

# CS 585

## *Natural Language Processing*

February 12, 2024

# Announcements / Reminders

- Please follow the Week 05 To Do List instructions (if you haven't already)
- Written Assignment #03: posted
- Programming Assignment #02 will be posted soon
- Midterm Exam: 02/26/2024
  - Section 02 – Make arrangements with Mr. Charles Scott

# Plan for Today

- Text Classification

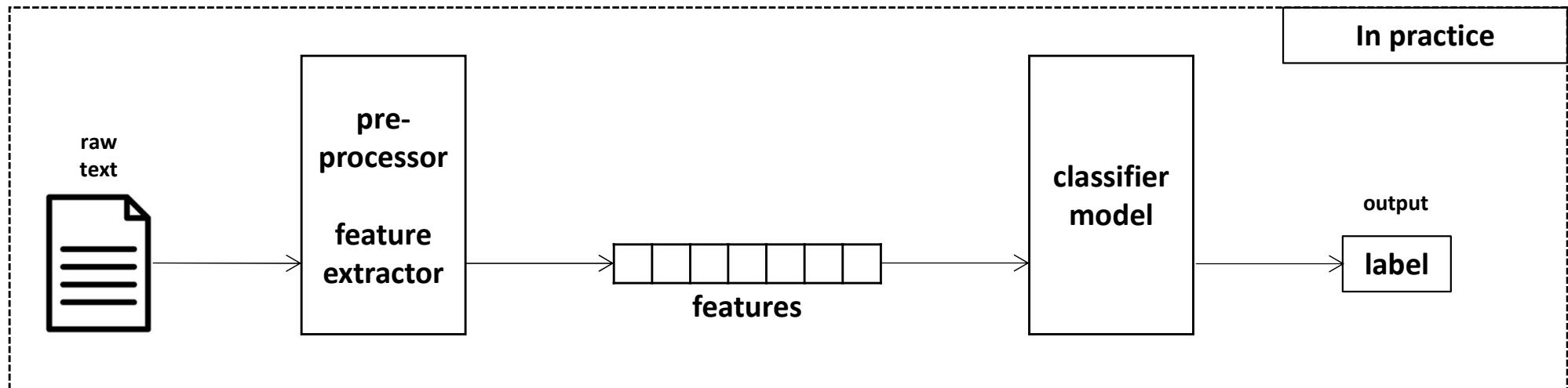
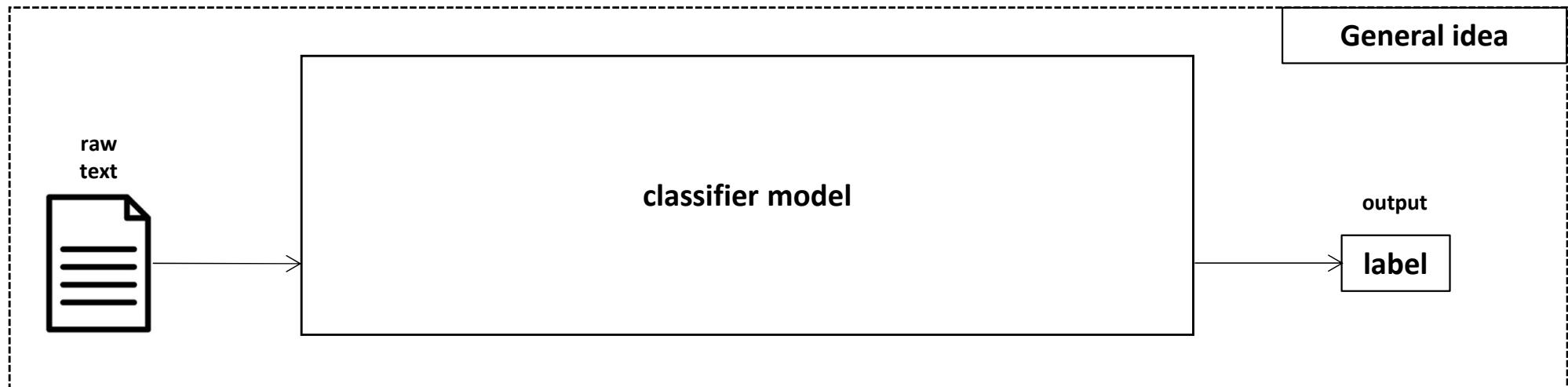
# Text Classification: Definition

*Input:*

- a document  $d$
- a fixed set of classes  $C = \{c_1, c_2, \dots, c_J\}$

*Output:* a predicted class  $c \in C$

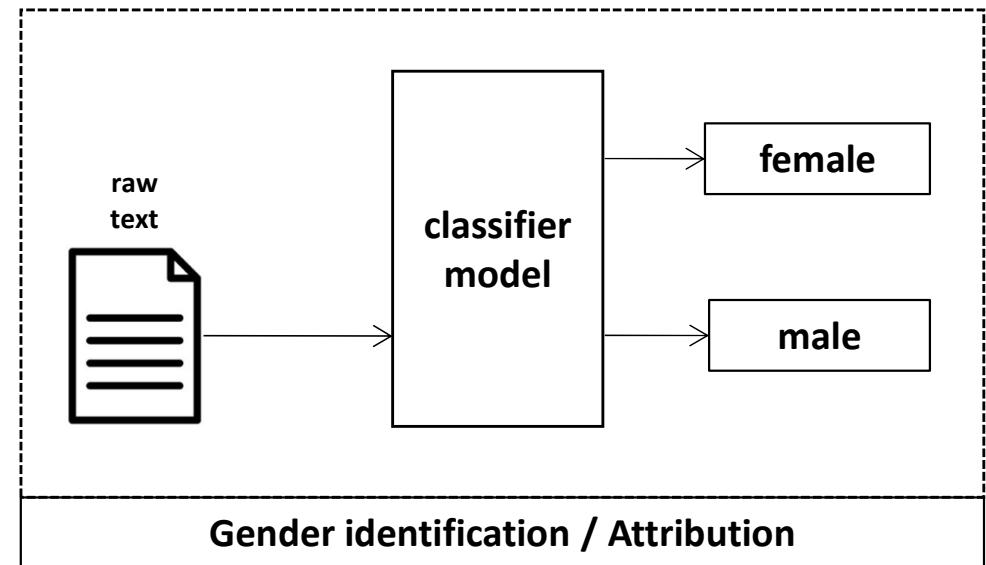
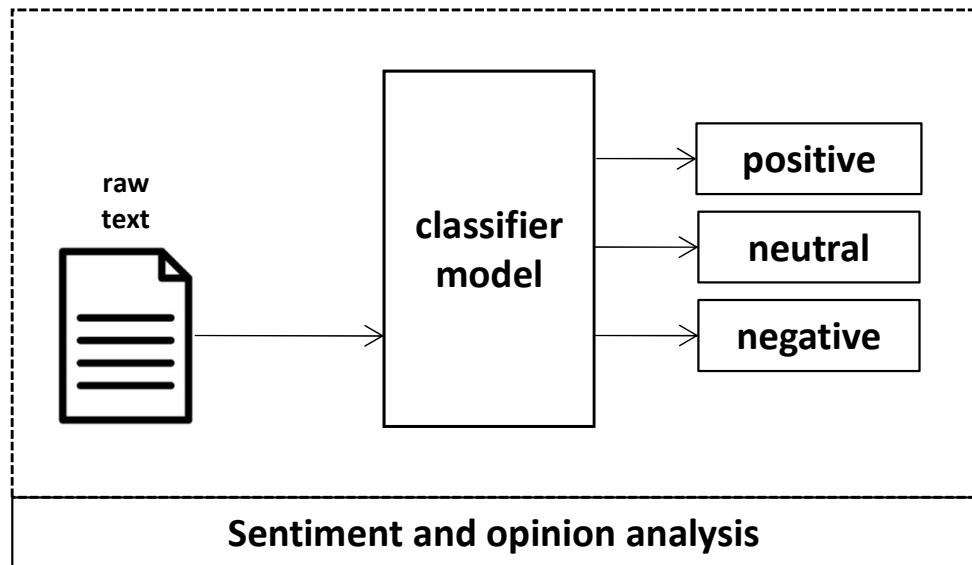
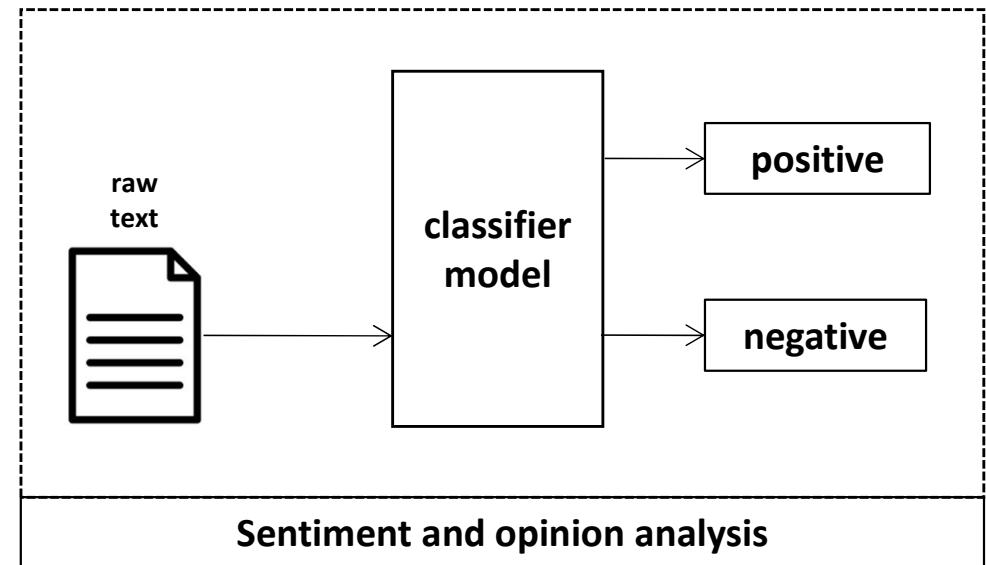
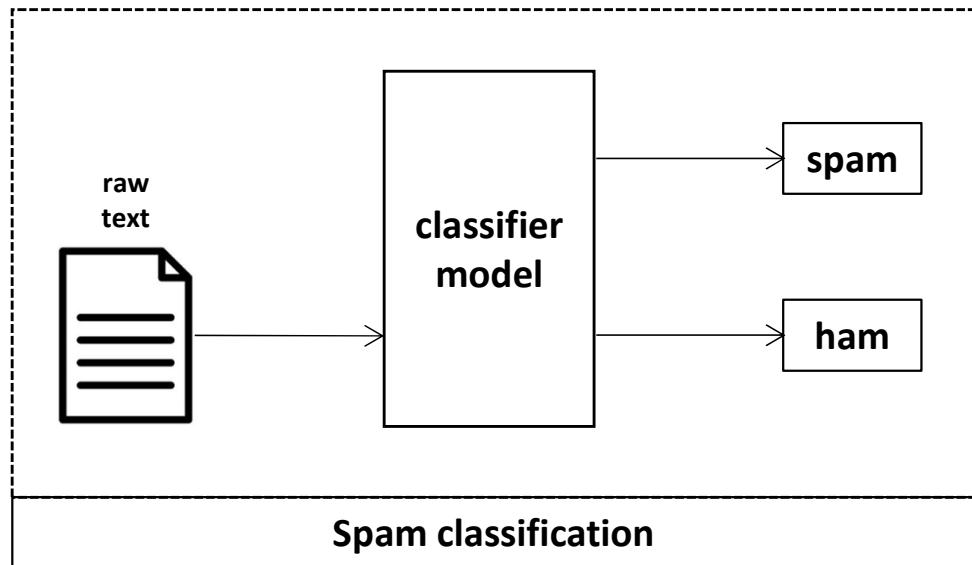
# Text Classification: the Idea



# **Text Classification: Applications**

- **Sentiment / opinion analysis**
- **Spam detection**
- **Gender identification**
- **Authorship identification**
- **Language identification**
- **Assigning subject categories, topics, or genres**
- ...

# Text Classification: Applications



# Text Classification: Rule-Based

- Rules based on combinations of words or other features
  - spam: black-list-address OR (“dollars” AND “you have been selected”)
- Accuracy can be high
  - If rules carefully refined by expert
- But building and maintaining these rules is expensive

# Text Classification: Supervised ML

*Input:*

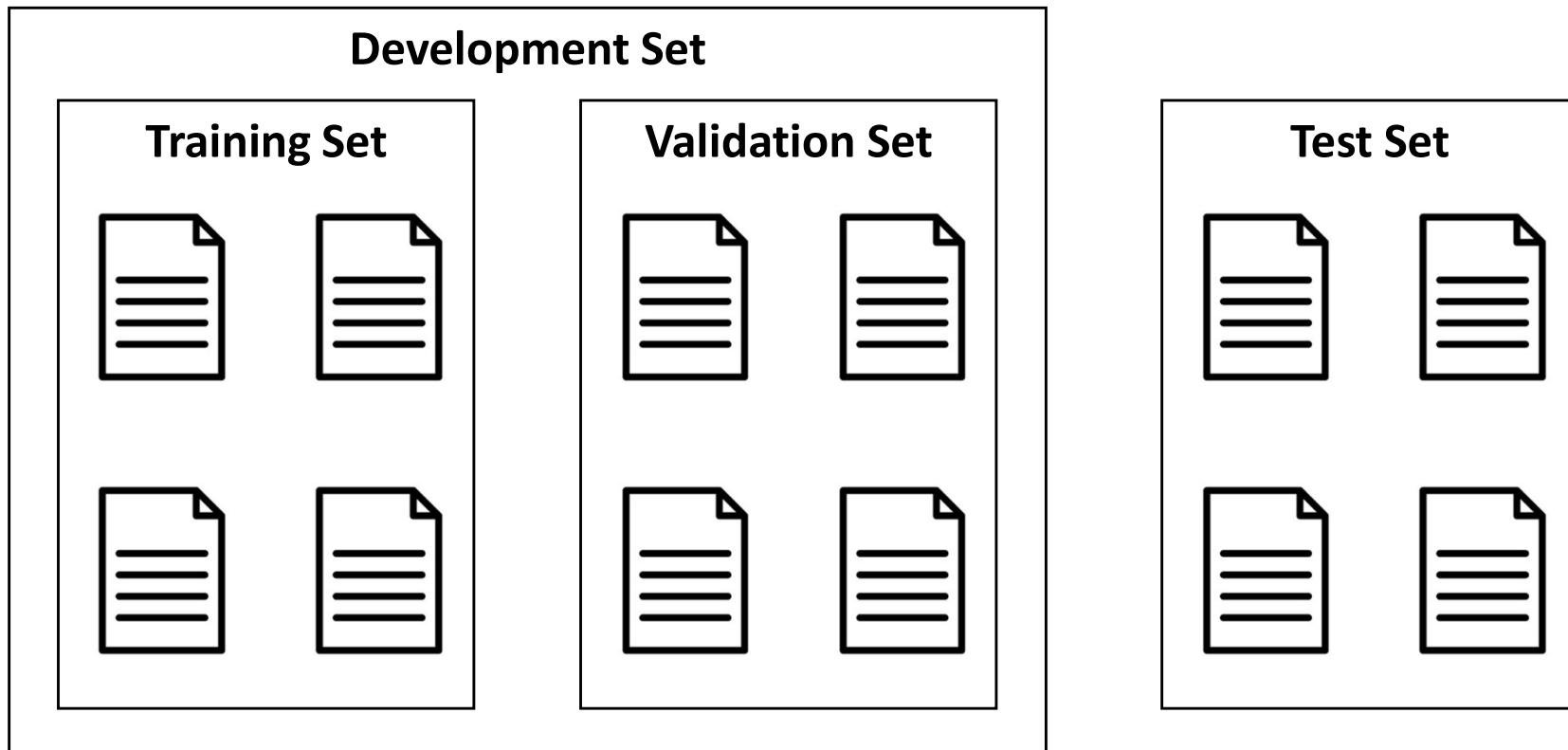
- a document  $d$
- a fixed set of classes  $C = \{c_1, c_2, \dots, c_J\}$
- a training set of  $m$  hand-labeled documents  $(d_1, c_1), \dots, (d_m, c_m)$

