

Text Classification and Naïve Bayes

Text Classification

- **Text Classification (Text Categorization)** is *the task of assigning a label or categorization category to an entire text or document.*
- Some of common text categorization tasks are:
 - Sentiment analysis
 - extraction of **sentiment**, positive or negative orientation that writer expresses toward an object.
 - Spam detection
 - binary classification task of assigning an email to one of the two classes spam or not-spam.
 - Authorship identification
 - determining a text's author.
 - Age/gender identification
 - determining a text's author characteristics like gender and age.
 - Language Identification
 - finding the language of a text.

Text Classification: Definition

- Text classification can be defined as follows:

Input:

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

Output:

- a predicted class $c \in C$

Classification Methods:

Hand-Coded Rules

- The goal of classification is to take a single observation, extract some useful features, and thereby classify the observation into one of a set of discrete classes.
- One method for classifying text is to use **hand-written rules**.
- Rules based on combinations of words or other features
 - spam: black-list-address OR (“dollars” AND “have been selected”)
- Accuracy can be high if rules carefully are refined by experts.
- But building and maintaining these **hand-written rules** can be expensive.
 - Rules can be fragile.
 - It may require domain knowledge.

Classification Methods:

Supervised Machine Learning

- Most cases of *text classification* in language processing are done via **supervised machine learning** methods.
- The goal of a **supervised machine learning algorithm** is to learn how to map from a new observation to a correct output.

Input:

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

Output:

- a learned classifier **model**: $d \rightarrow c$

Classification Methods:

Supervised Machine Learning

- Our goal is to learn a classifier that is capable of mapping from a *new document* d to its *correct class* $c \in C$.
- A **probabilistic classifier** additionally will tell us the probability of the observation being in the class.
- **Generative classifiers** like **Naive Bayes** build a model of how a class could generate some input data.
 - Given an observation, they return the class most likely to have generated the observation.
- **Discriminative classifiers** like **logistic regression** instead learn what features from the input are most useful to discriminate between the different possible classes.
- Some classifiers:
 - Naïve Bayes
 - Logistic regression
 - Support-vector machines
 - k-Nearest Neighbors, ...

Naïve Bayes Classifier

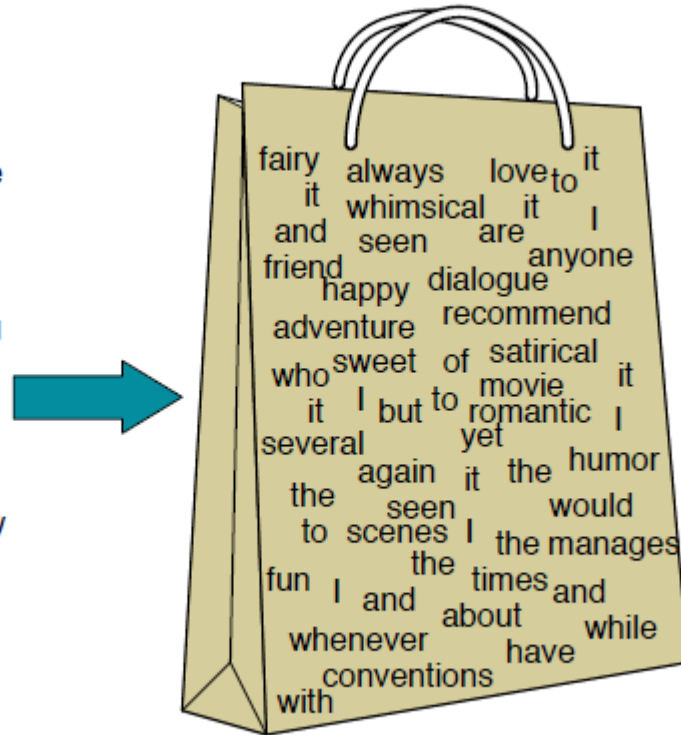
Naïve Bayes Classifier

- A **Naïve Bayes Classifier** is a *Bayesian classifier* (based on *Bayes rule*) that makes a simplifying (*naive*) assumption about how the features interact.
- A text document is represented as a **bag of words**.
 - i.e. a text document is represented as an unordered set of words with their position ignored, keeping only their frequency in the document.

