\text{Spam} \cup \text{Ham}) + 2k}$$

$$P(\text{Spam}) = \frac{3 + 2}{7 + 2(2)} = \frac{5}{11}$$

Spam Filter



Compute for $P(\text{Ham})$:

Since we have four emails in the spam dataset and a total of seven messages for both spam and ham dataset

$$P(\text{Ham}) = \frac{\text{count}(\text{Ham}) + k}{\text{count}(\text{Spam} \cup \text{Ham}) + 2k}$$

$$P(\text{Ham}) = \frac{4 + 2}{7 + 2(2)} = \frac{6}{11}$$

Compute for $P(\text{message}|\text{Spam})$:

$$P(\text{"secret"}|\text{Spam}) = \frac{\text{count}(\text{"secret"} \text{ in Spam}) + k}{\text{count}(\text{total no of words in Spam}) + (k * (dS + nW))}$$

$$P(\text{"place"}|\text{Spam}) = \frac{\text{count}(\text{"place"} \text{ in Spam}) + k}{\text{count}(\text{total no of words in Spam}) + (k * (dS + nW))}$$

$$P(\text{"haven"}|\text{Spam}) = \frac{\text{count}(\text{"haven"} \text{ in Spam}) + k}{\text{count}(\text{total no of words in Spam}) + (k * (dS + nW))}$$

Updated Formula



$$P(w_n|Spam) = \frac{\text{count}(w \text{ in Spam}) + k}{\text{count}(\text{total no of words in Spam}) + (k * (dSize + \text{count}(\text{new words})))}$$

dSize (or *dictionary size*) is the number of unique words found on both Spam and Ham dataset

new words refers to the words present in the message to be classified but does not exist in the Ham and Spam dictionary

Spam Filter



Spam and Ham dictionary

Unique Words	
secret	for
link	win
offers	reply
click	place
ticket	

Dictionary Size: 9

Message to classify:
secret place haven

Count(new words): 1

Spam Bag-of-Words

Index	Word	Frequency
0	secret	3
1	link	3
2	offers	3
3	click	1
4	win	2
5	ticket	1
6	for	1

TNOW:

14

DS:

7

Compute for $P(\text{message}|\text{Spam})$:

$$P(\text{secret}|\text{Spam}) = \frac{\text{count}(\text{"secret"}\text{in Spam}) + k}{\text{count}(\text{total no of words in Spam}) + (k * (dS + nW))} = \frac{3 + 2}{14 + (2 * (9 + 1))} = \frac{5}{34}$$

$$P(\text{place}|\text{Spam}) = \frac{\text{count}(\text{"place"}\text{in Spam}) + k}{\text{count}(\text{total no of words in Spam}) + (k * (dS + nW))} = \frac{0 + 2}{14 + (2 * (9 + 1))} = \frac{2}{34}$$

$$P(\text{haven}|\text{Spam}) = \frac{\text{count}(\text{"haven"}\text{in Spam}) + k}{\text{count}(\text{total no of words in Spam}) + (k * (dS + nW))} = \frac{0 + 2}{14 + (2 * (9 + 1))} = \frac{2}{34}$$

Ham Bag-of-Words

Index	Word	Frequency
0	link	3
1	for	4
2	ticket	3
3	offers	1
4	click	1

Index	Word	Frequency
5	win	2
6	reply	1
7	secret	1
8	place	1

TNOW:

17

DS:

9

Compute for $P(\text{message}|\text{Ham})$:

$$P(\text{secret}|\text{Ham}) = \frac{\text{count}(\text{"secret" in Ham}) + k}{\text{count}(\text{total no of words in Ham}) + (k * (dS + nW))} = \frac{1 + 2}{17 + (2 * (9 + 1))} = \frac{3}{37}$$

$$P(\text{place}|\text{Ham}) = \frac{\text{count}(\text{"place" in Ham}) + k}{\text{count}(\text{total no of words in Ham}) + (k * (dS + nW))} = \frac{1 + 2}{17 + (2 * (9 + 1))} = \frac{3}{37}$$

$$P(\text{haven}|\text{Ham}) = \frac{\text{count}(\text{"haven" in Ham}) + k}{\text{count}(\text{total no of words in Ham}) + (k * (dS + nW))} = \frac{0 + 2}{17 + (2 * (9 + 1))} = \frac{2}{37}$$