*Output:*

- a learned classifier  $\gamma: d \rightarrow c$

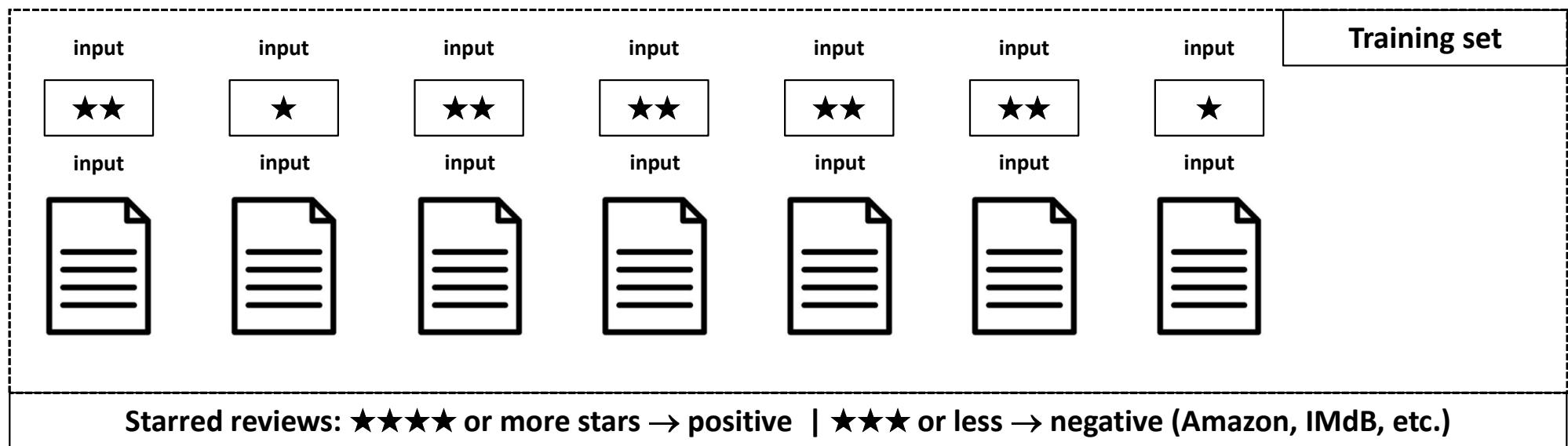
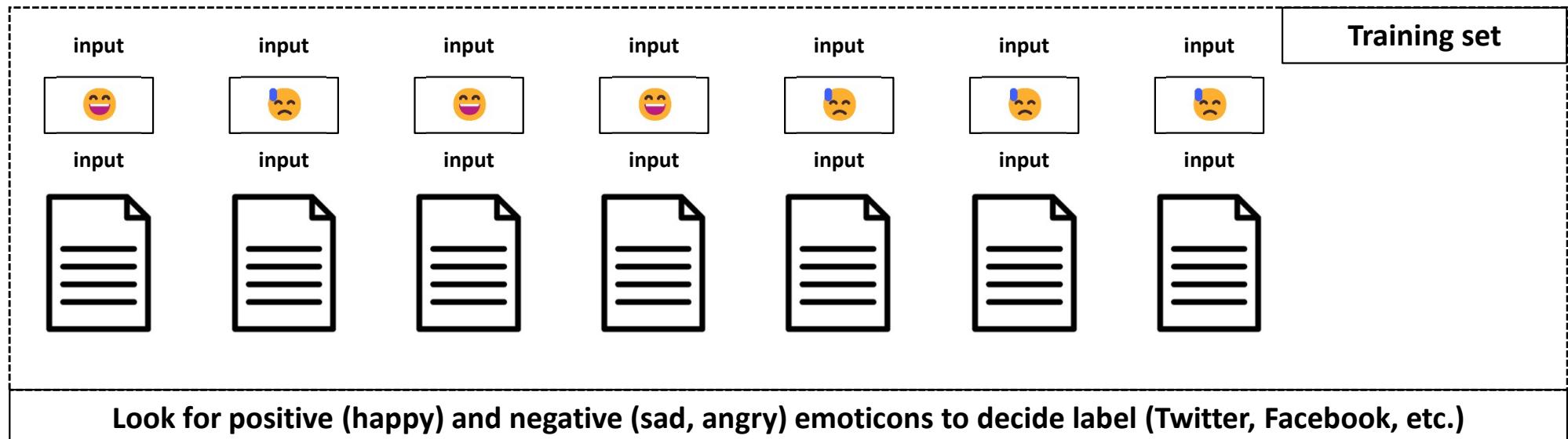
# Corpus: Training / Validation / Test

Corpus



Typical training / validation / test set split for a text corpora

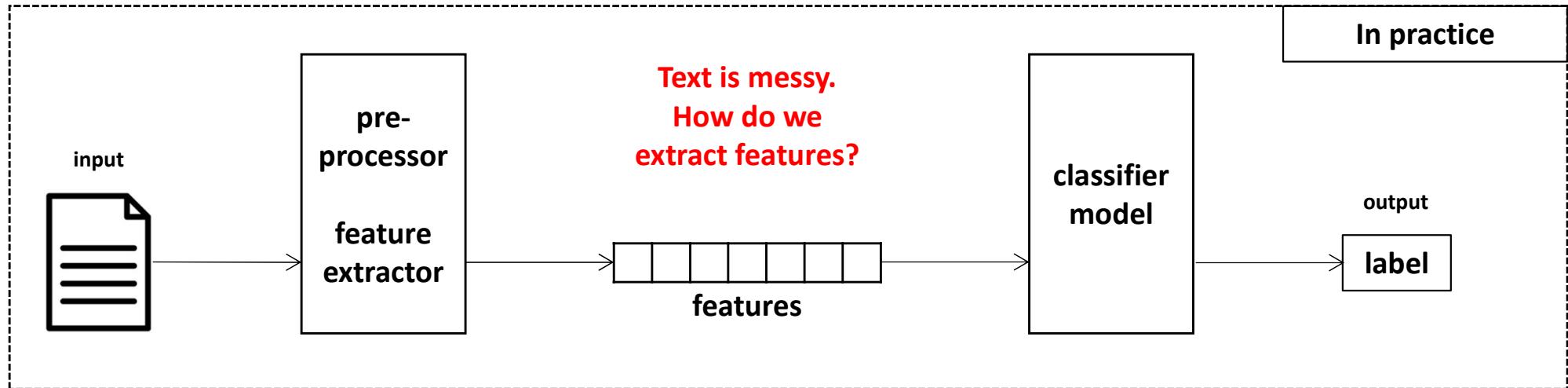
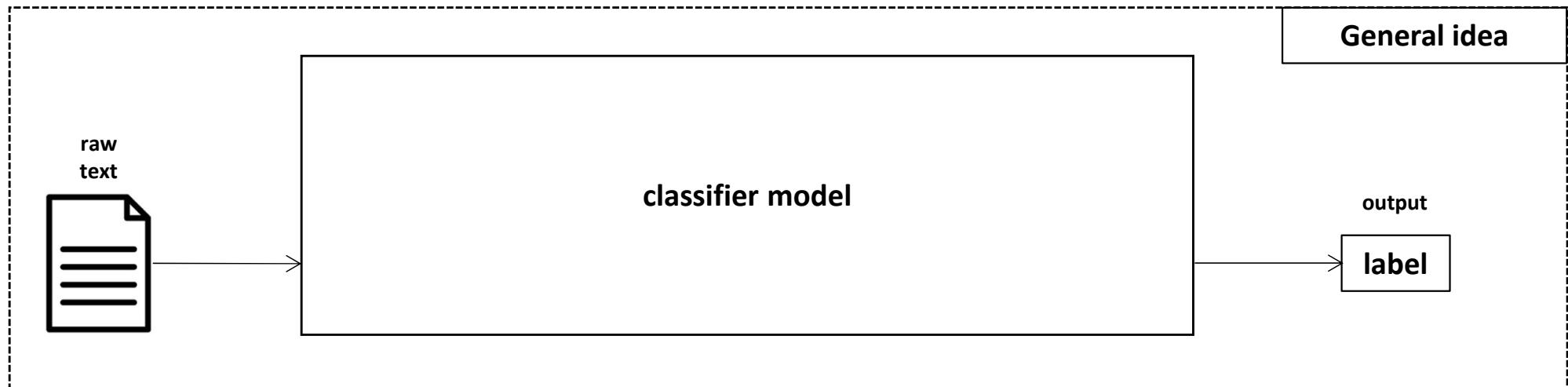
# Text Training Set (Auto) Labeling



# **Text Classification: Supervised ML**

- Various Machine Learning supervised learning classifier approaches can be employed:
  - Naïve Bayes
  - Logistic regression
  - Neural networks
  - k-Nearest Neighbors
  - etc.

# Text Classification: Feature Extraction



# Bag of Words: the Idea

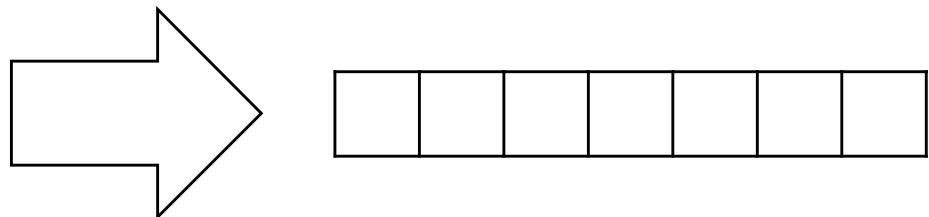
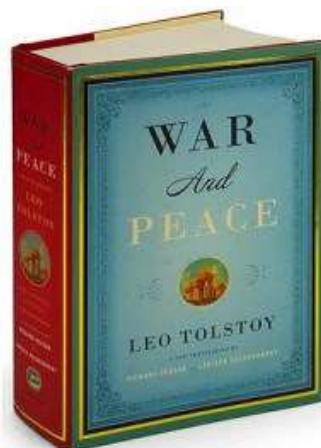


By Amy Bizzari 1st March 2022

Obtained from the autumnal flowering of the strawberry tree on the island of Sardinia, corbezzolo honey isn't sweet and has a history that dates back more than 2,000 years.

**C**orbezzolo honey tricks the palate. Instead of the sweetness one would expect, this extremely rare honey, born in the mountains of the Italian island of Sardinia, is surprisingly bitter, with notes of leather, liquorice and smoke. Nomadic beekeepers have been setting up beehives in the region to collect this aromatic treat – derived from the white, bell-shaped flowers of the wild strawberry tree – for more than 2,000 years.

Statesman, lawyer and philosopher Marcus Tullius Cicero (106-43 BCE) mentioned the honey in his defence of a Roman citizen accused of murder in Nora, Sardinia. "Omne quod Sardinia fert, homines et res, mala est! Etiam mel quod ea insula abundat, amarum est! (Everything that the island of Sardinia produces, men and things, is bad!)," he exclaimed. "Even the honey, abundant on that island, is bitter!"



**FIXED size**

**Feature vector**

# Bag of Words: the Idea

## Some document:

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

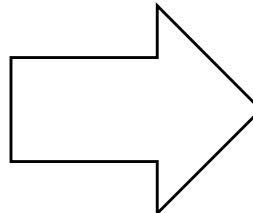


**Bag of words assumption:** word/token position does not matter.

# Bag of Words: the Idea

## Some document:

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



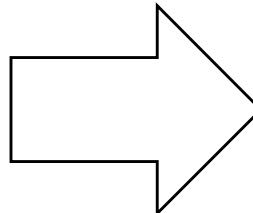
Word:	Frequency:
it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
whimsical	1
times	1
....	...

**Bag of words assumption:** word/token position does not matter.

# Bag of Words: Document Vector

## Some document:

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



Word:	Frequency:
it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
whimsical	1
times	1
....	...



vector

# Bag of Words: Document Vector

Pre-defined Vocabulary:

Word 1	Word 2	Word 3	Word 4	Word 5	Word 6	...	Word N
--------	--------	--------	--------	--------	--------	-----	--------

Document A **Binary** Vector [0-word absent | 1-word present]:

1	0	1	1	1	0	...	1
---	---	---	---	---	---	-----	---

Document B **Binary** Vector [0-word absent | 1-word present]:

1	1	0	0	1	0	...	1
---	---	---	---	---	---	-----	---

Document C **Binary** Vector [0-word absent | 1-word present]:

0	0	1	0	0	1	...	0
---	---	---	---	---	---	-----	---

Document vectors can be used to **compare documents**.

# Bag of Words: Document Vector

Pre-defined Vocabulary:

Word 1	Word 2	Word 3	Word 4	Word 5	Word 6	...	Word N
--------	--------	--------	--------	--------	--------	-----	--------

Document A Non-binary Vector [0-word absent | >0-word count]:

6	0	2	3	1	0	...	4
---	---	---	---	---	---	-----	---

Document B Non-binary Vector [0-word absent | >0-word count]:

4	2	0	0	5	0	...	1
---	---	---	---	---	---	-----	---

Document C Non-binary Vector [0-word absent | >0-word count]:

0	0	3	0	0	7	...	0
---	---	---	---	---	---	-----	---

Document vectors can be used to compare documents.

# Bag of Words: Document Vector

Pre-defined Vocabulary:

Word 1	Word 2	Word 3	Word 4	Word 5	Word 6	...	Word N
--------	--------	--------	--------	--------	--------	-----	--------

Document A **Binary** Vector [0-word absent | 1-word present]:

1	0	1	1	1	0	...	1
---	---	---	---	---	---	-----	---

Document B **Binary** Vector [0-word absent | 1-word present]:

1	1	0	0	1	0	...	1
---	---	---	---	---	---	-----	---

Document C **Binary** Vector [0-word absent | 1-word present]:

0	0	1	0	0	1	...	0
---	---	---	---	---	---	-----	---

Document vectors can be used to **compare documents**.

# Document Vector = Feature Vector

Pre-defined **Features**:

Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	...	Feature N
-----------	-----------	-----------	-----------	-----------	-----------	-----	-----------

Document A **Binary** Vector [0-word absent | 1-word present]:

1	0	1	1	1	0	...	1
---	---	---	---	---	---	-----	---

Document B **Binary** Vector [0-word absent | 1-word present]:

1	1	0	0	1	0	...	1
---	---	---	---	---	---	-----	---

Document C **Binary** Vector [0-word absent | 1-word present]:

0	0	1	0	0	1	...	0
---	---	---	---	---	---	-----	---

Document vectors can be used to **compare documents**.

# Bag of Words: Document Vector

Pre-defined Vocabulary:

she	want	to	walk	drive	fly	there	or
-----	------	----	------	-------	-----	-------	----

“She wants to walk there today”: Binary Document Vector

1	1	1	1	0	0	1	0
---	---	---	---	---	---	---	---

“She wants to drive there today”: Binary Document Vector

1	1	1	0	1	0	1	0
---	---	---	---	---	---	---	---

“She wants to fly or drive there today”: Binary Document Vector

1	1	1	0	1	1	1	1
---	---	---	---	---	---	---	---

Note: sentences lemmatized and lowercased.

# Bag of Bigrams: Document Vector

Pre-defined Bigrams:

w1, w2	w2, w3	w3, w4	w4, w5	w5, w6	w6,w7	...	wN-1,wN
--------	--------	--------	--------	--------	-------	-----	---------

Document A **Binary** Vector [0-word absent | 1-word present]:

1	0	1	1	1	0	...	1
---	---	---	---	---	---	-----	---

Document B **Binary** Vector [0-word absent | 1-word present]:

1	1	0	0	1	0	...	1
---	---	---	---	---	---	-----	---

Document C **Binary** Vector [0-word absent | 1-word present]:

0	0	1	0	0	1	...	0
---	---	---	---	---	---	-----	---

Document vectors can be used to **compare documents**.

# Bag of Words: Classification

category =  $h($

Learned Classifier model  
(hypothesis)

6
5
4
3
3
2
1
1
1
...

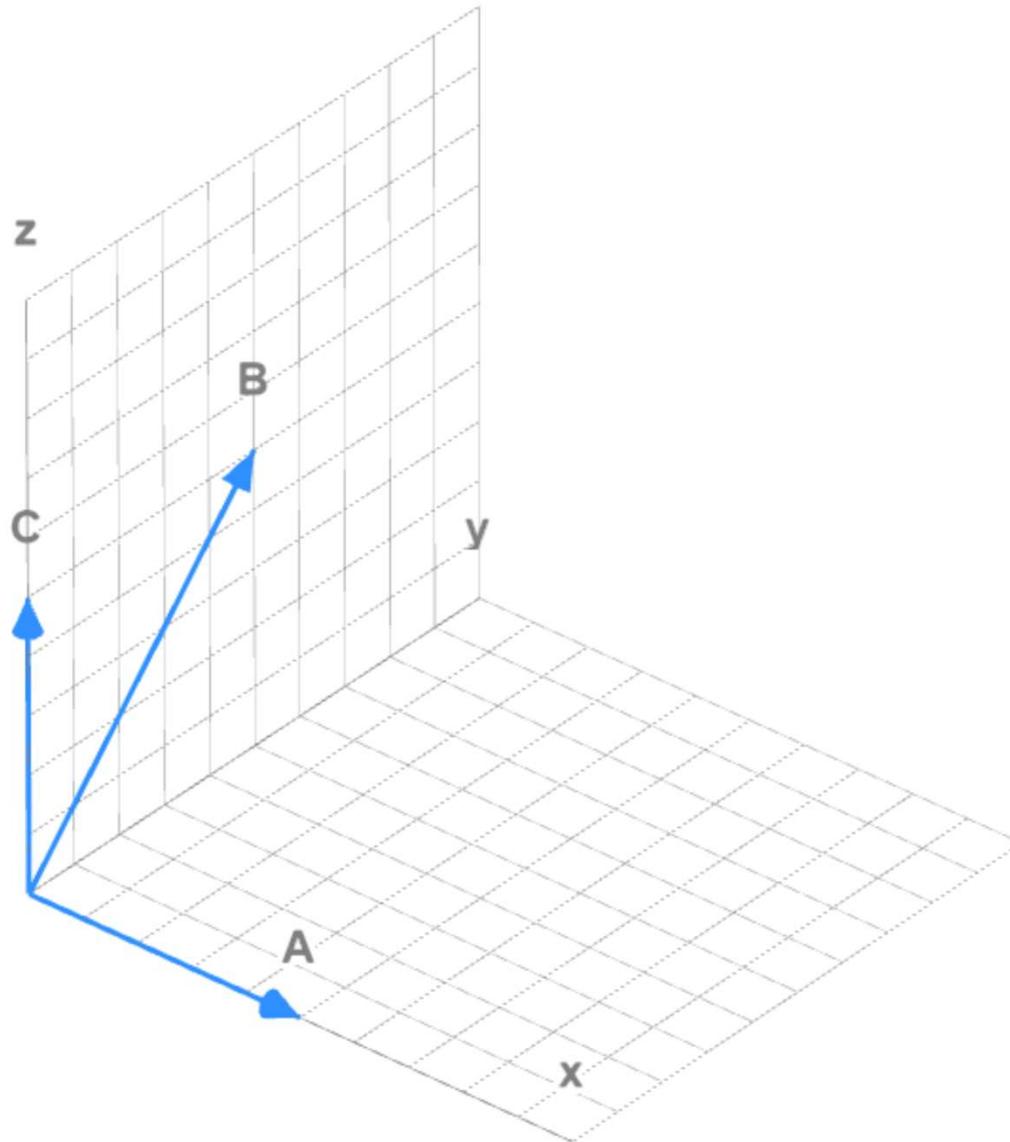
)

# **Similar Documents**

=

# **Similar Structure**

# Document Vectors in Vector Space



Note: vector space can be N-dimensional (N - feature vector length).