Bag of Words Representation

Bag of Words: *The position of the words is ignored (the **bag of words assumption**) and we make use of the frequency of each word.*

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!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

Bag of Words Representation

- A Naive Bayes classifier may use **all words in the text as features**.

classifier (`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.`) = **class**

Bag of Words Representation

- Or, a Naive Bayes classifier may use a subset of words in the text as **features**.

classifier (

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.

) = **class**

Naïve Bayes Classifier

- Naive Bayes is a **probabilistic classifier**, meaning that for a document **d**, out of all classes **c** ∈ **C** the classifier returns the class **ĉ** which has the **maximum posterior probability** given the document.
 - the hat notation **ĉ** to mean “**our estimate of the correct class**”.

- For a document **d** and a class **c**

$$P(\mathbf{c}|\mathbf{d}) = \frac{P(\mathbf{d}|\mathbf{c}) P(\mathbf{c})}{P(\mathbf{d})}$$

- Most likely class (MAP: maximum a posteriori):

$$\hat{\mathbf{c}} = \mathbf{c}_{\text{MAP}} = \underset{\mathbf{c} \in \mathbf{C}}{\text{argmax}} P(\mathbf{c}|\mathbf{d})$$

Naïve Bayes Classifier

$$\hat{c} = c_{\text{MAP}} = \underset{c \in C}{\operatorname{argmax}} P(c|d)$$

By Bayes Rule:

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} \frac{P(d|c) P(c)}{P(d)}$$

Dropping the denominator (it is same for all classes):

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(d|c) P(c)$$

**likelihood of
the document**

**prior probability
of the class**

Naïve Bayes Classifier

$$\hat{c} = \operatorname{argmax}_{c \in C} P(d|c) P(c)$$

- we can represent a document d as a set of features f_1, f_2, \dots, f_n :

$$\hat{c} = \operatorname{argmax}_{c \in C} P(f_1, f_2, \dots, f_n | c) P(c)$$

- Unfortunately, it is hard to compute $P(f_1, f_2, \dots, f_n | c)$ directly:
 - estimating probability of every possible combination of features would require huge numbers of parameters and impossibly large training sets.
- Naive Bayes classifiers make **two simplifying assumptions**.
 - **Bag of Words Assumption**: Assume position doesn't matter.
 - **Naive Bayes Assumption**: *conditional independence* assumption that the probabilities $P(f_i | c)$ are *independent* given the class c and hence they can be multiplied as follows:

$$P(f_1, \dots, f_n | c) = P(f_1 | c) * \dots * P(f_n | c)$$

Naïve Bayes Classifier

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(f_1, f_2, \dots, f_n | c) P(c)$$

- With **simplifying assumptions** (**bag of words assumption** and **Naive bayes assumption**) the class chosen by a **naive Bayes classifier** is:

$$c_{NB} = \underset{c \in C}{\operatorname{argmax}} P(c) \prod_{f \in F} P(f|c)$$

Applying Naive Bayes Classifier to Text Classification

- To apply the naive Bayes classifier to text, we need to consider word positions, by simply walking an index through every word position in the document:

positions ← all word positions in test document

$$c_{\text{NB}} = \operatorname{argmax}_{c \in \mathcal{C}} P(c) \prod_{i \in \text{positions}} P(w_i | c)$$

- Naive Bayes calculations can be done in log space, in order to avoid underflow and increase speed.

$$c_{\text{NB}} = \operatorname{argmax}_{c \in \mathcal{C}} \log(P(c)) + \sum_{i \in \text{positions}} \log(P(w_i | c))$$

Naive Bayes Classifier: Learning

$$\mathbf{c}_{\text{NB}} = \underset{\mathbf{c} \in \mathbf{C}}{\operatorname{argmax}} \mathbf{P}(\mathbf{c}) \prod_{i \in \text{positions}} \mathbf{P}(w_i | \mathbf{c})$$