Spam Filter



Compute for $P(\text{message}|\text{Ham})$ and $P(\text{message}|\text{Spam})$:

$$P(\text{message}|\text{Spam}) \\ = P(w_0|\text{Spam})P(w_1|\text{Spam}) \dots P(w_n|\text{Spam})$$

$$= \frac{5}{34} * \frac{2}{34} * \frac{2}{34}$$

$$= \frac{5}{9826}$$

$$P(\text{message}|\text{Ham}) \\ = P(w_0|\text{Ham})P(w_1|\text{Ham}) \dots P(w_n|\text{Ham})$$

$$= \frac{3}{37} * \frac{3}{37} * \frac{2}{37}$$

$$= \frac{12}{50653}$$



Too small decimal values?

Use log probabilities for
Naïve Bayes



Too small decimal values?

Use log probabilities for
Naïve Bayes

Compute for Total Spam = $P(\text{message}|\text{Spam})P(\text{Spam})$:

$$P(\text{message}|\text{Spam}) = P(w_0|\text{Spam})P(w_1|\text{Spam}) \dots P(w_n|\text{Spam}) P(\text{Spam})$$

$$= \ln \left(\frac{5}{34} * \frac{2}{34} * \frac{2}{34} * \frac{5}{11} \right)$$

$$= -8.371806661$$

Compute for Total Ham = $P(\text{message}|\text{Ham})P(\text{Ham})$:

$$P(\text{message}|\text{Ham}) = P(w_0|\text{Ham})P(w_1|\text{Ham}) \dots P(w_n|\text{Ham}) P(\text{Ham})$$

$$= \ln \left(\frac{3}{37} * \frac{3}{37} * \frac{2}{37} * \frac{6}{11} \right)$$

$$= -8.548517784$$

Compute for $P(\text{Spam}|\text{message})$:

$$\begin{aligned}P(\text{Spam}|\text{message}) &= \frac{P(\text{message}|\text{Spam})P(\text{Spam})}{P(\text{message}|\text{Spam})P(\text{Spam}) + P(\text{message}|\text{Ham})P(\text{Ham})} \\&= \frac{\exp(\text{total spam})}{\exp(\text{total spam}) + \exp(\text{total ham})} \\&= \frac{\exp(-8.371806661)}{\exp(-8.371806661) + \exp(-8.548517784)} \\&= 0.5440631776 = \mathbf{54.41\%}\end{aligned}$$

Spam Filter



Since $P(\text{Spam}|\text{message}) = 54.41\%$

The message “secret place haven”
is classified as a **spam** message!



Content

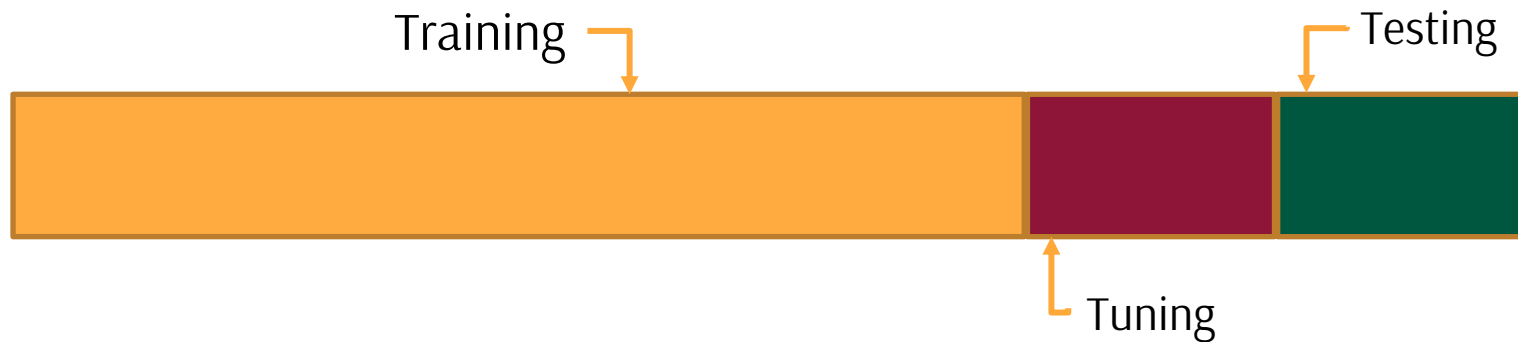
- I. Review on Naïve Bayes Spam Filter
- II. Augmenting the Spam Filter with Laplace Smoothing