**How similar are two documents?**

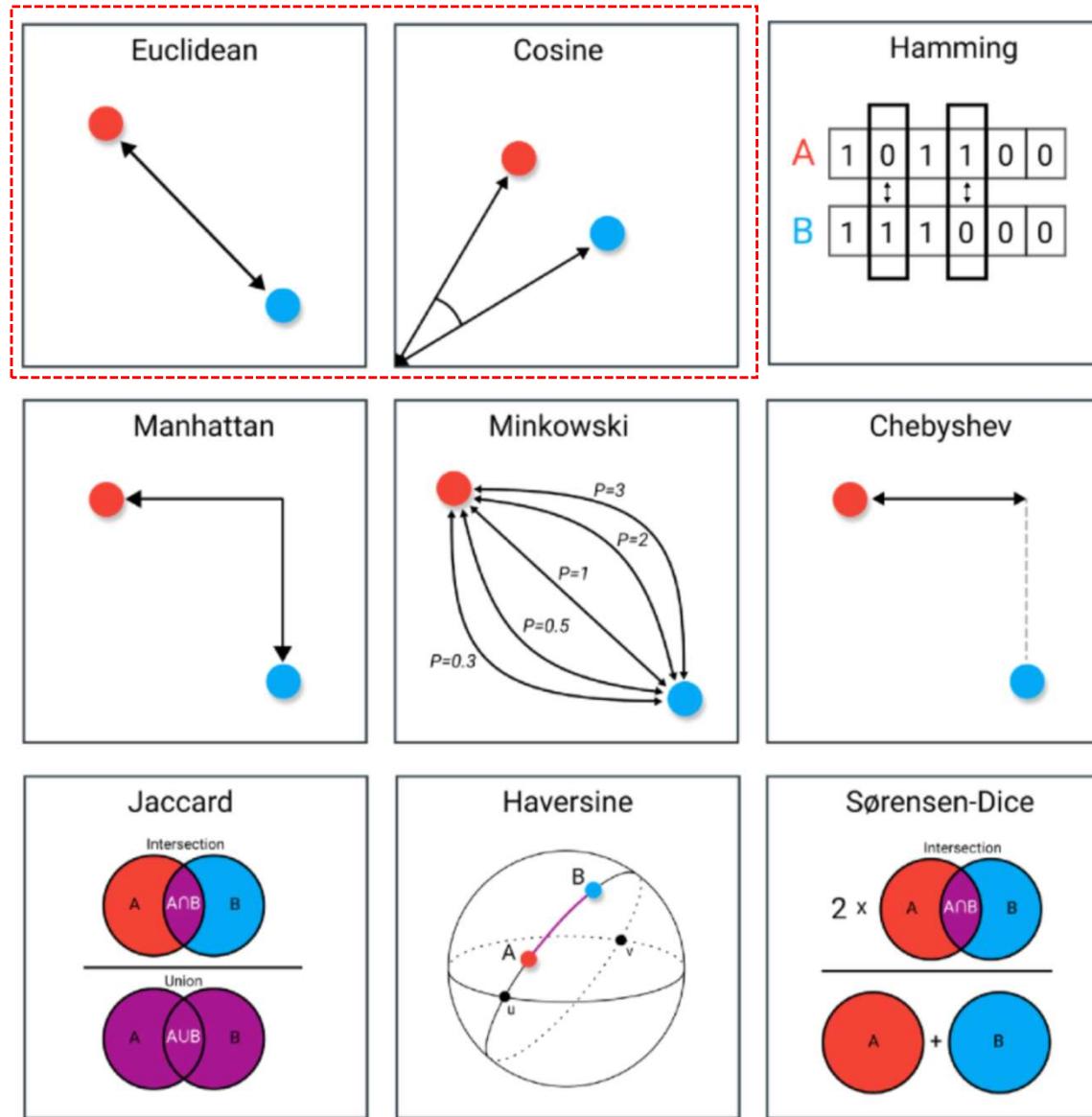
=

**How similar are their structures?**

=

**How close (in a vector space) are  
points defined by their document  
vectors**

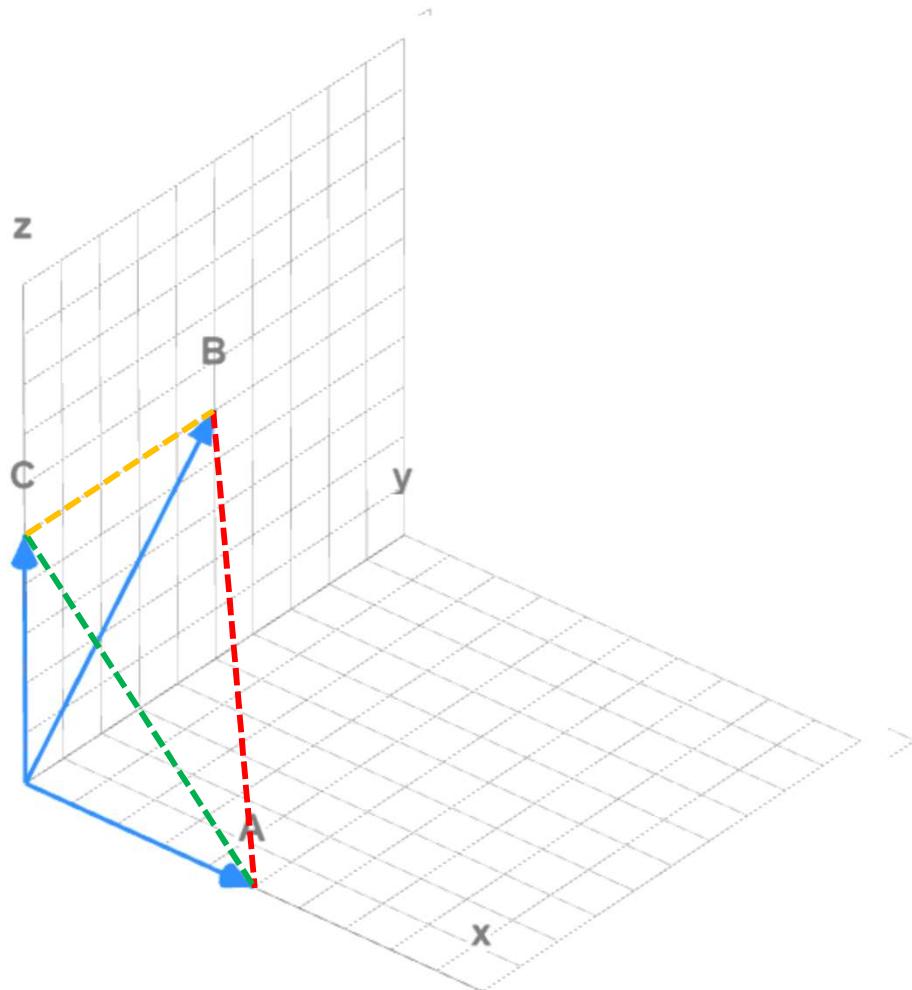
# Distance Measures



Source: <https://towardsdatascience.com/9-distance-measures-in-data-science-918109d069fa>

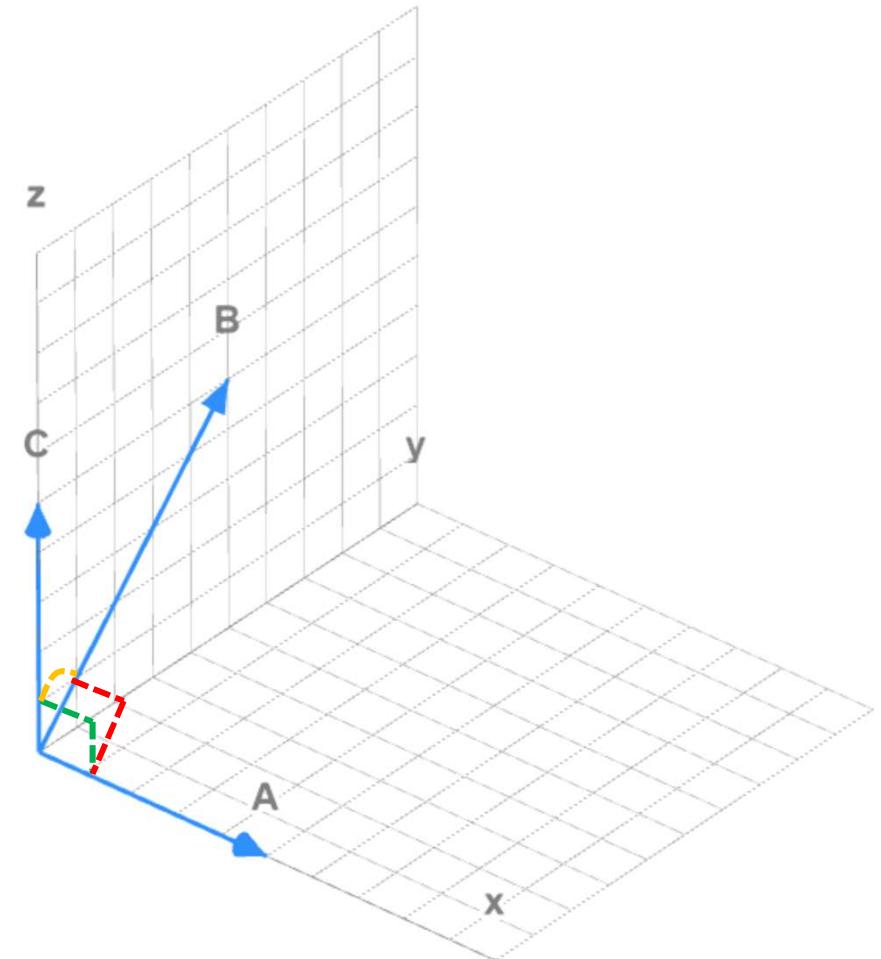
# Distance Measures

Euclidean distance



$$D(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Cosine similarity



$$D(x, y) = \cos(\theta) = \frac{x \cdot y}{\|x\| \|y\|}$$

# Bag of Words: Limitations

- Word locations ignored
- Semantics ignored
  - similar / synonymous words could become distinct features

Pre-defined Vocabulary (features):

Word 1	soccer	Word 3	Word 4	football	Word 6	...	Word N
--------	--------	--------	--------	----------	--------	-----	--------

- similar sentences will have different vectors

buy	old	desktop	purchase	used	PC	...	Word N
-----	-----	---------	----------	------	----	-----	--------

“Buy old desktop” Vector [0-word absent | 1-word present]:

1	1	1	0	0	0	...	0
---	---	---	---	---	---	-----	---

“Purchase used PC” vector [0-word absent | 1-word present]:

0	0	0	1	1	1	...	0
---	---	---	---	---	---	-----	---

- New / unknown words | vocabulary range

# Text Classification: Definition

*Input:*

- a document  $x$
- a fixed set of classes  $Y = \{y_1, y_2, \dots, y_J\}$

*Output:* a predicted class  $y \in Y$

# Classification: Key Question

Given a document (email, tweet, etc.):



which category / class does it belong to?

# Classification: Key Question

Given a document (email, tweet, etc.):



which category / class is the best  
**(predicted) match** for this document?

# Classification: Key Question

Given a document (email, tweet, etc.):



which category / class is **the most probable** (= lowest error) for this document?

# Classification: Key Question

Given a document (email, tweet, etc.):



which category / class has **the highest**

$$P(y = \text{class} \mid x = \text{document})?$$

# Classification: Key Question

Which category / class has **the highest**

$$P(y = \text{class}_1 \mid x = \text{document}) = ???$$

$$P(y = \text{class}_2 \mid x = \text{document}) = ???$$

...

$$P(y = \text{class}_j \mid x = \text{document}) = ???$$

Calculate all probabilities ...

# Classification: Key Question

Which category / class has **the highest**

$$P(y = \text{class}_1 \mid x = \text{document}) = 0.1$$

$$P(y = \text{class}_2 \mid x = \text{document}) = 0.3$$

...

$$P(y = \text{class}_j \mid x = \text{document}) = 0.2$$

... and pick the maximum P().

# Classification: Key Question

Which category / class has **the highest**

$$P(y = \text{class}_1 \mid x = \text{document}) = 0.1$$

$$P(y = \text{class}_2 \mid x = \text{document}) = 0.3$$

...

$$P(y = \text{class}_j \mid x = \text{document}) = 0.2$$

Corresponding class → most probable.

# Classification: Key Question

Which category / class has **the highest**

$$P(y = \text{class}_1 \mid x = \text{document}) = ???$$

$$P(y = \text{class}_2 \mid x = \text{document}) = ???$$

...

$$P(y = \text{class}_j \mid x = \text{document}) = ???$$

Calculate all probabilities ... **but how?**

# Bayes' Rule

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

# Bayes' Rule: Another Interpretation

Another way to think about Bayes' rule: it allows us to update the hypothesis  $H$  in light of some new data/evidence  $e$ .

$$P(H | e) = \frac{P(e | H) * P(H)}{P(e)}$$

$$P(\text{Hypothesis} | \text{evidence}) = \frac{P(\text{evidence} | \text{Hypothesis}) * P(\text{Hypothesis})}{P(\text{evidence})}$$

where:

- $P(H)$  - probability of the Hypothesis  $H$  being true BEFORE we see new data/evidence  $e$  (prior probability)
- $P(H | e)$  - probability of the Hypothesis  $H$  being true AFTER we see new data/evidence  $e$  (posterior probability)
- $P(e | H)$  - probability of new data/evidence  $e$  being true under the Hypothesis  $H$  (likelihood)
- $P(e)$  - probability of new data/evidence  $e$  being true under ANY hypothesis (normalizing constant)

# Bayes' Rule: Another Interpretation

Another way to think about Bayes' rule: it allows us to update the hypothesis  $H$  in light of some new data/evidence  $e$ .

$$P(H | e) = \frac{P(e | H) * P(H)}{P(e)}$$

$$P(\text{Hypothesis} | \text{evidence}) = \frac{P(\text{evidence} | \text{Hypothesis}) * P(\text{Hypothesis})}{P(\text{evidence})}$$

$$P(y | x) = \frac{P(x | y) * P(y)}{P(e)}$$

$$P(\text{class} | \text{document}) = \frac{P(\text{document} | \text{class}) * P(\text{class})}{P(\text{document})}$$

$$P(y | x_1, x_2, \dots, x_N) = \frac{P(x_1, x_2, \dots, x_N | y) * P(y)}{P(x_1, x_2, \dots, x_N)}$$

for example:

$$P(y = y_k | x_1 = 1, x_2 = 3, \dots, x_N = 0) = \frac{P(x_1 = 1, x_2 = 3, \dots, x_N = 0 | y = y_k) * P(y = y_k)}{P(x_1 = 1, x_2 = 3, \dots, x_N = 0)}$$

# Bayes' Rule

$$posterior = \frac{likelihood * prior}{evidence}$$

# Bayes' Rule

$$P(y | x) = \frac{P(x | y) * P(y)}{P(x)}$$

$$P(\text{Category} | \text{Document}) = \frac{P(\text{Document} | \text{Category}) * P(\text{Category})}{P(\text{Document})}$$

$$P(\text{Category} | \text{Instance}) = \frac{P(\text{Instance} | \text{Category}) * P(\text{Category})}{P(\text{Instance})}$$

$$P(\text{Category} | \text{Sample}) = \frac{P(\text{Sample} | \text{Category}) * P(\text{Category})}{P(\text{Sample})}$$

# Classification: Conditional Probability

$$P(y | x) = \frac{P(x | y) * P(y)}{P(x)}$$

$\mathbf{x} = x_1, x_2, \dots, x_N$ , so:

$$P(y | x_1 \wedge x_2 \wedge \dots \wedge x_N) = \frac{P(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * P(y)}{P(x_1 \wedge x_2 \wedge \dots \wedge x_N)}$$

# Classification: Conditional Probability

$$P(y | x) = \frac{P(x | y) * P(y)}{P(x)}$$

$\mathbf{x} = x_1, x_2, \dots, x_N$ , so:

How to calculate?

$$P(y | x_1 \wedge x_2 \wedge \dots \wedge x_N) = \frac{P(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * P(y)}{P(x_1 \wedge x_2 \wedge \dots \wedge x_N)}$$

constant

# Classifier

$$y_{MAP} = \underset{y \in Y}{argmax} (\mathbf{P}(y | \mathbf{x})) = \underset{y \in Y}{argmax} \left( \frac{\mathbf{P}(\mathbf{x} | y) * \mathbf{P}(y)}{\mathbf{P}(\mathbf{x})} \right)$$

$\mathbf{x} = x_1, x_2, \dots, x_N$ , so:

$$y_{MAP} = \underset{y \in Y}{argmax} \left( \frac{\mathbf{P}(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * \mathbf{P}(y)}{\mathbf{P}(x_1 \wedge x_2 \wedge \dots \wedge x_N)} \right)$$

constant | we can drop

$$y_{MAP} \propto \underset{y \in Y}{argmax} (\mathbf{P}(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * \mathbf{P}(y))$$

**MAP:** Maximum a posteriori (corresponds to the most likely class).

# Classifier

$$y_{MAP} \propto \underset{y \in Y}{argmax} \quad (P(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * P(y))$$

proportional

Prior

Likelihood

The diagram illustrates the Maximum a Posteriori (MAP) formula. The formula is shown as  $y_{MAP} \propto \underset{y \in Y}{argmax} (P(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * P(y))$ . A red dashed box encloses the term  $P(x_1 \wedge x_2 \wedge \dots \wedge x_N | y)$ , which is labeled 'Likelihood'. A green dashed box encloses the term  $P(y)$ , which is labeled 'Prior'. A black arrow points from the word 'proportional' up to the  $\propto$  symbol. Another black arrow points from the word 'Likelihood' down to the red box. A green dotted arrow points from the word 'Prior' down to the green box.

**MAP:** Maximum a posteriori (corresponds to the most likely class).

# Classifier

$$y_{MAP} \propto \underset{y \in Y}{argmax} (P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y))$$

How to  
calculate?

**MAP:** Maximum a posteriori (corresponds to the most likely class).

# Classifier

$$y_{MAP} \propto \underset{y \in Y}{argmax} (P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y))$$

How to  
calculate?