- How can we learn the probabilities $\mathbf{P}(\mathbf{c})$ and $\mathbf{P}(w_i | \mathbf{c})$?
- We learn these probabilities from our **training set: $(d_1, c_1), \dots, (d_m, c_m)$**
- $\mathbf{P}(\mathbf{c})$ can be computed (estimated) as follows:

$$\mathbf{P}(\mathbf{c}) = \frac{N_{\mathbf{c}}}{N_{\text{doc}}}$$

- $N_{\mathbf{c}}$ is the **number of documents in our training data with class \mathbf{c}** and N_{doc} is the **total number of documents**.

Naive Bayes Classifier: Learning

- In order to compute (estimate) $P(\mathbf{w}_i|\mathbf{c})$:
 - First, concatenate all documents with category \mathbf{c} into one big “category \mathbf{c} ” text.
 - Then, use the frequency of \mathbf{w}_i in this concatenated document to give a maximum likelihood estimate of the probability.

$$P(\mathbf{w}_i|\mathbf{c}) = \frac{\text{count}(\mathbf{w}_i,\mathbf{c})}{\sum_{\mathbf{w} \in V} \text{count}(\mathbf{w},\mathbf{c})}$$

- The vocabulary V consists of the union of all the word types in all classes, not just the words in one class \mathbf{c} .

Naive Bayes Classifier: Learning

Example: We can compute $P(\text{fantastic}|\text{positive})$ as follows:

- The probability of **fantastic** word given category **positive**.

$$P(\text{fantastic}|\text{positive}) = \frac{\text{count}(\text{fantastic},\text{positive})}{\sum_{w \in V} \text{count}(w,\text{positive})}$$

Naive Bayes Classifier: Learning

Smoothing

- In order to get rid of zero probability values for $\mathbf{P}(\mathbf{w}_i|\mathbf{c})$, we should use a **smoothing method**.
- The simplest solution is the **add-one (Laplace) smoothing**:

$$\mathbf{P}(\mathbf{w}_i|\mathbf{c}) = \frac{\text{count}(\mathbf{w}_i,\mathbf{c})+1}{\sum_{\mathbf{w} \in \mathbf{V}}(\text{count}(\mathbf{w},\mathbf{c})+1)} = \frac{\text{count}(\mathbf{w}_i,\mathbf{c})+1}{(\sum_{\mathbf{w} \in \mathbf{V}} \text{count}(\mathbf{w},\mathbf{c})) + |\mathbf{V}|}$$

Naive Bayes Classifier: Learning *unknown words*

- What do we do about words that occur in our test data but are not in our vocabulary at all because they did not occur in any training document in any class?
- The solution for such **unknown words** is to *ignore them* — *remove them from the test document and not include any probability for them at all.*
- Some systems choose to completely ignore another class of words:
 - **stop words**, very frequent words like the and a.

Naive Bayes Algorithm: Train

using add-1 smoothing

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

$bigdoc[c] \leftarrow$ 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 $logprior, loglikelihood, V$

Naive Bayes Algorithm: Test

using add-1 smoothing

```
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  $\operatorname{argmax}_c$  sum[c]
```

Naive Bayes Classifier

sentiment analysis example

- A *sentiment analysis domain* with the two classes **positive (+)** and **negative (-)**, and the following miniature **training and test documents** simplified from actual movie reviews.

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
	+	the most fun film of the summer
Test	?	predictable with no fun

- The prior $\mathbf{P(c)}$ for the two classes: $\mathbf{P(c)} = \frac{N_c}{N_{\text{doc}}}$
 $P(-) = \frac{3}{5}$ $P(+) = \frac{2}{5}$