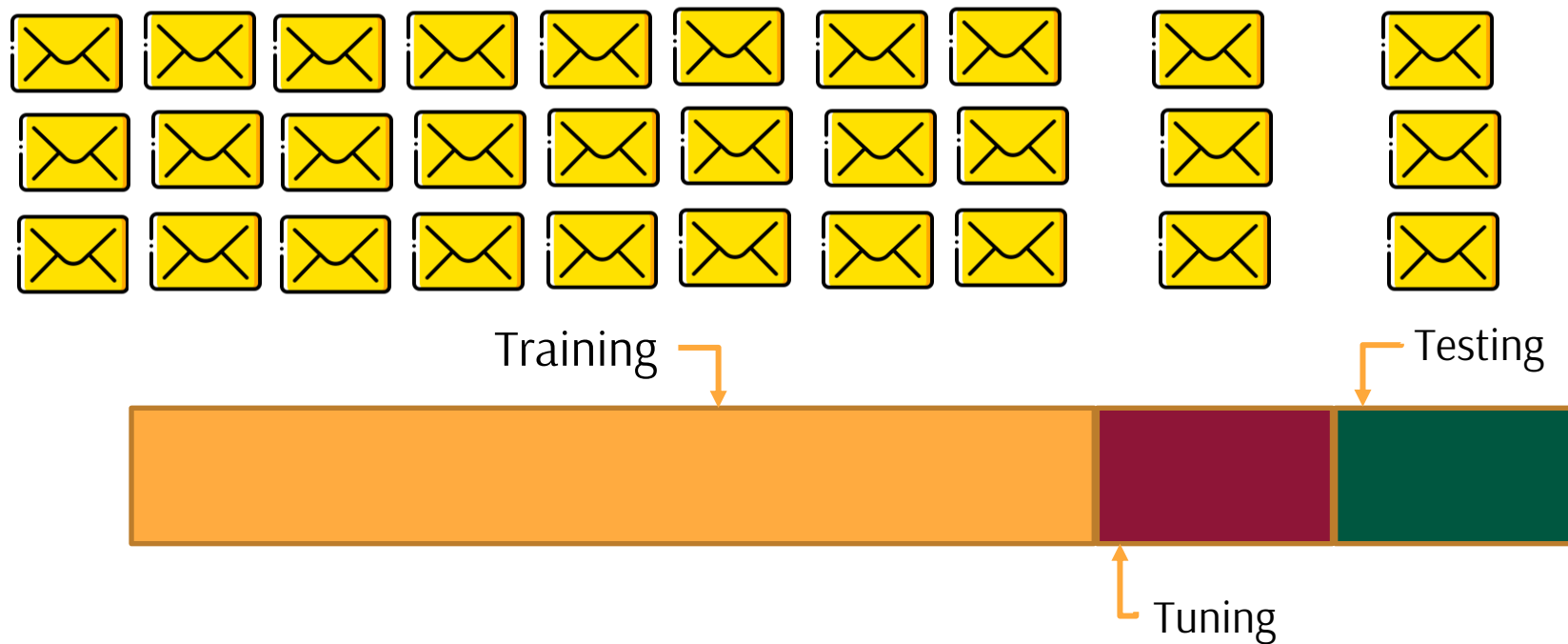
Exercise

Cross validation

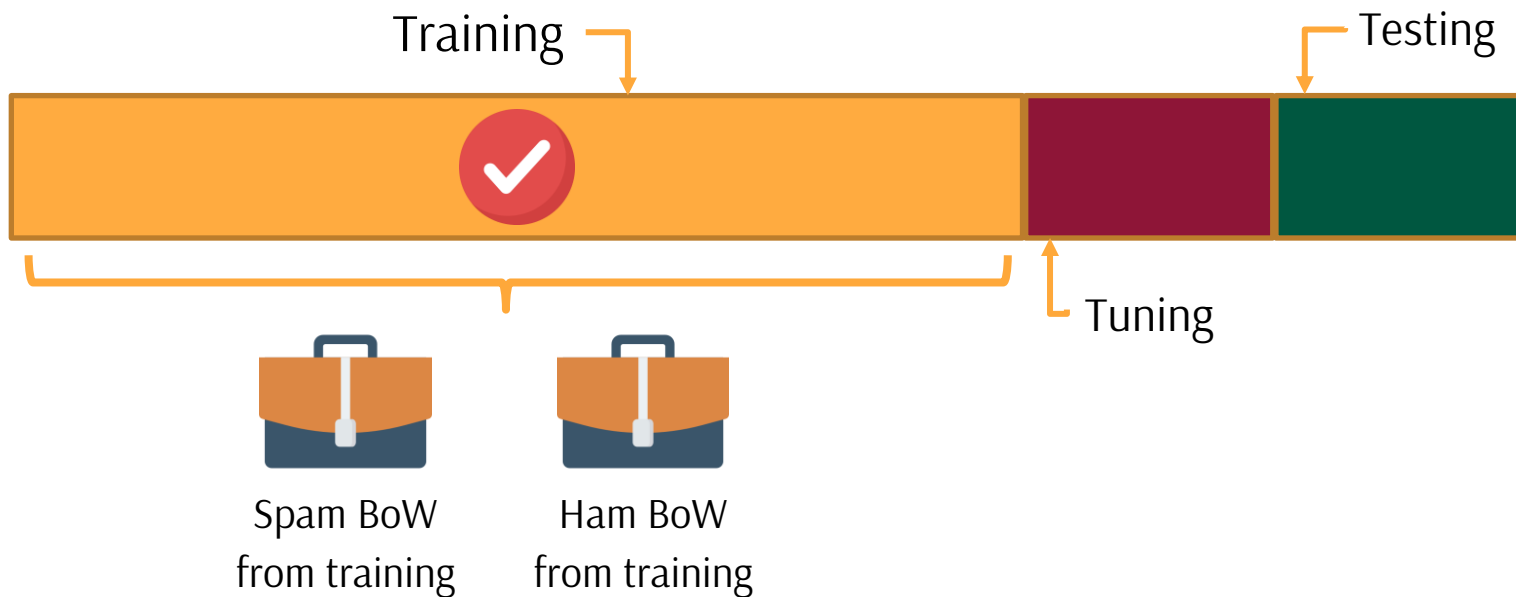
Spam Filter



