

INF4080 – 2020 FALL

NATURAL LANGUAGE PROCESSING

Jan Tore Lønning

(Mostly Text) Classification, Naive Bayes

Lecture 3, 31 Aug

Today - Classification

3

- Motivation
- Classification
- Naive Bayes classification
- NB for text classification
 - ▣ The multinomial model
 - ▣ The Bernoulli model
- Experiments: training, test and cross-validation
- Evaluation

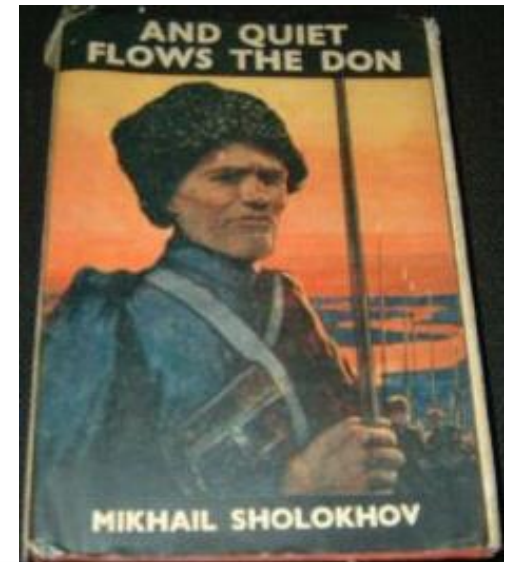
4

Motivation

Did Mikhail Sholokov write *And Quiet Flows the Don*?

5

- Sholokov, 1905-1984
- *And Quiet Flows the Don*
 - ▣ published 1928-1940
- Nobel prize, literature, 1965
- Authorship contested
 - ▣ e.g. Aleksandr Solzhenitsyn, 1974
- Geir Kjetsaa (UiO) et al, 1984
 - ▣ refuted the contestants
- Nils Lid Hjort, 2007, confirmed Kjetsaa by using sentence length and advanced statistics.
- https://en.wikipedia.org/wiki/Mikhail_Sholokhov



Kjetsaa according to Hjort
In addition to various linguistic analyses and several doses of detective work, quantitative data were gathered and organised, for example, relating to word lengths, frequencies of certain words and phrases, sentence lengths, grammatical characteristics, etc.

Positive or negative movie review?

6



□ *unbelievably disappointing*



□ *Full of zany characters and richly applied satire, and come great plot twists*



□ *this is the greatest screwball comedy ever filmed*

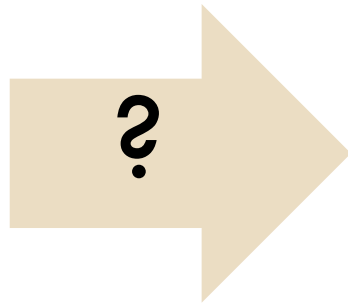


□ *It was pathetic. The worst part about it was the boxing scenes.*

From Jurafsky & Martin

What is the subject of this article?

MEDLINE Article



MeSH Subject Category Hierarchy

- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

8

Classification

Classification

9

- Can be rule-based, but mostly machine learned
- Text classification is a sub-class

- Text classification examples:
 - Spam detection
 - Genre classification
 - Language identification
 - Sentiment analysis:
 - Positive-negative

- Other types of classification:
 - Word sense disambiguation
 - Sentence splitting
 - Tagging
 - Named-entity recognition