**MAP:** Maximum a posteriori (corresponds to the most likely class).

# Conditional Probability (Product Rule)

$$P(A \wedge B) = P(A | B) * P(B)$$

so:

$$P(A | B) * P(B) = P(A \wedge B)$$

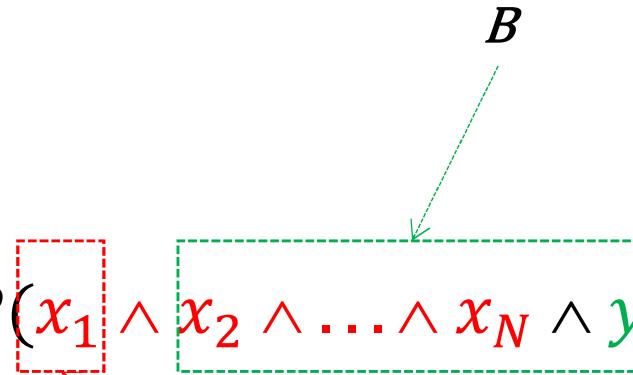
# Conditional Probability (Product Rule)

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y) = P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y)$$

so:

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y) = P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y)$$

# Conditional Probability (Product Rule)

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y) = P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y)$$


*and*

*A*

$$P(A \wedge B) = P(A \mid B) * P(B)$$

*so:*

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y) = P(x_1 \mid x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 \wedge \dots \wedge x_N \wedge y)$$

# Chain Rule

Conditional probabilities can be used to decompose conjunctions using the chain rule. For any events  $f_1, f_2, \dots, f_n$ :

$$P(f_1 \wedge f_2 \wedge \dots \wedge f_n) =$$

$$P(f_1) *$$

$$P(f_2 | f_1) *$$

$$P(f_3 | f_1 \wedge f_2) *$$

...

$$P(f_n | f_1 \wedge \dots \wedge f_{n-1}) =$$

$$= \prod_{i=1}^n P(f_i | f_1 \wedge \dots \wedge f_{i-1})$$

# Expansion

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y) =$$

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 \wedge \dots \wedge x_N \wedge y) =$$

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * P(x_3 \wedge \dots \wedge x_N \wedge y) =$$

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * P(x_3 | x_4 \wedge \dots \wedge x_N \wedge y) * P(x_4 \wedge \dots \wedge x_N \wedge y) =$$

...

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * \dots * P(x_N | y) * P(y)$$

# Independence

Assume that the knowledge of the truth of one proposition Y, does not affect the agent's belief in another proposition, X, in the context of other propositions Z. We say that X is **independent** of Y given Z.

# Conditional Independence

**Random variable X is conditionally independent of random variable Y given Z if for all  $x \in D_x$ , for all  $y \in D_y$ , and for all  $z \in D_z$ , such that**

$$P(Y = y \wedge Z = z) > 0 \text{ and } P(Y = y' \wedge Z = z) > 0$$

$$P(X = x \mid Y = y \wedge Z = z) = P(X = x \mid Y = y' \wedge Z = z)$$

**In other words, given a value of Z, knowing Y's value DOES NOT affect your belief in the value of X.**

# Conditional Independence

The following four statements are equivalent as long as conditional probabilities:

1. X is conditionally independent of Y given Z
2. Y is conditionally independent of X given Z
3.  $P(X | Y, Z) = P(X | Z)$
4.  $P(X, Y | Z) = P(X | Z) * P(Y | Z)$

# Naive Bayes Assumption

$$\begin{aligned} P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y) &= \\ P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 \wedge \dots \wedge x_N \wedge y) &= \\ P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * P(x_3 \wedge \dots \wedge x_N \wedge y) &= \\ P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * P(x_3 | x_4 \wedge \dots \wedge x_N \wedge y) * P(x_4 \wedge \dots \wedge x_N \wedge y) &= \\ \dots \\ P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * \dots * P(x_N | y) * P(y) \end{aligned}$$

Now let's assume that all events  $x_1, x_2, \dots, x_N$  are **mutually independent** (not true in reality) and **conditionally independent given  $y$**   $\rightarrow$  **Naive Bayes assumption.**

Under this assumption:

$$P(x_i | x_{i+1} \wedge \dots \wedge x_N \wedge y) = P(x_i | y)$$

# Naive Bayes Assumption

**Under Naive Bayes assumption:**

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y) =$$

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 \wedge \dots \wedge x_N \wedge y) =$$

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * P(x_3 \wedge \dots \wedge x_N \wedge y) =$$

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * P(x_3 | x_4 \wedge \dots \wedge x_N \wedge y) * P(x_4 \wedge \dots \wedge x_N \wedge y) =$$

...

$$P(x_1 | x_2 \wedge \dots \wedge x_N \wedge y) * P(x_2 | x_3 \wedge \dots \wedge x_N \wedge y) * \dots * P(x_N | y) * P(y)$$

**becomes:**

$$P(x_1 \wedge x_2 \wedge \dots \wedge x_N \wedge y) =$$

$$P(x_1 | y) * P(x_2 | y) * P(x_3 | y) * \dots * P(x_{N-1} | y) * P(x_N | y) * P(y) =$$

$$P(y) * [P(x_1 | y) * P(x_2 | y) * P(x_3 | y) * \dots * P(x_{N-1} | y) * P(x_N | y)] =$$

$$P(y) * \prod_{i=1}^N P(x_i | y)$$

# Naive Bayes Classifier

Under Naive Bayes assumption:

$$y_{MAP} \propto \underset{y \in Y}{argmax} (P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y))$$

becomes:

$$y_{MAP} \propto \underset{y \in Y}{argmax} \left( P(y) * \prod_{i=1}^N P(x_i \mid y) \right)$$

MAP: Maximum a posteriori (corresponds to the most likely class).

# Naive Bayes Classifier

Under Naive Bayes assumption:

$$y_{MAP} \propto \underset{y \in Y}{argmax} (P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y))$$

becomes:

$$y_{MAP} \propto \underset{y \in Y}{argmax} \left( P(y) * \prod_{i=1}^N P(x_i \mid y) \right)$$

MAP: Maximum a posteriori (corresponds to the most likely class).

# Naive Bayes Classifier

Under Naive Bayes assumption:

$$y_{MAP} \propto \underset{y \in Y}{argmax} (P(x_1 \wedge x_2 \wedge \dots \wedge x_N \mid y) * P(y))$$

becomes:

$$y_{MAP} \propto \underset{y \in Y}{argmax} \left( P(y) * \prod_{i=1}^N P(x_i \mid y) \right)$$

How to calculate?

MAP: Maximum a posteriori (corresponds to the most likely class).

# Text Classification: Supervised ML

*Input:*

- a document  $\mathbf{x}$
- a fixed set of classes  $Y = \{y_1, y_2, \dots, y_J\}$
- a training set of  $N$  hand-labeled documents  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$

*Output:*

- a learned classifier  $h: \mathbf{x} \rightarrow y$  ( $y = h(\mathbf{x})$ )

# Text Classification: Classifier

category/class =  $h(\text{document})$

Learned Classifier model  
(hypothesis)

# Text Classification: Classifier

$$y = h(x)$$

Learned Classifier model  
(hypothesis)

# Text Classification: Supervised ML

*Input:*

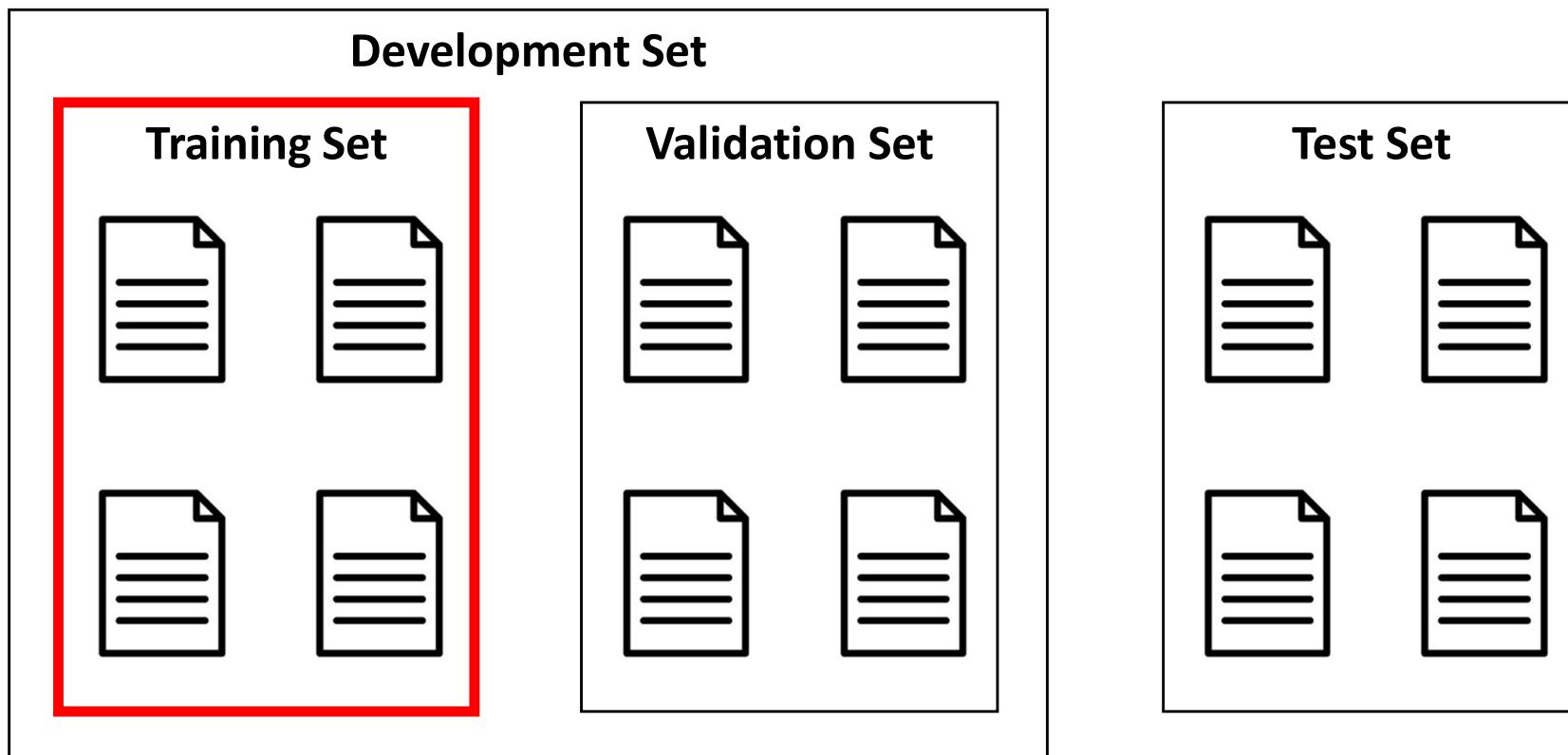
- a document  $\mathbf{x}$
- a fixed set of classes  $Y = \{y_1, y_2, \dots, y_J\}$
- a **training set** of  $N$  hand-labeled documents  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$

*Output:*

- a learned classifier  $h: \mathbf{x} \rightarrow y$  ( $y = h(\mathbf{x})$ )

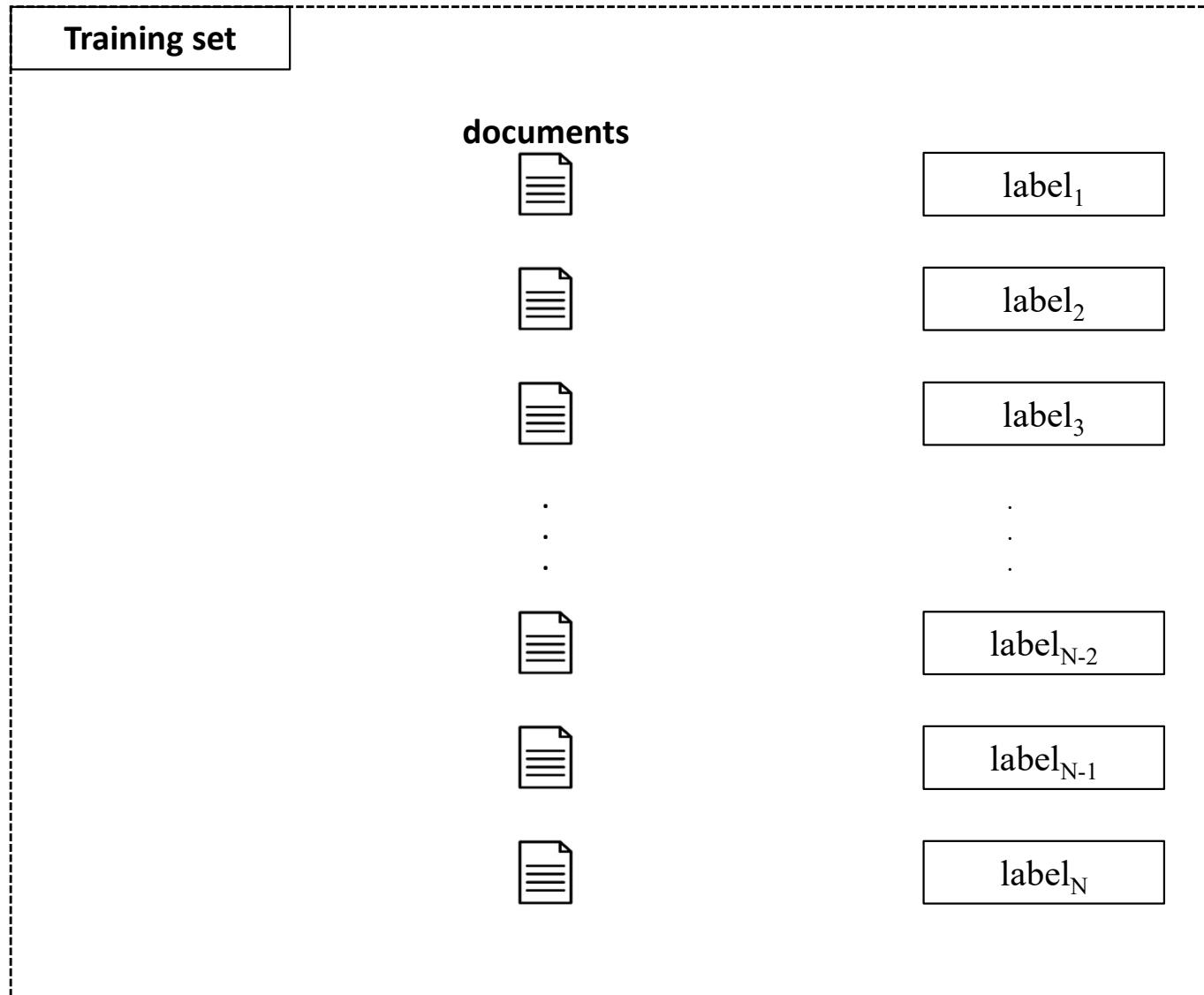
# Corpus: Training / Validation / Test

Corpus



Typical training / validation / test set split for a text corpora

# Text Classification: Training Set

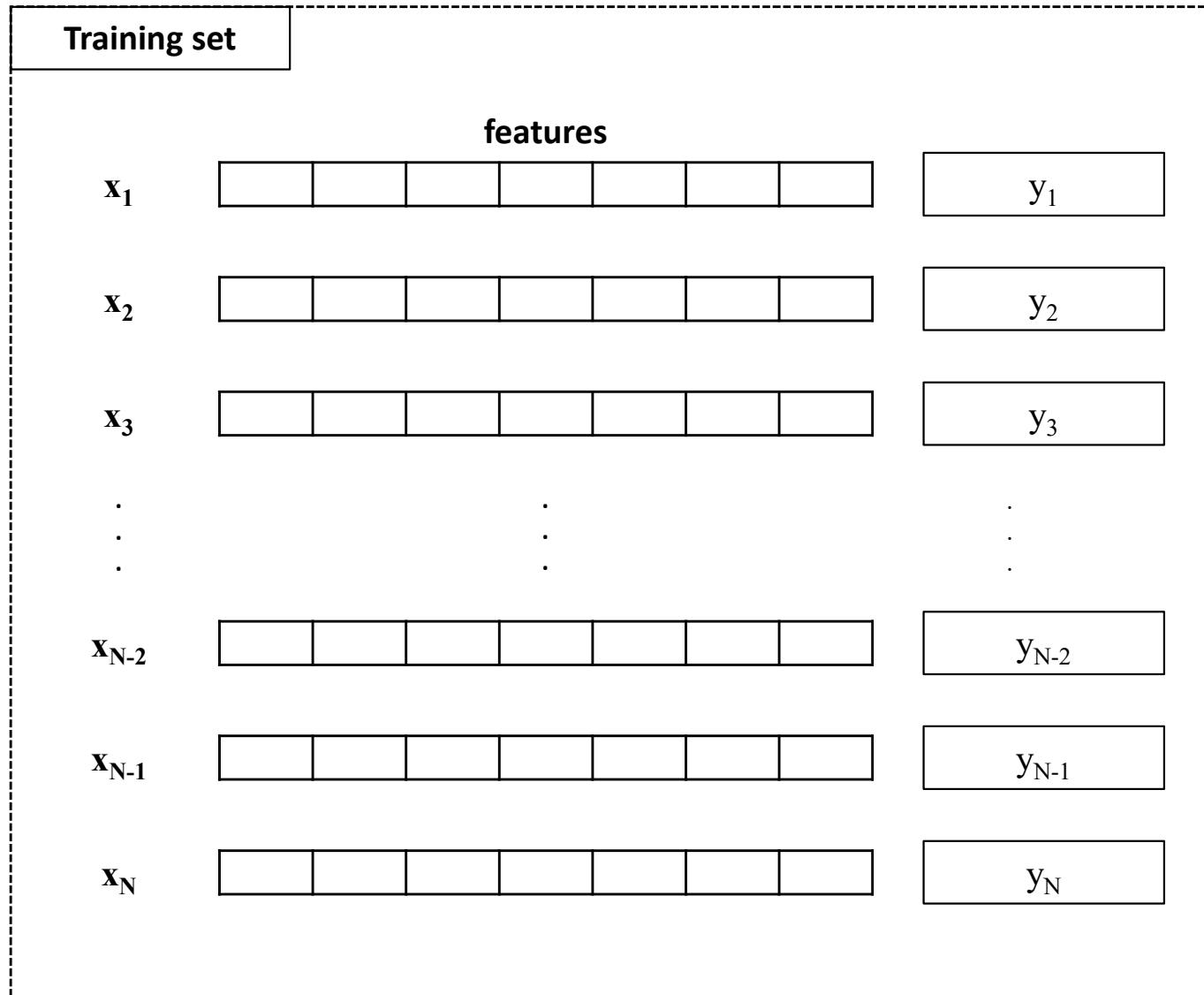


# Text Classification: Training Set

Training set							
features (bag of words)							
$x_1$							label <sub>1</sub>
$x_2$							label <sub>2</sub>
$x_3$							label <sub>3</sub>
.							.
.							.
.							.
$x_{N-2}$							label <sub>N-2</sub>
$x_{N-1}$							label <sub>N-1</sub>
$x_N$							label <sub>N</sub>