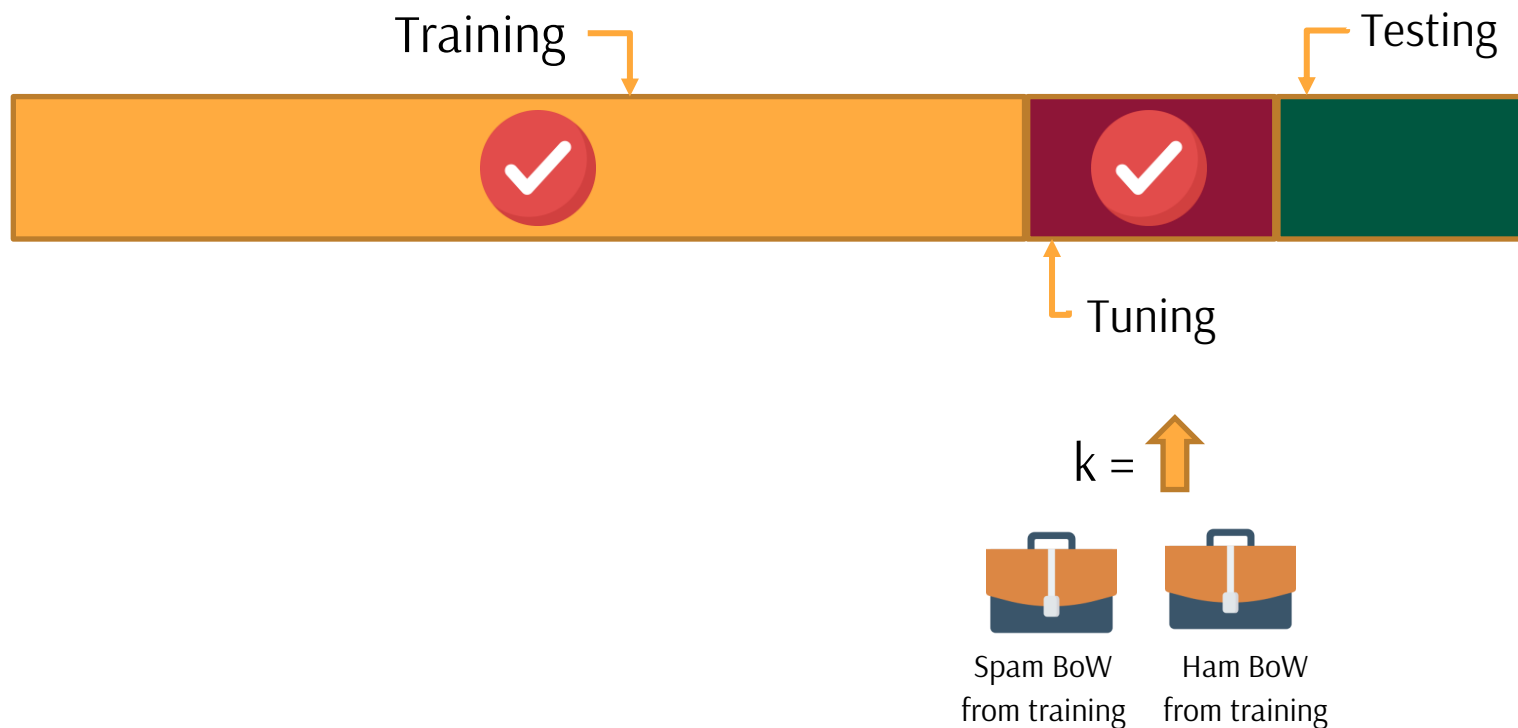
Spam Filter



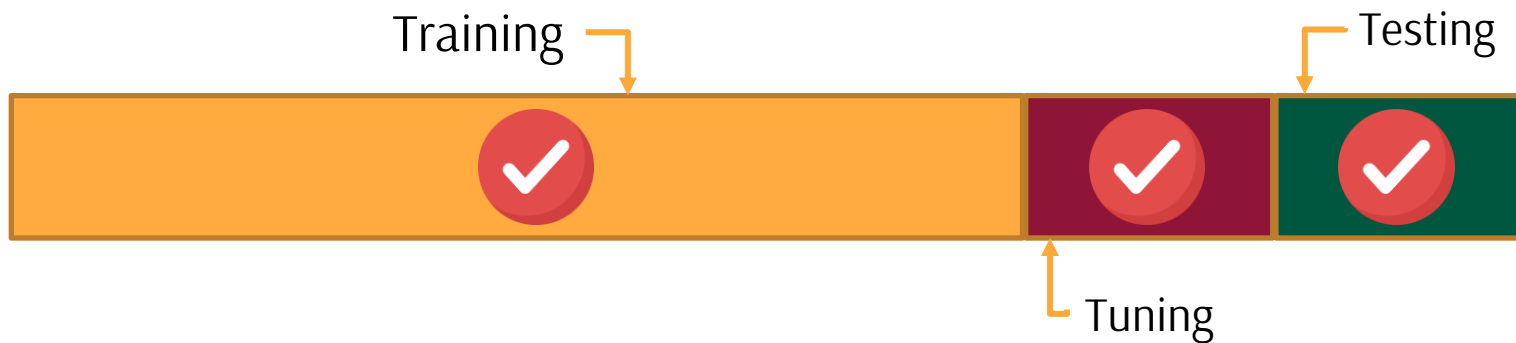
Spam Filter



Spam Filter



Spam Filter



best **k value** from tuning



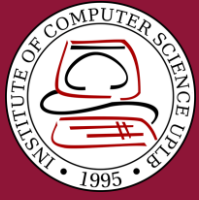
Spam BoW
from training



Ham BoW
from training



Confusion matrix



Keep safe!