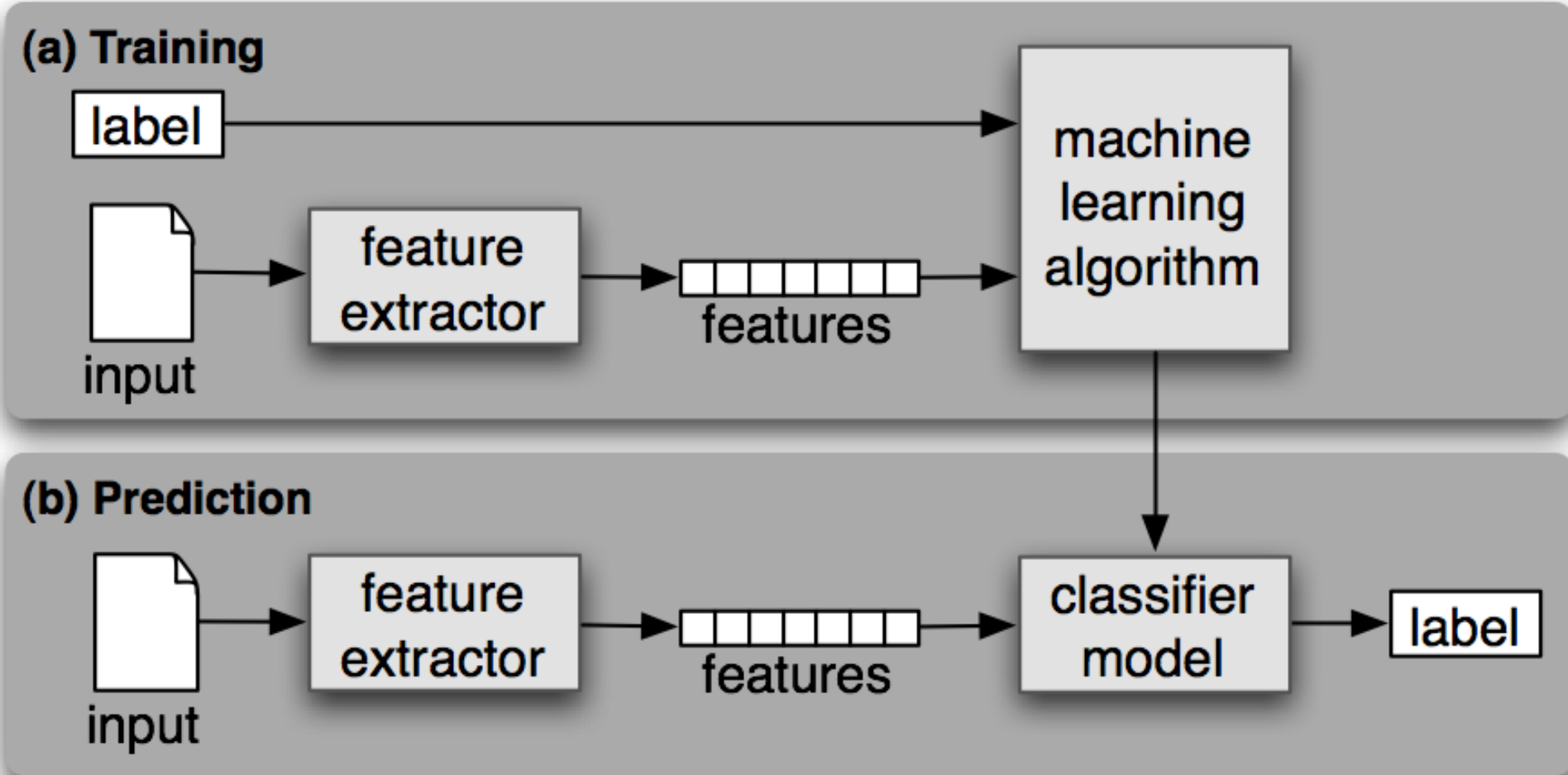
Machine learning

10

1. Supervised
 1. Classification (categorical)
 2. Regression (numerical)
2. Unsupervised
3. Semi-supervised
4. Reinforcement learning

- Supervised:
 - ▣ Given classes
 - ▣ Given examples of correct classes
- Unsupervised:
 - ▣ Construct classes

Supervised classification



Supervised classification

12

- Given
 - ▣ a well-defined set of observations, O
 - ▣ a given set of classes, $C = \{c_1, c_2, \dots, c_k\}$
- Goal: a classifier, γ , a mapping from O to C
- For supervised training one needs a set of pairs from $O \times C$

Task	O	C
Spam classification	E-mails	Spam, no-spam
Language identification	Pieces of text	Arabian, Chinese, English, Norwegian, ...
Word sense disambiguation	Occurrences of "bass"	Sense 1, ..., sense 8

Features

13

- To represent the objects in O , extract a set of features
- Be explicit:
 - ▣ Which features
 - ▣ For each feature
 - The type
 - Categorical
 - Numeric (Discrete/Continuous)
 - The value space

O: person
Features:

- height
- weight
- hair color
- eye color
- ...