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**)

# Text Classification: Training Set



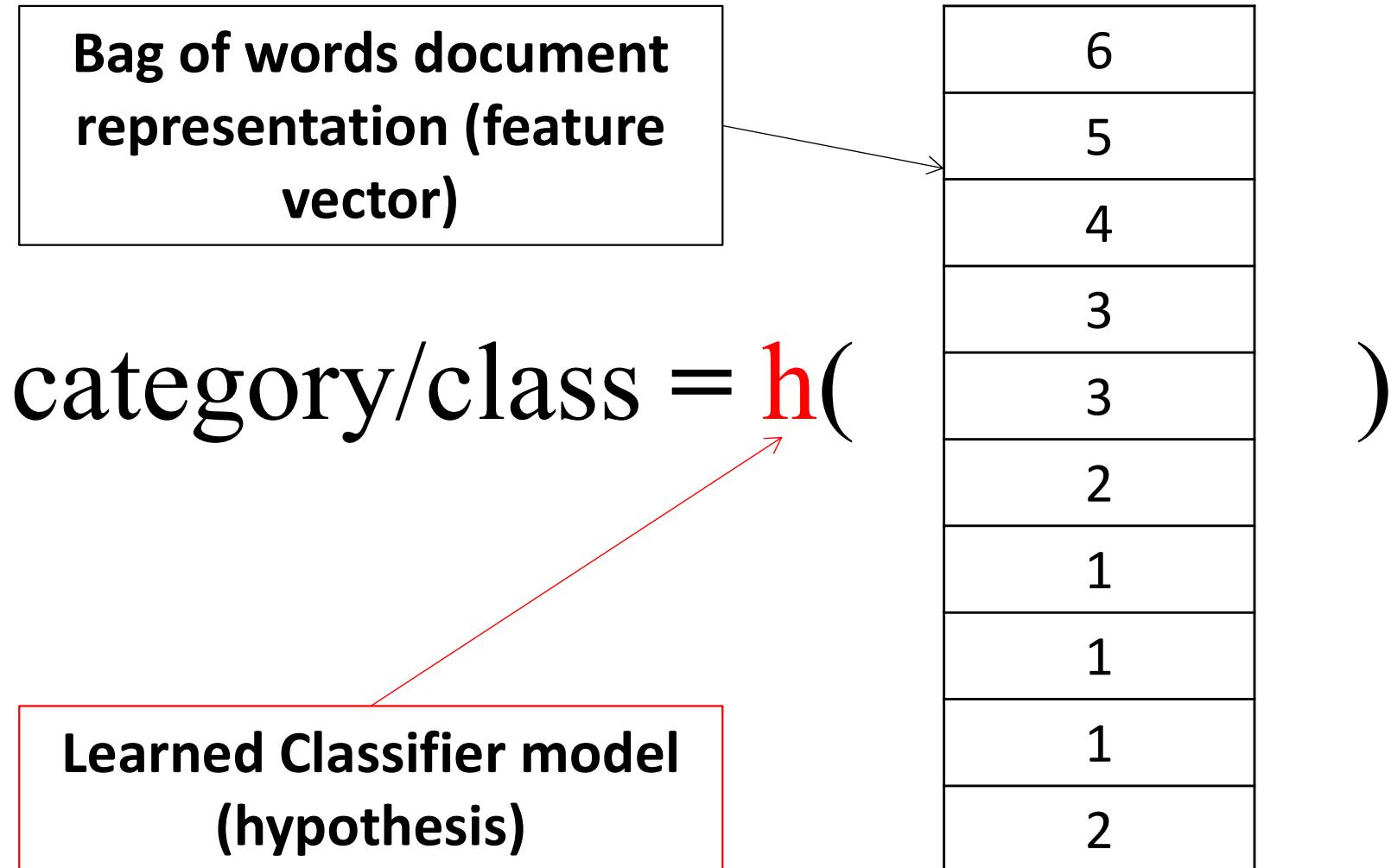
$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Training Set

Training set								
Vocabulary $V$								
$x_1$	word1	rolex	word3	replica	word5	word6	word7	
	0	0	1	0	1	1	1	
								$y_1 = \text{HAM}$
$x_2$	1	0	1	1	0	1	1	
								$y_2 = \text{HAM}$
$x_3$	0	1	0	1	0	1	1	
								$y_3 = \text{SPAM}$
.								.
.								.
.								.
$x_{N-2}$	1	1	1	1	0	1	1	
								$y_{N-2} = \text{HAM}$
$x_{N-1}$	1	1	0	1	0	0	1	
								$y_{N-1} = \text{SPAM}$
$x_N$	1	0	0	1	0	0	1	
								$y_N = \text{HAM}$

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Text Classification: Bag of Words



# Text Classification: Bag of Words

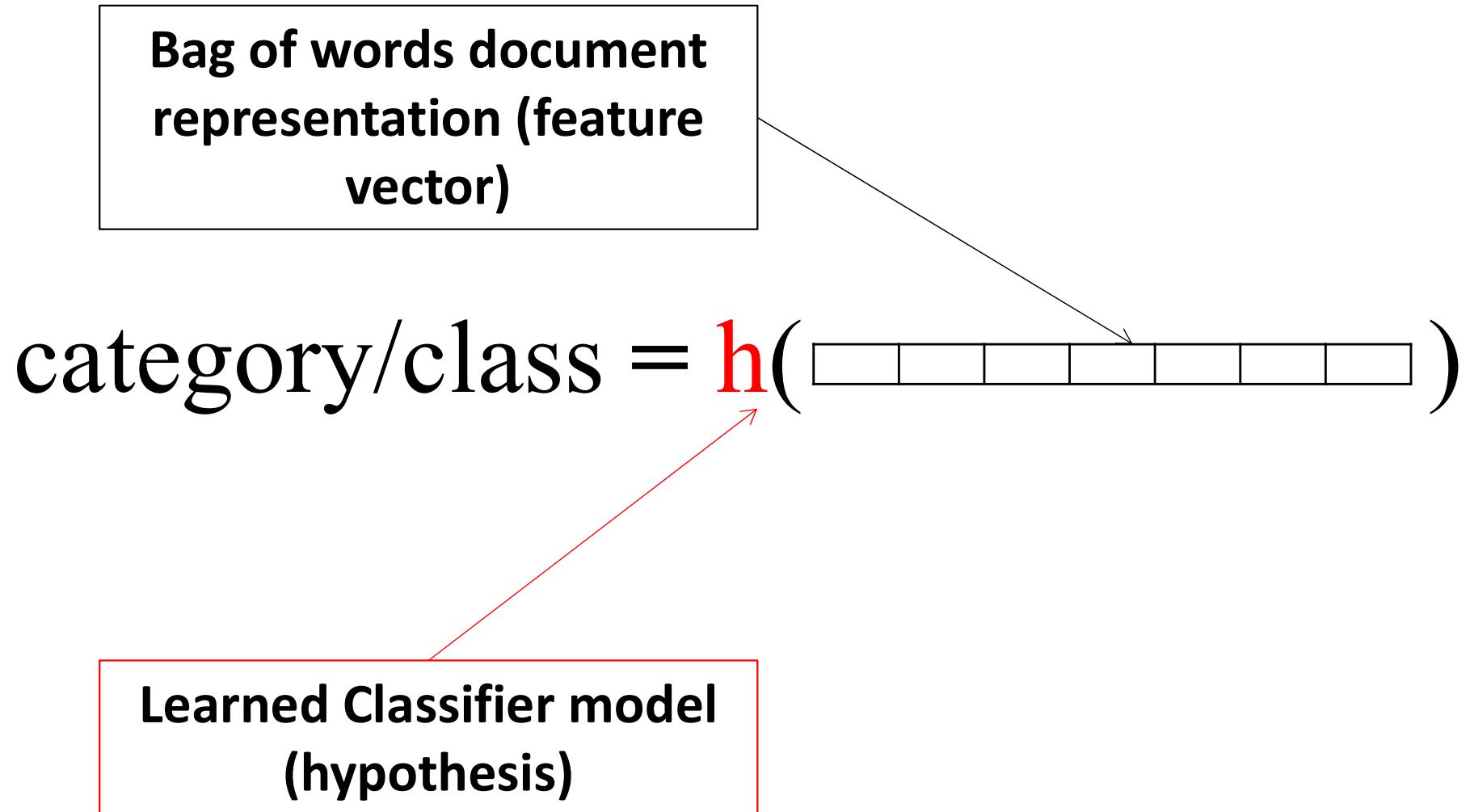
Bag of words **binary**  
document representation  
(feature vector)

category/class =  $h($

Learned Classifier model  
(hypothesis)

1
1
0
1
0
0
1
1
1
0

# Text Classification: Bag of Words



# Spam Detection: Learning

Training set

Vocabulary  $V$

	word1	rolex	word3	replica	word5	word6	word7
$x_1$	0	0	1	0	1	1	1

$y_1 = \text{HAM}$

$x_2$	1	0	1	1	0	1	1
-------	---	---	---	---	---	---	---

$y_2 = \text{HAM}$

$x_3$	0	1	0	1	0	1	1
-------	---	---	---	---	---	---	---

$y_3 = \text{SPAM}$

.	.	.	.	.	.	.	.
$x_{N-2}$	1	1	1	1	0	1	1

$y_{N-2} = \text{HAM}$

$x_{N-1}$	1	1	0	1	0	0	1
-----------	---	---	---	---	---	---	---

$y_{N-1} = \text{SPAM}$

$x_N$	1	0	0	1	0	0	1
-------	---	---	---	---	---	---	---

$y_N = \text{HAM}$

Learning

Naive Bayes Classifier:

$$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i | y) \right)$$

Probability estimates (Maximum Likelihood estimation):

$$P(y_k) = \frac{N_{\text{samples labeled } y_k}}{N}$$

$$P(x_i | y_k) = \frac{\text{count}(x_i, y_k)}{\sum_{x \in V} \text{count}(x, y_k)}$$

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Learning

Training set

Vocabulary  $V$

	word1	rolex	word3	replica	word5	word6	word7
$x_1$	0	0	1	0	1	1	1

$y_1 = \text{HAM}$

$x_2$	1	0	1	1	0	1	1
-------	---	---	---	---	---	---	---

$y_2 = \text{HAM}$

$x_3$	0	1	0	1	0	1	1
-------	---	---	---	---	---	---	---

$y_3 = \text{SPAM}$

$x_4$	1	1	1	1	0	0	0
-------	---	---	---	---	---	---	---

$y_4 = \text{HAM}$

$x_5$	1	1	1	1	0	1	1
-------	---	---	---	---	---	---	---

$y_5 = \text{HAM}$

$x_6$	1	1	0	1	0	0	1
-------	---	---	---	---	---	---	---

$y_6 = \text{SPAM}$

$x_7$	1	0	0	1	0	0	1
-------	---	---	---	---	---	---	---

$y_7 = \text{HAM}$

Learning

Naive Bayes Classifier:

$$y_{MAP} \propto \underset{y \in Y}{argmax} \left( P(y) * \prod_{i=1}^N P(x_i | y) \right)$$

Probability estimates (Maximum Likelihood estimation):

$$P(y = \text{HAM}) = \frac{N_{\text{samples labeled HAM}}}{N} = \frac{5}{7}$$

$$P(y = \text{SPAM}) = \frac{N_{\text{samples labeled SPAM}}}{N} = \frac{2}{7}$$

$$P(x_i = \text{rolex} | y = \text{SPAM}) = \frac{\text{count}(x_i = \text{rolex}, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{2}{8}$$

and so on...

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Learning

Training set

Vocabulary  $V$

	word1	rolex	word3	replica	word5	word6	word7
$x_1$	0	0	1	0	1	1	1

$y_1 = \text{HAM}$

$x_2$	1	0	1	1	0	1	1
-------	---	---	---	---	---	---	---

$y_2 = \text{HAM}$

$x_3$	0	1	0	1	0	1	1
-------	---	---	---	---	---	---	---

$y_3 = \text{SPAM}$

.	.	.	.	.	.	.	.
---	---	---	---	---	---	---	---

$x_{N-2}$	1	1	1	1	0	1	1
-----------	---	---	---	---	---	---	---

$y_{N-2} = \text{HAM}$

$x_{N-1}$	1	1	0	1	0	0	1
-----------	---	---	---	---	---	---	---

$y_{N-1} = \text{SPAM}$

$x_N$	1	0	0	1	0	0	1
-------	---	---	---	---	---	---	---

$y_N = \text{HAM}$

Learning

Naive Bayes Classifier:

$$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i | y) \right)$$

Probability estimates:

$$P(y_k) = \frac{N_{\text{samples labeled } y_k}}{N}$$

or

- **equiprobable** (all classes have equal probability)

$$P(y = \text{HAM}) = P(y = \text{SPAM}) = 0.5$$

- can be determined by experts in the area

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Building a Naïve Bayes Classifier

# Classifier

$$y_{MAP} = \underset{y \in Y}{argmax} (\mathbf{P}(y | \mathbf{x})) = \underset{y \in Y}{argmax} \left( \frac{\mathbf{P}(\mathbf{x} | y) * \mathbf{P}(y)}{\mathbf{P}(\mathbf{x})} \right)$$

$\mathbf{x} = x_1, x_2, \dots, x_N$ , so:

$$y_{MAP} = \underset{y \in Y}{argmax} \left( \frac{\mathbf{P}(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * \mathbf{P}(y)}{\mathbf{P}(x_1 \wedge x_2 \wedge \dots \wedge x_N)} \right)$$

constant | we can drop

$$y_{MAP} \propto \underset{y \in Y}{argmax} (\mathbf{P}(x_1 \wedge x_2 \wedge \dots \wedge x_N | y) * \mathbf{P}(y))$$

**MAP:** Maximum a posteriori (corresponds to the most likely class).

# Naive Bayes Classifier: Assumptions

- All events (words)  $x_1, x_2, \dots, x_N$  are **mutually independent**
  - Bag-of-words representation: the order of the words in a document  $d$  makes no difference (repetitions do)
- All events (words)  $x_1, x_2, \dots, x_N$  are **conditionally independent given  $y$**  (category / class)
  - **words appear independently of each other, given the document category / class  $y$**  (e.g. if you see word “car”, the word “drive” is no more likely to appear than if you saw “dog”)

# Naive Bayes Classifier

category/class =  $\text{h}(\text{document})$

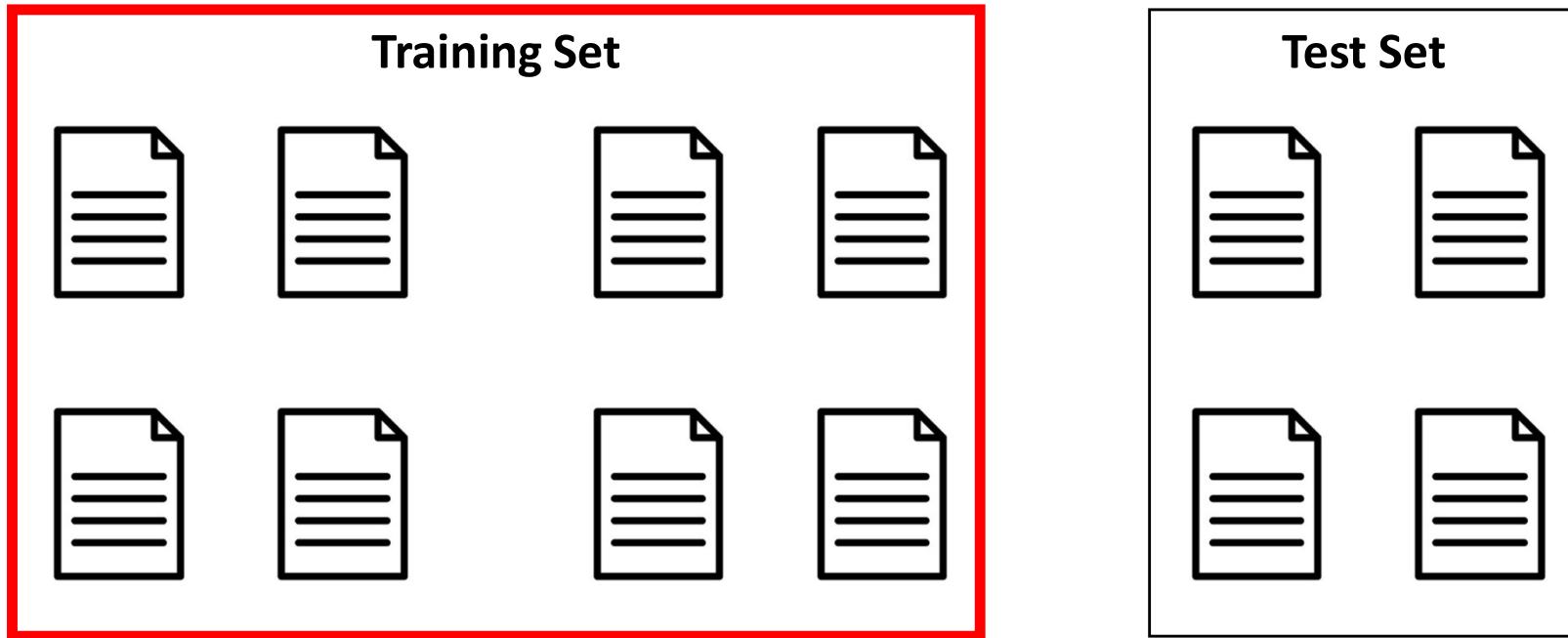
Finding model / hypothesis  $\text{h} \rightarrow$  Finding probabilities for  $y_{MAP}$

$$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i | y) \right)$$

MAP: Maximum a posteriori (corresponds to the most likely class).