Naive Bayes Classifier

sentiment analysis example

- The word **with** doesn't occur in the training set, so we drop it completely.
- The **likelihoods (conditional probabilities)** from the training set for the remaining three words “**predictable**”, “**no**”, and “**fun**” are:

$$P(\text{“predictable”}|-) = \frac{1+1}{14+20} \quad P(\text{“predictable”}|+) = \frac{0+1}{9+20}$$

$$P(\text{“no”}|-) = \frac{1+1}{14+20} \quad P(\text{“no”}|+) = \frac{0+1}{9+20}$$

$$P(\text{“fun”}|-) = \frac{0+1}{14+20} \quad P(\text{“fun”}|+) = \frac{1+1}{9+20}$$

$$P(w_i|c) = \frac{\text{count}(w_i,c)+1}{(\sum_{w \in V} \text{count}(w,c))+|V|} \quad |V| = 20$$

$$\sum_{w \in V} \text{count}(w, -) = 14$$

$$\sum_{w \in V} \text{count}(w, +) = 9$$

	Cat	Documents
Training	-	just plain boring
	-	entirely predictable and lacks energy
	-	no surprises and very few laughs
	+	very powerful
	+	the most fun film of the summer
Test	?	predictable with no fun

Naive Bayes Classifier

sentiment analysis example

- For the test sentence $S = \text{“predictable with no fun”}$, after removing the word ‘**with**’, the **chosen class** is **negative (-)**:

$$P(-)P(S|-) = \frac{3}{5} \times \frac{2 \times 2 \times 1}{34^3} = 6.1 \times 10^{-5}$$

$$P(+)P(S|+) = \frac{2}{5} \times \frac{1 \times 1 \times 2}{29^3} = 3.2 \times 10^{-5}$$

Naïve Bayes and Language Modeling

- When we use **all** of the words in the text (not a subset), **Naïve Bayes** has an important similarity to **language modeling**.
- Each class = a **unigram language model**
- Assigning each sentence: $P(s|c) = \prod P(\text{word}|c)$

<u>$P(w \text{pos})$</u>	<u>w</u>					
0.1	I	<u>I</u>	<u>love</u>	<u>this</u>	<u>fun</u>	<u>film</u>
0.1	love	0.1	0.1	0.01	0.05	0.1
0.01	this					
0.05	fun					
0.1	film					
...						

$$P(s|\text{pos}) = 0.0000005$$

Naïve Bayes as a Language Model

- Which class assigns the higher probability to s?

Model pos	
$P(w pos)$	<u>w</u>
0.1	<u>I</u>
0.1	love
0.01	this
0.05	fun
0.1	film

Model neg	
$P(w neg)$	<u>w</u>
0.2	<u>I</u>
0.001	love
0.01	this
0.005	fun
0.1	film

<u>I</u>	<u>love</u>	<u>this</u>	<u>fun</u>	<u>film</u>
0.1	0.1	0.01	0.05	0.1
0.2	0.001	0.01	0.005	0.1

$$P(s|pos) = 0.0000005$$

$$P(s|neg) = 0.000000001$$

$$P(s|pos) > P(s|neg)$$

Text Classification: Evaluation

Contingency Table, Precision, Recall, F-measure

- In order to evaluate *how good is our classifier*, we can use different **evaluation metrics**.
- In evaluation, we compare the **test set results** of our classifier with **gold labels** (*the human labels for the test set documents*).
- As a result of this comparison, first we build a **contingency table** before calculate our **evaluation metrics**.

Contingency Table

		<i>gold standard labels</i>		
		gold positive	gold negative	
<i>system output labels</i>	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

Text Classification: Evaluation

Accuracy, Precision, Recall

- **Accuracy** is *percentage of all the observations our system labeled correctly*.
 - Although accuracy might seem a natural metric, accuracy doesn't work well when the classes are unbalanced.

$$\text{Accuracy} = \frac{\text{tp} + \text{tn}}{\text{tp} + \text{fp} + \text{tn} + \text{fn}}$$

- **Precision** measures *percentage of items that system detected that are in fact positive*.

$$\text{Precision} = \frac{\text{tp}}{\text{tp} + \text{fp}}$$

- **Recall** measures *percentage of items actually present in the input that were correctly identified by the system*.