O: email
Features:

- length
- sender
- contained words
- language
- ...

Cf. First lecture
Classes and features are both attributes of the observations

Supervised classification

- A given set of classes, $C = \{c_1, c_2, \dots, c_k\}$
 - A well defined class of observations, O
-
- Some features f_1, f_2, \dots, f_n
 - For each feature: a set of possible values V_1, V_2, \dots, V_n
 - The set of feature vectors: $V = V_1 \times V_2 \times \dots \times V_n$
 - Each observation in O is represented by some member of V :
 - Written $(f_1=v_1, f_2=v_2, \dots, f_n=v_n)$, or
 - (v_1, v_2, \dots, v_n) , if we have decided the order
 - A classifier, γ , can be considered a mapping from V to C

A variety of ML classifiers

15

- k-Nearest Neighbors
- Rocchio
- Naive Bayes
- Logistic regression (Maximum entropy)
- Support Vector Machines
- Decision Trees
- Perceptron
- Multi-layered neural nets ("Deep learning")

16

Naïve Bayes

Example: Jan. 2021

Professor, do
you think I will
enjoy
IN3050?



I can give you a
scientific answer
using machine
learning.



Baseline

- Survey
- Asked all the students of 2020
- 200 answered:
 - ▣ 130 yes
 - ▣ 70 no
- Baseline classifier:
 - ▣ Choose the majority class
 - ▣ Accuracy $0.65=65\%$
 - ▣ (With two classes, always ≥ 0.5)

Yes,
you will like it.



Example: one year from now, Jan. 2021

Professor, do
you think I will
enjoy
IN3050?



To answer that, I
have to ask you
some questions.



The 2020 survey (imaginary)

Ask each of the 200 students:

- Did you enjoy the course?
 - ▣ Yes/no
- Do you like mathematics?
 - ▣ Yes/no
- Do you have programming experience?
 - ▣ None/some/good (= 3 or more courses)
- Have you taken advanced machine learning courses?
 - ▣ Yes/no
- And many more questions, but we have to simplify here

Results of the 2020 survey: a data set

Student no	Enjoy maths	Programming	Adv. ML	Enjoy
1	Y	Good	N	Y
2	Y	Some	N	Y
3	N	Good	Y	N
4	N	None	N	N
5	N	Good	N	Y
6	N	Good	Y	Y
....				

Summary of the 2020 survey

	A	B	C	D	E	F
1	programing	AdvML-course	Like maths	enjoyed	not enjoye	sum
2	good	yes	yes	3	10	13
3	good	yes	no	7	4	11
4	good	no	yes	50	4	54
5	good	no	no	40	4	44
6	some	yes	yes	4	1	5
7	some	yes	no	0	0	0
8	some	no	yes	11	9	20
9	some	no	no	10	24	34
10	none	yes	yes	1	2	3
11	none	yes	no	0	0	0
12	none	no	yes	2	5	7
13	none	no	no	2	7	9
14				130	70	200

Our new student

	A	B	C	D	E	F
1	programing	AdvML-course	Like maths	enjoyed	not enjoye	sum
2	good	yes	yes	3	10	13
3	good	yes	no	7	4	11
4	good	no	yes	50	4	54
5	good	no	no	40	4	44
6	some	yes	yes	4	1	5
7	some	yes	no	0	0	0
8	some	no	yes	11	9	20
9	some	no	no	10	24	34
10	none	yes	yes	1	2	3
11	none	yes	no	0	0	0
12	none	no	yes	2	5	7
13	none	no	no	2	7	9
14				130	70	200

But what should we say to a student with some programming background, and adv. ML course who does not like maths.?