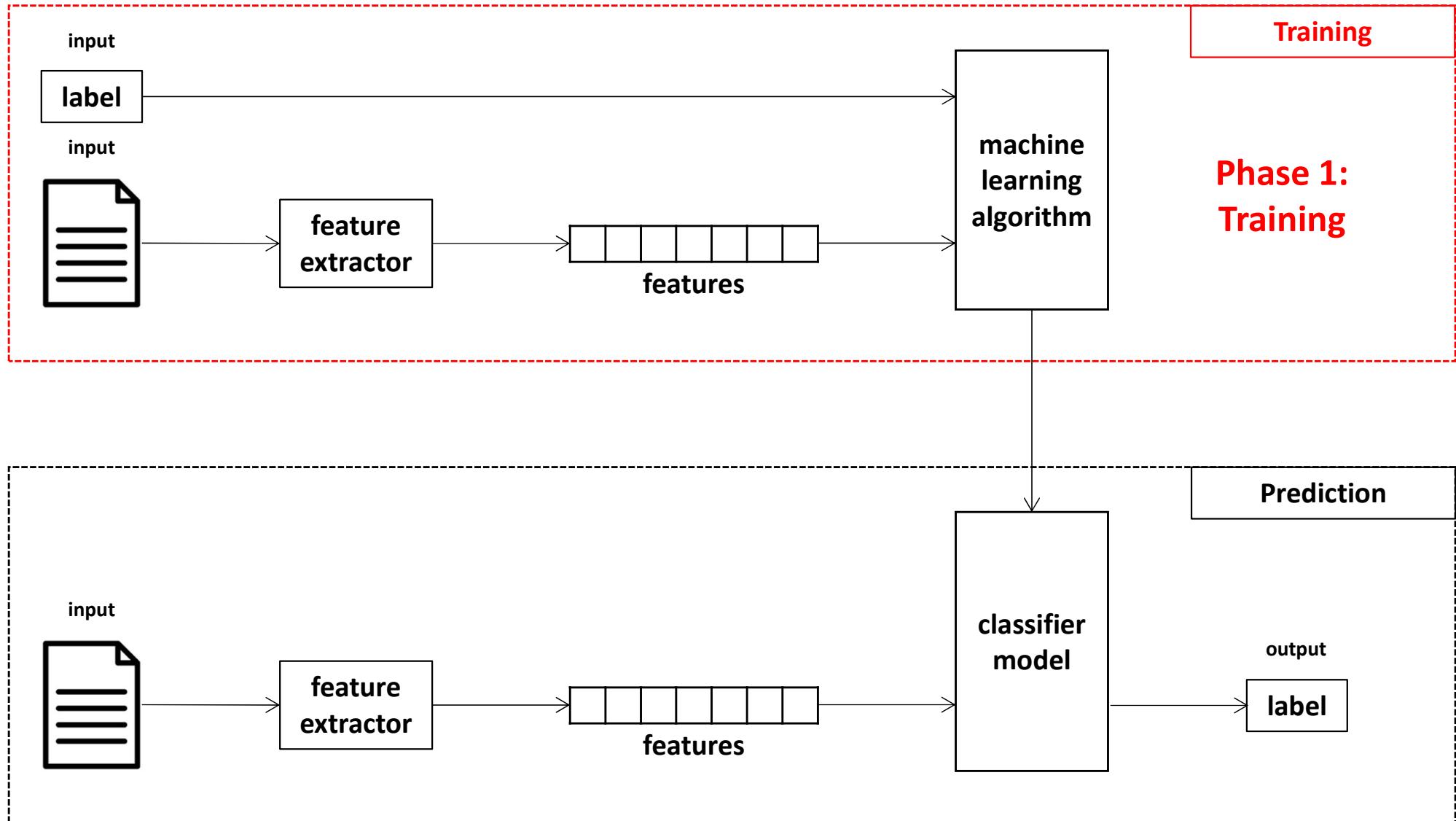
# Corpus: Training / Test

Corpus



Typical training / test set split for a text corpora [I will ignore validation set for the sake of an example]

# Supervised Learning with ML



# Spam Detection: Training Set

Training set							
Vocabulary $V$							
	I	rolex	own	replica	watch	buy	cheap
$x_1$	1	1	1	0	1	0	0
I own rolex watch							$y_1 = \text{HAM}$
$x_2$	1	0	1	0	1	0	0
I own watch							$y_2 = \text{HAM}$
$x_3$	0	1	0	1	0	1	1
buy cheap rolex replica							$y_3 = \text{SPAM}$
$x_4$	1	0	1	0	0	0	0
I own							$y_4 = \text{HAM}$
$x_5$	1	0	1	1	0	0	1
I own cheap replica							$y_5 = \text{HAM}$
$x_6$	0	1	0	1	0	0	1
cheap rolex replica							$y_6 = \text{SPAM}$
$x_7$	1	0	0	0	1	0	0
I watch							$y_7 = \text{HAM}$

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Learning

Training set

Vocabulary  $V$

	1	rolex	own	replica	watch	buy	cheap
$x_1$	1	1	1	0	1	0	0

$y_1 = \text{HAM}$

$x_2$	1	0	1	0	1	0	0
-------	---	---	---	---	---	---	---

$y_2 = \text{HAM}$

$x_3$	0	1	0	1	0	1	1
-------	---	---	---	---	---	---	---

$y_3 = \text{SPAM}$

$x_4$	1	0	1	0	0	0	0
-------	---	---	---	---	---	---	---

$y_4 = \text{HAM}$

$x_5$	1	0	1	1	0	0	1
-------	---	---	---	---	---	---	---

$y_5 = \text{HAM}$

$x_6$	0	1	0	1	0	0	1
-------	---	---	---	---	---	---	---

$y_6 = \text{SPAM}$

$x_7$	1	0	0	0	1	0	0
-------	---	---	---	---	---	---	---

$y_7 = \text{HAM}$

Learning

Probability estimates:

Naive Bayes Classifier:

$$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i | y) \right)$$

Probability estimates (Maximum Likelihood estimation):

$$P(y = \text{HAM}) = \frac{N_{\text{samples labeled HAM}}}{N}$$

$$P(y = \text{SPAM}) = \frac{N_{\text{samples labeled SPAM}}}{N}$$

$$\begin{aligned} P(x_i = \text{word} | y = \text{CLASS}) &= \\ &= \frac{\text{count}(x_i = \text{word}, y = \text{CLASS})}{\sum_{x \in V} \text{count}(x, y = \text{CLASS})} \end{aligned}$$

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Learning

Training set		Learning							
Vocabulary $V$									
$x_1$		1	rolex	own	replica	watch	buy	cheap	$y_1 = \text{HAM}$
		1	1	1	0	1	0	0	
$x_2$		1	0	1	0	1	0	0	$y_2 = \text{HAM}$
		1	0	1	0	1	0	0	
$x_3$		0	1	0	1	0	1	1	$y_3 = \text{SPAM}$
		0	1	0	1	0	1	1	
$x_4$		1	0	1	0	0	0	0	$y_4 = \text{HAM}$
		1	0	1	0	0	0	0	
$x_5$		1	0	1	1	0	0	1	$y_5 = \text{HAM}$
		1	0	1	1	0	0	1	
$x_6$		0	1	0	1	0	0	1	$y_6 = \text{SPAM}$
		0	1	0	1	0	0	1	
$x_7$		1	0	0	0	1	0	0	$y_7 = \text{HAM}$
		1	0	0	0	1	0	0	

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Learning

Training set

Vocabulary  $V$

	I	rolex	own	replica	watch	buy	cheap
$x_1$	1	1	1	0	1	0	0

$y_1 = \text{HAM}$

$x_2$	1	0	1	0	1	0	0
-------	---	---	---	---	---	---	---

$y_2 = \text{HAM}$

$x_3$	0	1	0	1	0	1	1
-------	---	---	---	---	---	---	---

$y_3 = \text{SPAM}$

$x_4$	1	0	1	0	0	0	0
-------	---	---	---	---	---	---	---

$y_4 = \text{HAM}$

$x_5$	1	0	1	1	0	0	1
-------	---	---	---	---	---	---	---

$y_5 = \text{HAM}$

$x_6$	0	1	0	1	0	0	1
-------	---	---	---	---	---	---	---

$y_6 = \text{SPAM}$

$x_7$	1	0	0	0	1	0	0
-------	---	---	---	---	---	---	---

$y_7 = \text{HAM}$

Learning

$$P(x_i = I \mid y = \text{HAM}) = \frac{\text{count}(x_i = I, y = \text{HAM})}{\sum_{x \in V} \text{count}(x, y = \text{HAM})} = \frac{5}{15}$$

$$P(x_i = \text{rolex} \mid y = \text{HAM}) = \frac{\text{count}(x_i = \text{rolex}, y = \text{HAM})}{\sum_{x \in V} \text{count}(x, y = \text{HAM})} = \frac{1}{15}$$

$$P(x_i = \text{own} \mid y = \text{HAM}) = \frac{\text{count}(x_i = \text{own}, y = \text{HAM})}{\sum_{x \in V} \text{count}(x, y = \text{HAM})} = \frac{4}{15}$$

$$P(x_i = \text{replica} \mid y = \text{HAM}) = \frac{\text{count}(x_i = \text{replica}, y = \text{HAM})}{\sum_{x \in V} \text{count}(x, y = \text{HAM})} = \frac{1}{15}$$

$$P(x_i = \text{watch} \mid y = \text{HAM}) = \frac{\text{count}(x_i = \text{watch}, y = \text{HAM})}{\sum_{x \in V} \text{count}(x, y = \text{HAM})} = \frac{3}{15}$$

$$P(x_i = \text{buy} \mid y = \text{HAM}) = \frac{\text{count}(x_i = \text{buy}, y = \text{HAM})}{\sum_{x \in V} \text{count}(x, y = \text{HAM})} = \frac{0}{15}$$

$$P(x_i = \text{cheap} \mid y = \text{HAM}) = \frac{\text{count}(x_i = \text{cheap}, y = \text{HAM})}{\sum_{x \in V} \text{count}(x, y = \text{HAM})} = \frac{1}{15}$$

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Learning

Training set

Vocabulary  $V$

	I	rolex	own	replica	watch	buy	cheap
$x_1$	1	1	1	0	1	0	0
$x_2$	1	0	1	0	1	0	0
$x_3$	0	1	0	1	0	1	1
$x_4$	1	0	1	0	0	0	0
$x_5$	1	0	1	1	0	0	1
$x_6$	0	1	0	1	0	0	1
$x_7$	1	0	0	0	1	0	0

$y_1 = \text{HAM}$

$y_2 = \text{HAM}$

$y_3 = \text{SPAM}$

$y_4 = \text{HAM}$

$y_5 = \text{HAM}$

$y_6 = \text{SPAM}$

$y_7 = \text{HAM}$

Learning

$$P(x_i = I \mid y = \text{SPAM}) = \frac{\text{count}(x_i = I, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{0}{7} = 0$$

$$P(x_i = \text{rolex} \mid y = \text{SPAM}) = \frac{\text{count}(x_i = \text{rolex}, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{2}{7}$$

$$P(x_i = \text{own} \mid y = \text{SPAM}) = \frac{\text{count}(x_i = \text{own}, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{0}{7} = 0$$

$$P(x_i = \text{replica} \mid y = \text{SPAM}) = \frac{\text{count}(x_i = \text{replica}, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{2}{7}$$

$$P(x_i = \text{watch} \mid y = \text{SPAM}) = \frac{\text{count}(x_i = \text{watch}, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{0}{7} = 0$$

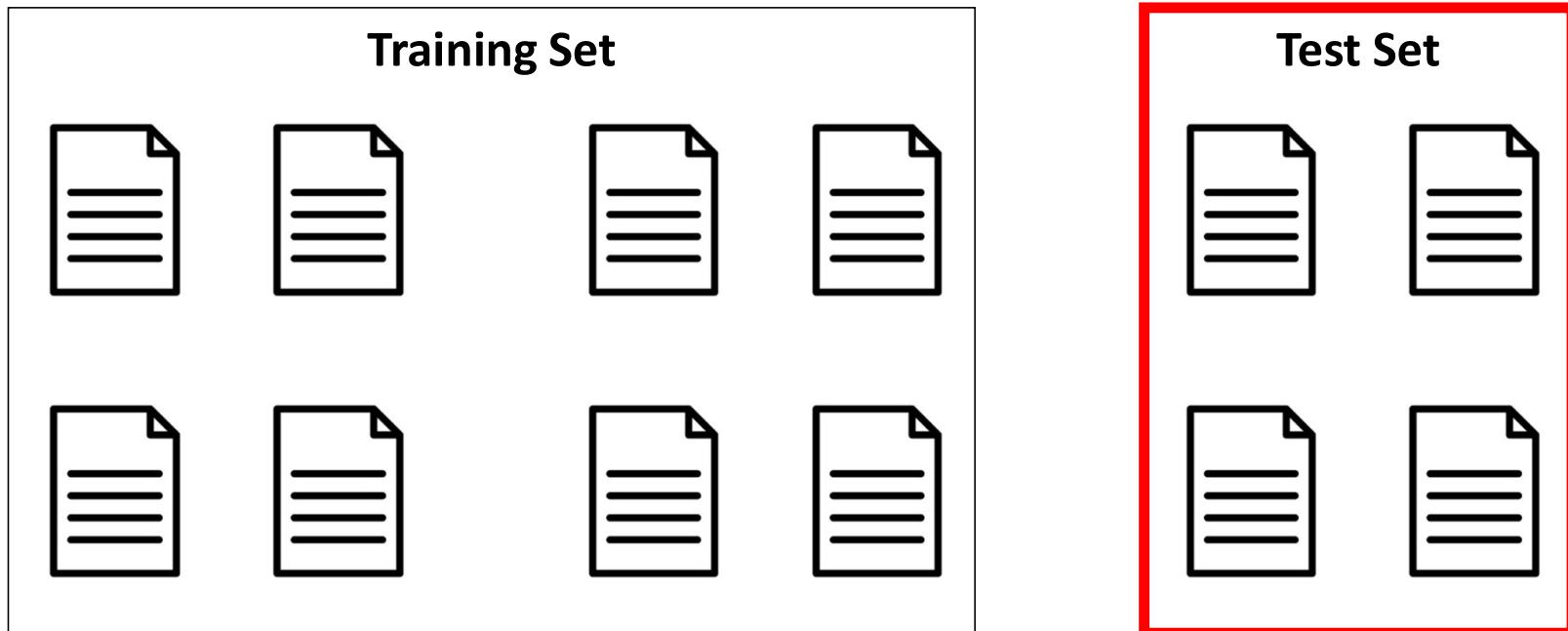
$$P(x_i = \text{buy} \mid y = \text{SPAM}) = \frac{\text{count}(x_i = \text{buy}, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{1}{7}$$

$$P(x_i = \text{cheap} \mid y = \text{SPAM}) = \frac{\text{count}(x_i = \text{cheap}, y = \text{SPAM})}{\sum_{x \in V} \text{count}(x, y = \text{SPAM})} = \frac{2}{7}$$

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

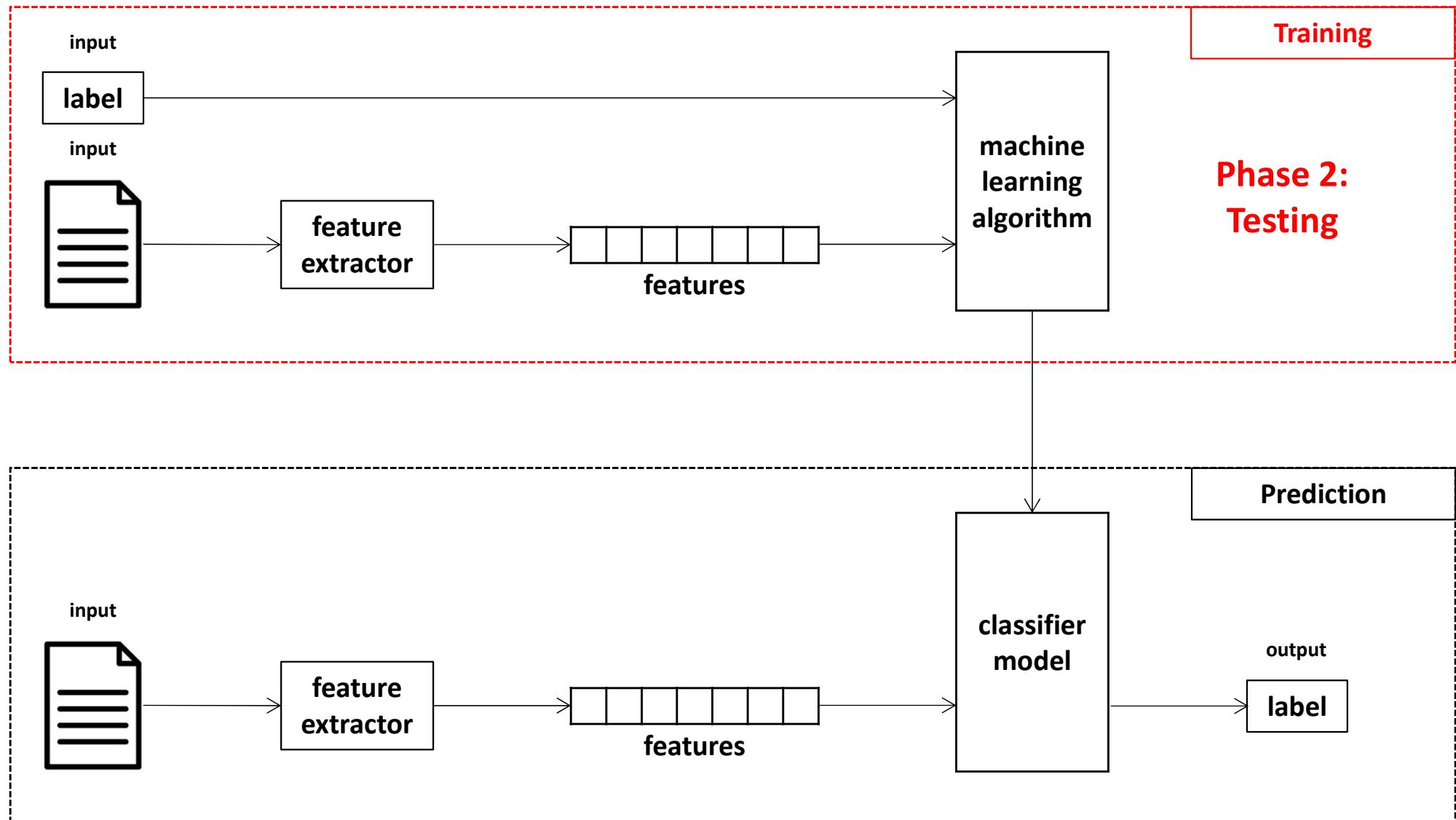
# Corpus: Training / Test

Corpus



Typical training / test set split for a text corpora [I will ignore validation set for the sake of an example]