$$\text{Recall} = \frac{\text{tp}}{\text{tp} + \text{fn}}$$

Text Classification: Evaluation

A combined measure: F-measure

- **Precision** and **Recall**, unlike Accuracy, emphasize *true positives*.
 - Looking only one of them can be misleading.
 - tp=1 fp=0 fn=99 → Precision = %100 (while Recall=%1)
 - tp=1 fp=99 fn=0 → Recall= %100 (while Precision=%1)
- **F-measure** is a single metric that incorporates aspects of both **precision** and **recall**.

$$F_{\beta} = \frac{(\beta^2 + 1) * P * R}{\beta^2 * P + R}$$

- The β parameter differentially weights the importance of recall and precision.
 - values of $\beta > 1$ favor recall, while values of $\beta < 1$ favor precision.
- The most frequently used metric, and is called $F_{\beta=1}$ or just F_1 :

$$F_1 = \frac{2 * P * R}{P + R}$$

Text Classification *with more than two classes*

- Many of classification tasks in language processing have more than two classes.
 - For sentiment analysis we generally have 3 classes (positive, negative, neutral).
- There are two kinds of multi-class classification tasks:
 - **any-of** (or **multi-label**) **classification**,
 - **one-of** (or **multinomial**) **classification**

any-of (or **multi-label**) **classification**:

- *Each document can be assigned more than one label.*
- We can solve **any-of classification** by building separate binary classifiers for each class c , trained on positive examples labeled c and negative examples not labeled c .
- Given a test document d , then each classifier makes their decision independently, and we may assign multiple labels to d .

Text Classification *with more than two classes*

one-of (or **multinomial**) **classification**:

- *classes are mutually exclusive and each document appears in exactly one class.*
- We again build a separate binary classifier trained on positive examples from c and negative examples from all other classes.
- Given a test document d , we run all the classifiers and *choose the label from the classifier with the highest score.*

Text Classification: Confusion matrix

with more than two classes

- Confusion matrix for a three-class categorization task,
 - for each pair (c_1, c_2) , how many documents from c_1 were (in)correctly assigned to c_2

		<i>gold labels</i>			
		urgent	normal	spam	
<i>system output</i>	urgent	8	10	1	$\text{precision}_u = \frac{8}{8+10+1}$
	normal	5	60	50	$\text{precision}_n = \frac{60}{5+60+50}$
	spam	3	30	200	$\text{precision}_s = \frac{200}{3+30+200}$
		$\text{recall}_u = \frac{8}{8+5+3}$	$\text{recall}_n = \frac{60}{10+60+30}$	$\text{recall}_s = \frac{200}{1+50+200}$	

Text Classification: class evaluation measures

with more than two classes

- **Accuracy:** Fraction of documents classified correctly (= 1 - error rate).

$$\text{Accuracy} = \frac{\sum_i c_{ii}}{\sum_j \sum_i c_{ij}}$$

- **Precision:** Fraction of documents assigned class **i** that are actually about class **i**:

$$\text{Precision} = \frac{\sum_i c_{ii}}{\sum_j c_{ji}}$$

- **Recall:** Fraction of documents in class **i** classified correctly:

$$\text{Recall} = \frac{\sum_i c_{ii}}{\sum_j c_{ij}}$$

Microaveraging and Macroaveraging

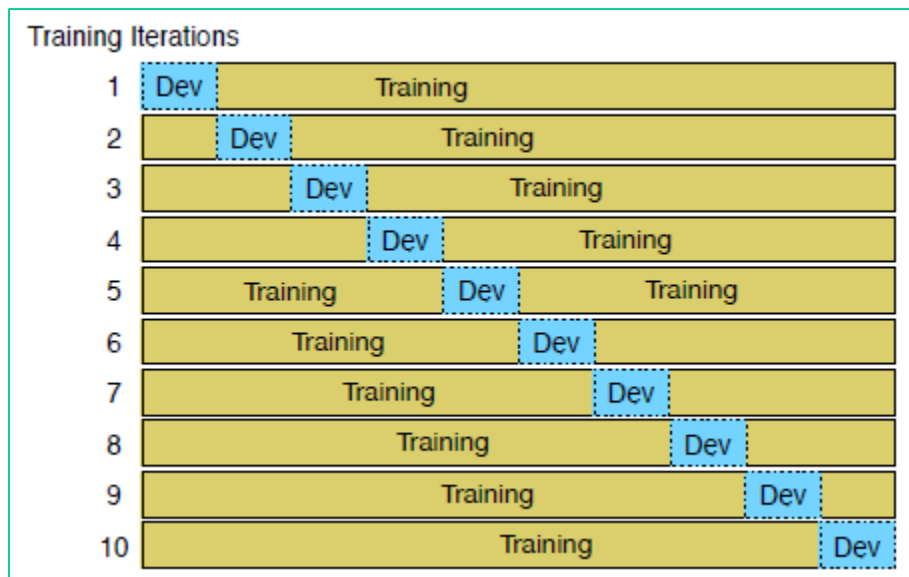
- In order to derive a single metric that tells us how well the system is doing, we can combine these values in two ways.
 - In **macroaveraging**, compute performance for each class, and then average over classes.
 - In **microaveraging**, collect decisions for all classes into a single contingency table, and then compute precision and recall from that table.

	Class 1: Urgent	Class 2: Normal	Class 3: Spam	Pooled																																				
	<table border="1"> <thead> <tr> <th></th> <th>true urgent</th> <th>true not</th> </tr> </thead> <tbody> <tr> <th>system urgent</th> <td>8</td> <td>11</td> </tr> <tr> <th>system not</th> <td>8</td> <td>340</td> </tr> </tbody> </table>		true urgent	true not	system urgent	8	11	system not	8	340	<table border="1"> <thead> <tr> <th></th> <th>true normal</th> <th>true not</th> </tr> </thead> <tbody> <tr> <th>system normal</th> <td>60</td> <td>55</td> </tr> <tr> <th>system not</th> <td>40</td> <td>212</td> </tr> </tbody> </table>		true normal	true not	system normal	60	55	system not	40	212	<table border="1"> <thead> <tr> <th></th> <th>true spam</th> <th>true not</th> </tr> </thead> <tbody> <tr> <th>system spam</th> <td>200</td> <td>33</td> </tr> <tr> <th>system not</th> <td>51</td> <td>83</td> </tr> </tbody> </table>		true spam	true not	system spam	200	33	system not	51	83	<table border="1"> <thead> <tr> <th></th> <th>true yes</th> <th>true no</th> </tr> </thead> <tbody> <tr> <th>system yes</th> <td>268</td> <td>99</td> </tr> <tr> <th>system no</th> <td>99</td> <td>635</td> </tr> </tbody> </table>		true yes	true no	system yes	268	99	system no	99	635
	true urgent	true not																																						
system urgent	8	11																																						
system not	8	340																																						
	true normal	true not																																						
system normal	60	55																																						
system not	40	212																																						
	true spam	true not																																						
system spam	200	33																																						
system not	51	83																																						
	true yes	true no																																						
system yes	268	99																																						
system no	99	635																																						
	$\text{precision} = \frac{8}{8+11} = .42$	$\text{precision} = \frac{60}{60+55} = .52$	$\text{precision} = \frac{200}{200+33} = .86$	$\text{microaverage precision} = \frac{268}{268+99} = .73$																																				
	$\text{macroaverage precision} = \frac{.42+.52+.86}{3} = .60$																																							

Cross-validation

cross-validation:

- we randomly choose a training and test set division of our data, train our classifier, and then compute the error rate on the test set.
- Then we repeat with a different randomly selected training set and test set.
- We do sampling process 10 times and average these 10 runs to get an average error rate.
- This is called **10-fold cross-validation**.



Text Classification: Summary

- Text categorization assigns an entire text to a class from a finite set,
- Some text categorization tasks are sentiment analysis, spam detection, language identification, and authorship attribution.
- Naive Bayes is a generative model that make the bag of words assumption (position doesn't matter) and the conditional independence assumption (words are conditionally independent of each other given the class).
- Classifiers are evaluated based on precision and recall.
- Classifiers are trained using distinct training, dev, and test sets, including the use of cross-validation in the training set.