- We ask our incoming new student the same three question
- From the table we can see e.g. that if:
 - ▣ she has good programming
 - ▣ no AdvML-course
 - ▣ does not like maths
- There is a $\frac{40}{44}$ chance she will enjoy the course

A little more formal

- What we do is that we consider
 - ▣ $P(\text{enjoy} = \text{yes} \mid \text{prog} = \text{good}, \text{AdvML} = \text{no}, \text{Maths} = \text{no})$ and
 - ▣ $P(\text{enjoy} = \text{not} \mid \text{prog} = \text{good}, \text{AdvML} = \text{no}, \text{Maths} = \text{no})$
- and decide on the class which has the largest probability, in symbols
 - ▣ $\operatorname{argmax}_{y \in \{\text{yes}, \text{no}\}} P(\text{enjoy} = y \mid \text{prog} = \text{good}, \text{AdvML} = \text{no}, \text{Maths} = \text{no})$

- But, there may be many more features.
 - ▣ An exponential growth in possible combinations
 - ▣ We might not have seen all combinations, or they may be rare
- Therefore we apply Bayes theorem, and we make a simplifying assumption

Naive Bayes: Decision

25

- Given an observation

- $\langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle$

- Consider

- $P(s_m | \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle)$ for each class s_m

- Choose the class with the largest value, in symbols

$$\arg \max_{s_m \in S} P(s_m | \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle)$$

- i.e. choose the class for which the observation is most likely

Naive Bayes: Model

26

□ Bayes formula

$$\square P(s_m | \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle) = \frac{P(\langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle | s_m) P(s_m)}{P(\langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle)}$$

□ Sparse data, we may not even have seen

$$\square \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle$$

□ We assume (wrongly) independence

$$\square P(\langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle | s_m) \approx \prod_{i=1}^n P(f_i = v_i | s_m)$$

□ Putting together, choose

$$\square \arg \max_{s_m \in S} P(s_m | \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle) \approx \arg \max_{s_m \in S} P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m)$$

Naive Bayes, Training 1

27

- Maximum Likelihood

- $$\hat{P}(s_m) = \frac{C(s_m, o)}{C(o)}$$

- ▣ where $C(s_m, o)$ are the number of occurrences of observations o in class s_m

- Observe what we are doing:

- ▣ We are looking for the true probability $P(s_m)$

- ▣ $\hat{P}(s_m)$ is an approximation to this, our best guess from a set of observations

- ▣ Maximum likelihood means that it is the model which makes the set of observations we have seen, most likely

Naive Bayes: Training 2

28

□ Maximum Likelihood

$$\square \hat{P}(f_i = v_i | s_m) = \frac{C(f_i = v_i, s_m)}{C(s_m)}$$

□ where $C(f_i = v_i, s_m)$ is the number of observations o

■ where the observation o belongs to class s_m

■ and the feature f_i takes the value v_i

□ $C(s_m)$ is the number of observations belonging to class s_m

Back to example

29

16	programming	enjoyed	num		$\hat{P}(x \mid \text{yes})$
17	good	yes	100	/130=	0,7692308
18	some	yes	25	/130=	0,1923077
19	none	yes	5	/130=	0,0384615
20					$\hat{P}(x \mid \text{no})$
21	good	no	22	/70=	0,3142857
22	some	no	34	/70=	0,4857143
23	none	no	14	/70=	0,2
24					
25	advanced	enjoyed	num		$\hat{P}(x \mid \text{yes})$
26	yes	yes	15	/130=	0,1153846
27	no	yes	115	/130=	0,8846154
28					$\hat{P}(x \mid \text{no})$
29	yes	no	17	/70=	0,2428571
30	no	no	53	/70=	0,7571429
31					
32	like maths	enjoyed	num		$\hat{P}(x \mid \text{yes})$
33	yes	yes	71	/130=	0,5461538
34	no	yes	59	/130=	0,4538462
35					$\hat{P}(x \mid \text{no})$
36	yes	no	31	/70=	0,4428571
37	no	no	39	/70=	0,5571429