# Supervised Learning with ML



# Spam Detection: Test Set

Test set							
Vocabulary $V$							
	I	rolex	own	replica	watch	buy	cheap
$x_8$	0	1	0	1	0	1	1
$y_8 = \text{SPAM}$							
buy cheap rolex replica rolex							
$x_9$	1	1	1	0	1	0	1
$y_9 = \text{HAM}$							
I own cheap rolex watch							
$x_{10}$	0	0	0	1	0	0	0
$y_{10} = \text{HAM}$							
replica							

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> buy cheap rolex replica rolex	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> I own cheap rolex watch	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> replica	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$ <p>category/class = <math>h(\text{document})</math></p>															
		$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$															
	<p>category/class = <math>h(x_8)</math></p>															
$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$																
$P(x_i = I   y = \text{HAM}) = \frac{5}{15}$																
$P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$																
$P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$																
$P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$																
$P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$																
$P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$																
$P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$																
$P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$																
$P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$																
$P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$																
$P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$																
$P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$																
$P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$																
$P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$																

$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$															
	<p>category/class = <math>h(\text{buy cheap rolex replica rolex})</math></p>															
$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$ - feature vectors (in <b>bold</b> )   $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$ - labels																

$$\begin{aligned}
 P(y = \text{HAM}) &= \frac{5}{7} & P(y = \text{SPAM}) &= \frac{2}{7} \\
 P(x_i = I | y = \text{HAM}) &= \frac{5}{15} \\
 P(x_i = \text{rolex} | y = \text{HAM}) &= \frac{1}{15} \\
 P(x_i = \text{own} | y = \text{HAM}) &= \frac{4}{15} \\
 P(x_i = \text{replica} | y = \text{HAM}) &= \frac{1}{15} \\
 P(x_i = \text{watch} | y = \text{HAM}) &= \frac{3}{15} \\
 P(x_i = \text{buy} | y = \text{HAM}) &= \frac{0}{15} \\
 P(x_i = \text{cheap} | y = \text{HAM}) &= \frac{1}{15} \\
 P(x_i = I | y = \text{SPAM}) &= \frac{0}{7} \\
 P(x_i = \text{rolex} | y = \text{SPAM}) &= \frac{2}{7} \\
 P(x_i = \text{own} | y = \text{SPAM}) &= \frac{0}{7} \\
 P(x_i = \text{replica} | y = \text{SPAM}) &= \frac{2}{7} \\
 P(x_i = \text{watch} | y = \text{SPAM}) &= \frac{0}{7} \\
 P(x_i = \text{buy} | y = \text{SPAM}) &= \frac{1}{7} \\
 P(x_i = \text{cheap} | y = \text{SPAM}) &= \frac{2}{7}
 \end{aligned}$$

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> buy cheap rolex replica rolex	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> I own cheap rolex watch	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> replica	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$															
	category/class = $h(\quad 0 \quad 1 \quad 0 \quad 1 \quad 0 \quad 1 \quad 1 \quad )$	$P(y = \text{HAM}) = \frac{5}{7} \quad P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing																
	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_8) \propto P(y = \text{HAM}) * \prod_{i=1}^5 P(x_i   y = \text{HAM})$	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_8) \propto P(y = \text{HAM}) * \prod_{i=1}^5 P(x_i   y = \text{HAM}) =$ $P(y = \text{HAM}) * P(x_1 = \text{buy}   y = \text{HAM}) * P(x_2 = \text{cheap}   y = \text{HAM}) * P(x_3 = \text{rolex}   y = \text{HAM}) * P(x_4 = \text{replica}   y = \text{HAM}) * P(x_5 = \text{rolex}   y = \text{HAM})$															
	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_1 = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_1 = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$															

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing																
	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_8) \propto P(y = \text{HAM}) * \prod_{i=1}^5 P(x_i   y = \text{HAM}) =$ $\frac{5}{7} * \frac{0}{15} * \frac{1}{15} * \frac{1}{15} * \frac{1}{15} * \frac{1}{15} = 0$	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{SPAM}   x_8) \propto P(y = \text{SPAM}) * \prod_{i=1}^5 P(x_i   y = \text{SPAM}) =$ $P(y = \text{SPAM}) * P(x_1 = \text{buy}   y = \text{SPAM}) * P(x_2 = \text{cheap}   y = \text{SPAM}) * P(x_3 = \text{rolex}   y = \text{SPAM}) * P(x_4 = \text{replica}   y = \text{SPAM}) * P(x_5 = \text{rolex}   y = \text{SPAM})$															
	$P(x_i = I   y = \text{HAM}) = \frac{5}{7}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$															

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

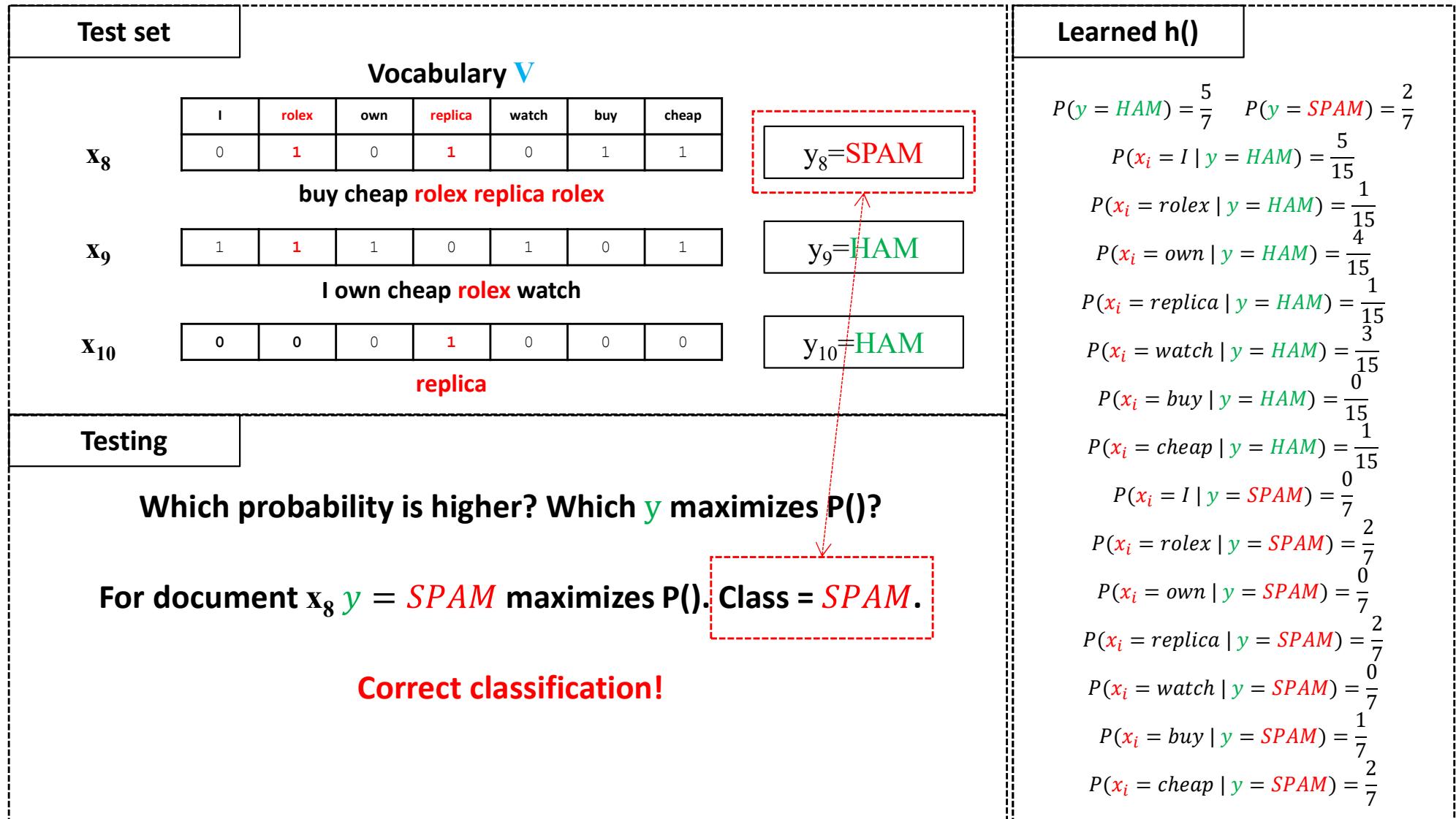
Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> buy cheap rolex replica rolex	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> I own cheap rolex watch	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> replica	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{SPAM}   x_8) \propto P(y = \text{SPAM}) * \prod_{i=1}^5 P(x_i   y = \text{SPAM}) =$ $\frac{2}{7} * \frac{1}{7} * \frac{2}{7} * \frac{2}{7} * \frac{2}{7} * \frac{2}{7} \approx 0.00027$															
	$P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$															

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_8) = 0$ $P(y = \text{SPAM}   x_8) \approx 0.00027$															
	For document $x_8$ $y = \text{SPAM}$ maximizes $P()$ . Class = $\text{SPAM}$ .															
$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$ - feature vectors (in <b>bold</b> )   $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$ - labels		$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

# Spam Detection: Testing Classifier



$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$ <p>category/class = <math>h(x_9)</math></p>															
		$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$															
	<p>category/class = <math>h(\text{I own cheap rolex watch})</math></p>															
$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$ - feature vectors (in <b>bold</b> )   $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$ - labels																

$$\begin{aligned}
 P(y = \text{HAM}) &= \frac{5}{7} & P(y = \text{SPAM}) &= \frac{2}{7} \\
 P(x_i = I | y = \text{HAM}) &= \frac{5}{15} \\
 P(x_i = \text{rolex} | y = \text{HAM}) &= \frac{1}{15} \\
 P(x_i = \text{own} | y = \text{HAM}) &= \frac{4}{15} \\
 P(x_i = \text{replica} | y = \text{HAM}) &= \frac{1}{15} \\
 P(x_i = \text{watch} | y = \text{HAM}) &= \frac{3}{15} \\
 P(x_i = \text{buy} | y = \text{HAM}) &= \frac{0}{15} \\
 P(x_i = \text{cheap} | y = \text{HAM}) &= \frac{1}{15} \\
 P(x_i = I | y = \text{SPAM}) &= \frac{0}{7} \\
 P(x_i = \text{rolex} | y = \text{SPAM}) &= \frac{2}{7} \\
 P(x_i = \text{own} | y = \text{SPAM}) &= \frac{0}{7} \\
 P(x_i = \text{replica} | y = \text{SPAM}) &= \frac{2}{7} \\
 P(x_i = \text{watch} | y = \text{SPAM}) &= \frac{0}{7} \\
 P(x_i = \text{buy} | y = \text{SPAM}) &= \frac{1}{7} \\
 P(x_i = \text{cheap} | y = \text{SPAM}) &= \frac{2}{7}
 \end{aligned}$$

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$																
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> buy cheap rolex replica rolex	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$		
I	rolex	own	replica	watch	buy	cheap												
0	1	0	1	0	1	1												
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>1</th><th>1</th><th>0</th><th>1</th><th>0</th><th>1</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> I own cheap rolex watch	I	1	1	0	1	0	1	1	1	1	0	1	0	1	$y_9 = \text{HAM}$		
I	1	1	0	1	0	1												
1	1	1	0	1	0	1												
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>0</th><th>0</th><th>0</th><th>1</th><th>0</th><th>0</th><th>0</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> replica	I	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	$y_{10} = \text{HAM}$
I	0	0	0	1	0	0	0											
0	0	0	1	0	0	0	0											
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$																	
	category/class = $h(\quad 1 \quad \textcolor{red}{1} \quad 1 \quad 0 \quad 1 \quad 0 \quad 1 \quad )$																	
$P(y = \text{HAM}) = \frac{5}{7} \quad P(y = \text{SPAM}) = \frac{2}{7}$																		
$P(x_i = I   y = \text{HAM}) = \frac{5}{15}$																		
$P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$																		
$P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$																		
$P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$																		
$P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$																		
$P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$																		
$P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$																		
$P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$																		
$P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$																		
$P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$																		
$P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$																		
$P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$																		
$P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$																		
$P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$																		

$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $\mathbf{V}$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_9) \propto P(y = \text{HAM}) * \prod_{i=1}^7 P(x_i   y = \text{HAM})$ $P(y = \text{SPAM}   x_9) \propto P(y = \text{SPAM}) * \prod_{i=1}^7 P(x_i   y = \text{SPAM})$	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing																
	Which probability is higher? Which $y$ maximizes $P()$ ?	$P(y = \text{HAM}   x_9) \propto P(y = \text{HAM}) * \prod_{i=1}^5 P(x_i   y = \text{HAM}) =$ $P(y = \text{HAM}) * P(x_1 = I   y = \text{HAM}) * P(x_2 = \text{own}   y = \text{HAM}) * P(x_3 = \text{cheap}   y = \text{HAM}) * P(x_4 = \text{rolex}   y = \text{HAM}) * P(x_5 = \text{watch}   y = \text{HAM})$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_9) \propto P(y = \text{HAM}) * \prod_{i=1}^5 P(x_i   y = \text{HAM}) =$ $\frac{5}{7} * \frac{5}{15} * \frac{4}{15} * \frac{1}{15} * \frac{1}{15} * \frac{3}{15} \approx 0.000056$															
	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$															

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{SPAM}   x_9) \propto P(y = \text{SPAM}) * \prod_{i=1}^5 P(x_i   y = \text{SPAM}) =$ $P(y = \text{SPAM}) * P(x_1 = I   y = \text{SPAM}) * P(x_2 = \text{own}   y = \text{SPAM}) * P(x_3 = \text{cheap}   y = \text{SPAM}) * P(x_4 = \text{rolex}   y = \text{SPAM}) * P(x_5 = \text{watch}   y = \text{SPAM})$															
	$P(y = \text{HAM}) = \frac{5}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$															

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{SPAM}   x_9) \propto P(y = \text{SPAM}) * \prod_{i=1}^5 P(x_i   y = \text{SPAM}) =$ $\frac{2}{7} * \frac{0}{7} * \frac{0}{7} * \frac{2}{7} * \frac{2}{7} * \frac{0}{7} = 0$	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_9) \approx 0.000056$ $P(y = \text{SPAM}   x_9) = 0$															
	For document $x_8$ $y = \text{HAM}$ maximizes $P()$ . Class = $\text{HAM}$ .															
$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$																

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	<p>Which probability is higher? Which <math>y</math> maximizes <math>P()</math>?</p> <p>For document <math>x_9</math>, <math>y = \text{HAM}</math> maximizes <math>P()</math>. Class = <math>\text{HAM}</math>.</p> <p>Correct classification!</p>															
		$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing		$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$ <p>category/class = <math>h(x_{10})</math></p>														
		$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $\mathbf{V}$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing		$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$ <p>category/class = <math>h(\text{replica})</math></p>														
		$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> buy cheap rolex replica rolex	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> I own cheap rolex watch	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> replica	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( P(y) * \prod_{i=1}^N P(x_i   y) \right)$															
	category/class = $h(\quad 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0 \quad )$	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_{N-2}, \mathbf{x}_{N-1}, \mathbf{x}_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_{10}) \propto P(y = \text{HAM}) * \prod_{i=1}^1 P(x_i   y = \text{HAM})$ $P(y = \text{SPAM}   x_{10}) \propto P(y = \text{SPAM}) * \prod_{i=1}^1 P(x_i   y = \text{SPAM})$															
	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$															

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_{10}) \propto P(y = \text{HAM}) * \prod_{i=1}^1 P(x_i   y = \text{HAM}) =$ $P(y = \text{HAM}) * P(x_1 = \text{replica}   y = \text{HAM}) = \frac{5}{7} * \frac{1}{15} \approx 0.048$	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing																
	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{SPAM}   x_{10}) \propto P(y = \text{SPAM}) * \prod_{i=1}^1 P(x_i   y = \text{SPAM}) =$ $P(y = \text{SPAM}) * P(x_1 = \text{replica}   y = \text{SPAM}) = \frac{2}{7} * \frac{2}{7} \approx 0.082$	$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $\mathbf{V}$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	Which probability is higher? Which $y$ maximizes $P()$ ?															
	$P(y = \text{HAM}   x_{10}) \approx 0.048$ $P(y = \text{SPAM}   x_{10}) \approx 0.082$															
	For document $x_{10}$ $y = \text{SPAM}$ maximizes $P()$ . Class = $\text{SPAM}$ .															
$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$																

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Spam Detection: Testing Classifier