- Collect the numbers
- Estimate the probabilities

	A	B	C	D	E	F
1	programing	AdvML-course	Like maths	enjoyed	not enjoye	sum
2	good	yes	yes	3	10	13
3	good	yes	no	7	4	11
4	good	no	yes	50	4	54
5	good	no	no	40	4	44
6	some	yes	yes	4	1	5
7	some	yes	no	0	0	0
8	some	no	yes	11	9	20
9	some	no	no	10	24	34
10	none	yes	yes	1	2	3
11	none	yes	no	0	0	0
12	none	no	yes	2	5	7
13	none	no	no	2	7	9
14				130	70	200

Back to example

30

			num		$\hat{P}(x \text{yes})$
16	programming	enjoyed			
17	good	yes	100	/130=	0,7692308
18	some	yes	25	/130=	0,1923077
19	none	yes	5	/130=	0,0384615
20					$\hat{P}(x \text{no})$
21	good	no	22	/70=	0,3142857
22	some	no	34	/70=	0,4857143
23	none	no	14	/70=	0,2
24					
25	advanced	enjoyed	num		$\hat{P}(x \text{yes})$
26	yes	yes	15	/130=	0,1153846
27	no	yes	115	/130=	0,8846154
28					$\hat{P}(x \text{no})$
29	yes	no	17	/70=	0,2428571
30	no	no	53	/70=	0,7571429
31					
32	like maths	enjoyed	num		$\hat{P}(x \text{yes})$
33	yes	yes	71	/130=	0,5461538
34	no	yes	59	/130=	0,4538462
35					$\hat{P}(x \text{no})$
36	yes	no	31	/70=	0,4428571
37	no	no	39	/70=	0,5571429

- $\operatorname{argmax}_{c_m \in C} P(c_m) \prod_{i=1}^n P(f_i = v_i | c_m)$
- $P(\text{yes}) \times P(\text{good}|\text{yes}) \times P(A:\text{no}|\text{yes}) \times P(M:\text{no}|\text{yes}) = \frac{130}{200} \times \frac{100}{130} \times \frac{115}{130} \times \frac{59}{130} = 0.2$
- $P(\text{no}) \times P(\text{good}|\text{no}) \times P(A:\text{no}|\text{no}) \times P(M:\text{no}|\text{no}) = \frac{70}{200} \times \frac{22}{70} \times \frac{53}{70} \times \frac{39}{70} = 0.046$
- So we predict that the student will most probably enjoy the class
- Accuracy on training data: 75%
 - ▣ Compare to Baseline: 65%
 - ▣ Best classifier: 80%

Laplace-smoothing

31



- MLE-estimate: $P(w_i) = \frac{c_i}{N}$
- Laplace-estimate: $P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$
- Lidstone-smoothing: add k , e.g. 0.5: $\hat{P}(w_i) = \frac{c_i + k}{N + kV}$
- `nltk.NaiveBayesClassifier` uses Lidstone (0.5) as default

Laplace applied to example

$$P_{\text{Laplace}}(w_i) = \frac{c_i + 1}{N + V}$$

32

16	programming	enjoyed	num		$\hat{P}(x \text{yes})$
17	good	yes	100	/130=	0,7692308
18	some	yes	25	/130=	0,1923077
19	none	yes	5	/130=	0,0384615
20					$\hat{P}(x \text{no})$
21	good	no	22	/70=	0,3142857
22	some	no	34	/70=	0,4857143
23	none	no	14	/70=	0,2
24					
25	advanced	enjoyed	num		$\hat{P}(x \text{yes})$
26	yes	yes	15	/130=	0,1153846
27	no	yes	115	/130=	0,8846154
28					$\hat{P}(x \text{no})$
29	yes	no	17	/70=	0,2428571
30	no	no	53	/70=	0,7571429
31					
32	like maths	enjoyed	num		$\hat{P}(x \text{yes})$
33	yes	yes	71	/130=	0,5461538
34	no	yes	59	/130=	0,4538462
35					$\hat{P}(x \text{no})$
36	yes	no	31	/70=	0,4428571
37	no	no	39	/70=	0,5571429

- $\hat{P}(\text{prog} = \text{good} | \text{yes}) = \frac{100+1}{130+3}$
- $\hat{P}(\text{prog} = \text{some} | \text{yes}) = \frac{25+1}{130+3}$
- $\hat{P}(\text{prog} = \text{none} | \text{yes}) = \frac{5+1}{130+3}$