Test set	Vocabulary $V$	Learned $h()$														
$x_8$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td></tr> </tbody> </table> <p>buy cheap rolex replica rolex</p>	I	rolex	own	replica	watch	buy	cheap	0	1	0	1	0	1	1	$y_8 = \text{SPAM}$
I	rolex	own	replica	watch	buy	cheap										
0	1	0	1	0	1	1										
$x_9$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>1</td><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>1</td></tr> </tbody> </table> <p>I own cheap rolex watch</p>	I	rolex	own	replica	watch	buy	cheap	1	1	1	0	1	0	1	$y_9 = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
1	1	1	0	1	0	1										
$x_{10}$	<table border="1"> <thead> <tr> <th>I</th><th>rolex</th><th>own</th><th>replica</th><th>watch</th><th>buy</th><th>cheap</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> </tbody> </table> <p>replica</p>	I	rolex	own	replica	watch	buy	cheap	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$
I	rolex	own	replica	watch	buy	cheap										
0	0	0	1	0	0	0										
Testing	<p>Which probability is higher? Which <math>y</math> maximizes <math>P()</math>?</p> <p>For document <math>x_{10}</math> <math>y = \text{SPAM}</math> maximizes <math>P()</math>. Class = <math>\text{SPAM}</math>.</p> <p><u>Incorrect classification! Misclassification.</u></p>															
		$P(y = \text{HAM}) = \frac{5}{7}$ $P(y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = I   y = \text{HAM}) = \frac{5}{15}$ $P(x_i = \text{rolex}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{own}   y = \text{HAM}) = \frac{4}{15}$ $P(x_i = \text{replica}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = \text{watch}   y = \text{HAM}) = \frac{3}{15}$ $P(x_i = \text{buy}   y = \text{HAM}) = \frac{0}{15}$ $P(x_i = \text{cheap}   y = \text{HAM}) = \frac{1}{15}$ $P(x_i = I   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{rolex}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{own}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{replica}   y = \text{SPAM}) = \frac{2}{7}$ $P(x_i = \text{watch}   y = \text{SPAM}) = \frac{0}{7}$ $P(x_i = \text{buy}   y = \text{SPAM}) = \frac{1}{7}$ $P(x_i = \text{cheap}   y = \text{SPAM}) = \frac{2}{7}$														

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Classifier Evaluation: Confusion Matrix

		Predicted class		
		Positive	Negative	
Actual class	Positive	True Positive (TP)	False Negative (FN) Type II Error	Sensitivity (Recall) $\frac{TP}{TP+FN}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	Specificity $\frac{TN}{TN+FP}$
		Precision $\frac{TP}{TP+FP}$	Negative Predictive Value $\frac{TN}{TN+FN}$	Accuracy $\frac{TP+TN}{TP+TN+FP+FN}$

# Classifier Evaluation: Confusion Matrix

		Predicted class		
		SPAM	HAM	
Actual class	SPAM	True Positive (TP)	False Negative (FN) Type II Error	Sensitivity (Recall) $\frac{TP}{TP+FN}$
	HAM	False Positive (FP) Type I Error	True Negative (TN)	Specificity $\frac{TN}{TN+FP}$
	Precision $\frac{TP}{TP+FP}$	Negative Predictive Value $\frac{TN}{TN+FN}$	Accuracy $\frac{TP+TN}{TP+TN+FP+FN}$	

# Spam Detection: Evaluating Classifier

Test set	Vocabulary $V$							Testing results					
$x_8$	1	rolex	own	replica	watch	buy	cheap	$y_8 = \text{SPAM}$					
	0	1	0	1	0	1	1						
	buy cheap rolex replica rolex												
$x_9$	1	1	1	0	1	0	1	$y_9 = \text{HAM}$					
	I own cheap rolex watch												
$x_{10}$	0	0	0	1	0	0	0	$y_{10} = \text{HAM}$					
	replica												
Evaluation								Confusion matrix					
	$y_8 = \text{SPAM}$	$y_8 = \text{SPAM}$	true positive										
	$y_9 = \text{HAM}$	$y_9 = \text{HAM}$	true negative										
	$y_{10} = \text{HAM}$	$y_{10} = \text{SPAM}$	false positive										
No false negatives in this example.													
$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$ - feature vectors (in <b>bold</b> )   $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$ - labels													

		SPAM	HAM
SPAM	True Positive (TP) Type II Error		
HAM	False Positive (FP) Type I Error		True Negative (TN)

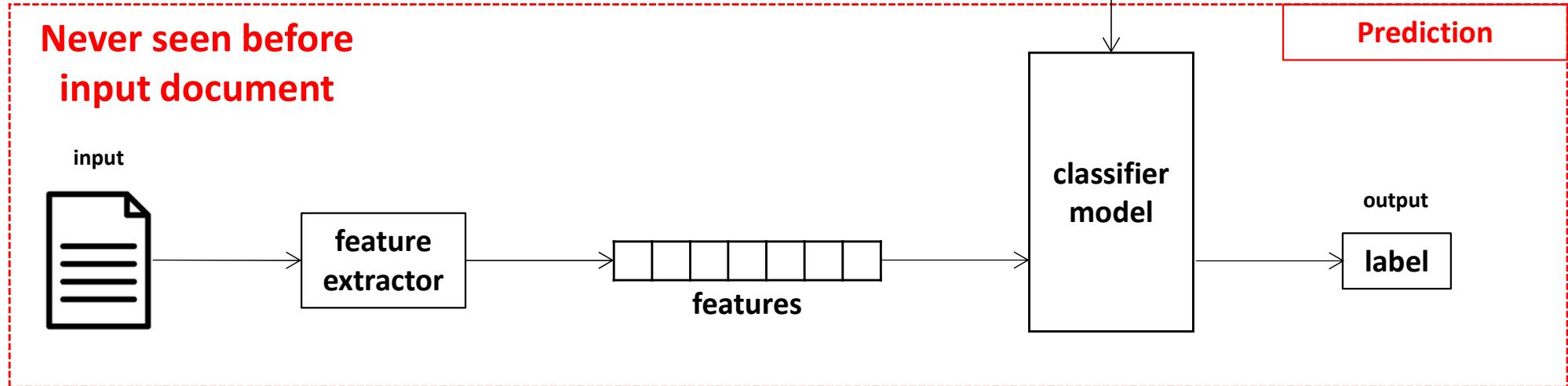
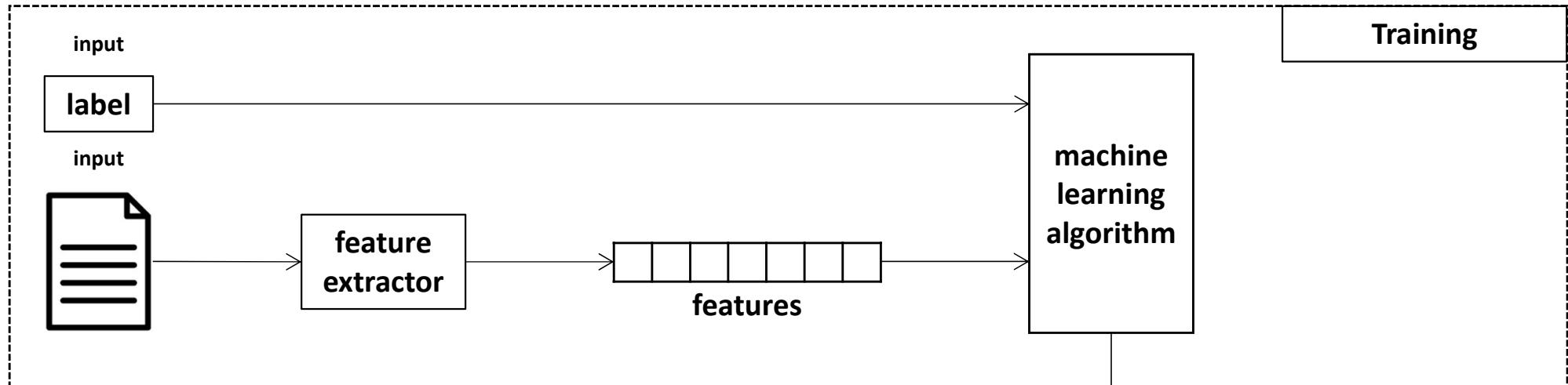
# Classifier Evaluation: Confusion Matrix

		Predicted class		
		SPAM	HAM	
Actual class	SPAM	TP = 1 ( $x_8$ )	FN = 0	Sensitivity (Recall) $\frac{1}{1+0} = 1.0$
	HAM	FP = 1 ( $x_{10}$ )	TN = 1 ( $x_9$ )	Specificity $\frac{1}{1+1} = 0.5$
		Precision $\frac{1}{1+1} = 0.5$	Negative Predictive Value $\frac{1}{1+0} = 1.0$	Accuracy $\frac{1+1}{1+1+1+0} = \frac{2}{3}$

# Confusion Matrix Explained

- Accuracy  $(TP+TN)/(TP+TN+FP+FN)$ :
  - Overall, how often is the classifier correct?
- Misclassification rate [Error Rate]  $(FP+FN)/(TP+TN+FP+FN)$ :
  - Overall, how often is the classifier incorrect?
- Sensitivity [Recall | True Positive Rate]  $(TP)/(TP+FN)$ :
  - When it's actually yes, how often does it predict yes?
- Specificity [True Negative Rate]  $(TN)/(TN+FP)$ 
  - When it's actually no, how often does it predict no?
- Precision  $(TP)/(TP+FP)$ 
  - When it predicts yes, how often is it correct?
- Negative Predictive Value  $(TN)/(TN+FN)$ 
  - When it predicts no, how often is it correct?

# Supervised Learning with ML



# Spam Detection: Prediction

Unseen  $x$

Vocabulary  $V$

I	rolex	own	replica	watch	buy	cheap
$x_?$	0	1	0	0	0	1

buy rolex

Label needs to be decided

$y_? = ????$

Prediction

Which probability is higher? Which  $y$  maximizes  $P()$ ?

$$P(y = \text{HAM} | x_?) \propto P(y = \text{HAM}) * \prod_{i=1}^2 P(x_i | y = \text{HAM}) =$$

$$P(y = \text{HAM}) * P(x_1 = \text{buy} | y = \text{HAM}) * P(x_2 = \text{rolex} | y = \text{HAM})$$

$$= \frac{5}{7} * \frac{0}{15} * \frac{1}{15} = 0$$

$$P(y = \text{SPAM} | x_?) \propto P(y = \text{SPAM}) * \prod_{i=1}^2 P(x_i | y = \text{SPAM}) =$$

$$P(y = \text{SPAM}) * P(x_1 = \text{buy} | y = \text{SPAM}) * P(x_2 = \text{rolex} | y = \text{SPAM})$$

$$= \frac{2}{7} * \frac{1}{7} * \frac{2}{7} \approx 0.012$$

For document  $x_?$ ,  $y = \text{SPAM}$  maximizes  $P()$ . Class =  $\text{SPAM}$

Learned  $h()$

$$P(y = \text{HAM}) = \frac{5}{7} \quad P(y = \text{SPAM}) = \frac{2}{7}$$

$$P(x_1 = I | y = \text{HAM}) = \frac{5}{15}$$

$$P(x_1 = \text{rolex} | y = \text{HAM}) = \frac{1}{15}$$

$$P(x_1 = \text{own} | y = \text{HAM}) = \frac{4}{15}$$

$$P(x_1 = \text{replica} | y = \text{HAM}) = \frac{1}{15}$$

$$P(x_1 = \text{watch} | y = \text{HAM}) = \frac{3}{15}$$

$$P(x_1 = \text{buy} | y = \text{HAM}) = \frac{0}{15}$$

$$P(x_1 = \text{cheap} | y = \text{HAM}) = \frac{1}{15}$$

$$P(x_1 = I | y = \text{SPAM}) = \frac{0}{7}$$

$$P(x_1 = \text{rolex} | y = \text{SPAM}) = \frac{2}{7}$$

$$P(x_1 = \text{own} | y = \text{SPAM}) = \frac{0}{7}$$

$$P(x_1 = \text{replica} | y = \text{SPAM}) = \frac{2}{7}$$

$$P(x_1 = \text{watch} | y = \text{SPAM}) = \frac{0}{7}$$

$$P(x_1 = \text{buy} | y = \text{SPAM}) = \frac{1}{7}$$

$$P(x_1 = \text{cheap} | y = \text{SPAM}) = \frac{2}{7}$$

$x_1, x_2, x_3, \dots, x_{N-2}, x_{N-1}, x_N$  - feature vectors (in **bold**) |  $y_1, y_2, y_3, \dots, y_{N-2}, y_{N-1}, y_N$  - labels

# Classifier Problems: Zero Counts

$$P(\textcolor{red}{x}_i = \text{word} \mid \textcolor{green}{y} = \textcolor{blue}{CLASS}) = \frac{\text{count}(\textcolor{red}{x}_i = \text{word}, \textcolor{green}{y} = \textcolor{blue}{CLASS})}{\sum_{\textcolor{brown}{x} \in V} \text{count}(\textcolor{brown}{x}, \textcolor{green}{y} = \textcolor{blue}{CLASS})}$$

- Unseen words:
  - $\text{count}(\textcolor{red}{x}_i = \text{word}, \textcolor{green}{y} = \textcolor{blue}{CLASS})$  can be zero
- Words NOT present in samples for one class (see: example):
  - $\text{count}(\textcolor{red}{x}_i = \text{word}, \textcolor{green}{y} = \textcolor{blue}{CLASS})$  can be zero
- Solution: smoothing (e.g. Laplace smoothing)

$$P(\textcolor{red}{x}_i = \text{word} \mid \textcolor{green}{y} = \textcolor{blue}{CLASS}) = \frac{\text{count}(\textcolor{red}{x}_i = \text{word}, \textcolor{green}{y} = \textcolor{blue}{CLASS}) + \alpha}{\sum_{\textcolor{brown}{x} \in V} \text{count}(\textcolor{brown}{x}, \textcolor{green}{y} = \textcolor{blue}{CLASS}) + \alpha * |V|}$$

where:  $\alpha$  - pseudo-occurrence (typically “add 1”),  $|V|$  - vocabulary size

# Classifier Problems: Underflow

$$P(y | x) \propto P(y) * \prod_{i=1}^N P(x_i | y)$$

- N can be large (100 and more):
  - long, “wordy”, documents
- some  $P(x_i | y)$  can be very small ( $< 0.1$ )
  - the product  $\prod_{i=1}^N P(x_i | y)$  may lead to underflow
- Solution: use logarithms

$$\log(P(y | x)) \propto \log(P(y)) + \sum_{i=1}^N \log(P(x_i | y))$$

# Naive Bayes Classifier

category/class =  $\text{h}(\text{document})$

Finding model / hypothesis  $\text{h} \rightarrow$  Finding probabilities for  $y_{MAP}$

$$y_{MAP} \propto \underset{y \in Y}{argmax} \left( \log(P(y)) + \sum_{i=1}^N \log(P(x_i | y)) \right)$$

MAP: Maximum a posteriori (corresponds to the most likely class).

# Naive Bayes Classifier

category/class =  $\text{h}(\text{document})$

Finding model / hypothesis  $\text{h} \rightarrow$  Finding probabilities for  $y_{MAP}$

$$y_{MAP} \propto \underset{y \in Y}{\operatorname{argmax}} \left( \log(P(y)) + \sum_{i=1}^N \log(P(x_i | y)) \right)$$

- **Taking log doesn't change the ranking of classes!**
  - The class with highest probability also has highest log probability!
- **It's a linear model:**
  - Just a max of a sum of weights: a linear function of the inputs
  - So Naive Bayes is a linear classifier

# Naive Bayes: Training/Testing

**function** TRAIN NAIVE BAYES(D, C) **returns**  $\log P(c)$  and  $\log P(w|c)$

**for each** class  $c \in C$  # Calculate  $P(c)$  terms

$N_{doc}$  = number of documents in D

$N_c$  = number of documents from D in class c

$logprior[c] \leftarrow \log \frac{N_c}{N_{doc}}$

$V \leftarrow$  vocabulary of D

$bigdoc[c] \leftarrow \text{append}(d)$  **for**  $d \in D$  **with** class c

**for each** word  $w$  in V # Calculate  $P(w|c)$  terms

$count(w,c) \leftarrow$  # of occurrences of w in  $bigdoc[c]$

$loglikelihood[w,c] \leftarrow \log \frac{count(w,c) + 1}{\sum_{w' \text{ in } V} (count(w',c) + 1)}$

**return**  $logprior, loglikelihood, V$

**function** TEST NAIVE BAYES( $testdoc, logprior, loglikelihood, C, V$ ) **returns** best  $c$

**for each** class  $c \in C$

$sum[c] \leftarrow logprior[c]$

**for each** position  $i$  in  $testdoc$

$word \leftarrow testdoc[i]$

**if**  $word \in V$

$sum[c] \leftarrow sum[c] + loglikelihood[word,c]$

**return**  $\text{argmax}_c sum[c]$

# Naive Bayes: Summary

- Pros:

- Very fast and easy-to-implement
- Well-understood formally & experimentally
  - see “Naive (Bayes) at Forty”, Lewis, ECML98

- Cons:

- Seldom gives the very best performance (baseline)
- “Probabilities”  $P(y | x)$  are not accurate
- Probabilities tend to be close to zero or one

# Naive Bayes: Stop Words

- Some systems ignore stop words
  - Stop words: very frequent words like the and a.
    - Sort the vocabulary by word frequency in training set
    - Call the top 10 or 50 words the stopword list.
    - Remove all stop words from both training and test sets
      - As if they were never there!
- But removing stop words doesn't usually help
  - So **in practice** most NB algorithms use all words and don't use stopword lists

# Naive Bayes: Unknown Words

- What about unknown words
  - that appear in our test data
  - but not in our training data or vocabulary?
- We **ignore** them
  - Remove them from the test document!
  - Pretend they weren't there!
  - Don't include any probability for them at all!
- Why don't we build an unknown word model?
  - It doesn't help: knowing which class has more unknown words is not generally helpful!

# Naive Bayes: More Than Two Classes

- Dealing with **any-of** or **multivalue** classification
  - A document can belong to 0, 1, or more than 1 classes.
- For each class  $c \in C$ 
  - Build a classifier  $h_c$  to distinguish  $c$  from all other classes  $c' \in C$
- Given test document  $d$ ,
  - Evaluate it for membership in each class using each  $h_c$
  - $d$  belongs to **any** class for which  $h_c$  returns true

# Naive Bayes: More Than Two Classes

- Dealing with **one-of** or **multinomial** classification
  - Classes are mutually exclusive: each document in exactly one class
- For each class  $c \in C$ 
  - Build a classifier  $h_c$  to distinguish  $c$  from all other classes  $c' \in C$
- Given test document  $d$ ,
  - Evaluate it for membership in each class using each  $h_c$
  - $d$  belongs to the **one** class with maximum score