- $\hat{P}(\text{adv} = \text{yes} | \text{yes}) = \frac{15+1}{130+2}$

Naive Bayes: Calculation

33

$$\square \arg \max_{s_m \in \mathcal{S}} P(s_m | \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle) \approx \arg \max_{s_m \in \mathcal{S}} P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m)$$

□ For calculations

□ avoid underflow, use logarithms

$$\square \arg \max_{s_m \in \mathcal{S}} P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m) =$$

$$\arg \max_{s_m \in \mathcal{S}} \left(\log \left(P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m) \right) \right)$$

$$= \arg \max_{s_m \in \mathcal{S}} \left(\log(P(s_m)) + \sum_{i=1}^n \log(P(f_i = v_i | s_m)) \right)$$

Properties of Naive Bayes

34

- A probabilistic classifier
- A multi-class classifier:
 - ▣ i.e. can handle more than two classes
- Categorical features natively
 - ▣ Can be adopted to numeric features
- NLTK contains an implementation
- The independence assumption is in general: **wrong!**
 - ▣ $P(v_1, v_2, \dots, v_n|c)$ is far from
 - ▣ $P(v_1|c) \times P(v_2|c) \cdots \times P(v_n|c)$
- Still NB works reasonably well as a classifier (discriminator)
- It is not prone to overfitting
- Other classifiers may work better

35

Text classification with NB

Text classification with NB

- Naive Bayes may be applied to various NLP tasks
- Text classification:
 - ▣ Goal: classify the text on the basis of the words in the text
 - ▣ What are the features?
 - ▣ What are the possible values.
- Two possible answers:
 - ▣ The Multinomial model
 - ▣ The Bernoulli model

1. Multinomial NB text classification

37

$$\arg \max_{s_m \in \mathcal{S}} P(s_m | \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle) \approx \arg \max_{s_m \in \mathcal{S}} P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m)$$

- f_i refers to position i in the text
- v_i is the word occurring in this position
- n is the number of tokens in the text

- Simplifying assumption: a word is equally likely in all positions
- Hence we count how many times each word occurs in the text

$$\arg \max_{s_m \in \mathcal{S}} P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m) = \arg \max_{s_m \in \mathcal{S}} P(s_m) \prod_{i=1}^n P(v_i | s_m)$$

Multinomial NB: Training

38

- $\hat{P}(s_m) = \frac{C(s_m, o)}{C(o)}$
 - ▣ where $C(s_m, o)$ is the number of occurrences of observations o in class s_m
- $\hat{P}(w_i | s_m) = \frac{C(w_i, s_m)}{\sum_j C(w_j, s_m)}$
 - ▣ where $C(w_i, s_m)$ is the number of occurrences of word w_i in all texts in class s_m
 - ▣ $\sum_j C(w_j, s_m)$ is the total number of words in all texts in class s_m

Example: Movie reviews corpus (NLTK)

39

- 2000 documents
 - ▣ (a subset of a larger corpus)
- Two classes: 'neg', 'pos', 1000 doc.s in each class

```
> from nltk.corpus import movie_reviews  
  
> documents = [(list(movie_reviews.words(fileid)), category)  
                for category in movie_reviews.categories()  
                for fileid in movie_reviews.fileids(category)]
```

Example: movie reviews, multinomial

40

- ▣ Considered 1900 doc.s for training
- ▣ 'pitt' occurs in 15 'pos' and 6 'neg' reviews
- ▣ 'pitt' occurs 31 times in the 'pos' reviews and 25 times in the negative reviews
- ▣ There are 798,742 words in the 'pos' reviews and 705,726 in the 'neg' reviews
- ▣ $\hat{P}(w = pitt|pos) = \frac{31}{798\,742}$ $\hat{P}(w = pitt|neg) = \frac{25}{705\,726}$

Example: more features

41

□ In [63]: pos_docs['pitt']

□ Out[63]: 15

□ In [64]: neg_docs['pitt']

□ Out[64]: 6

□ In [65]: neg_docs['spacey']

□ Out[65]: 4

□ In [66]: pos_docs['spacey']

□ Out[66]: 17

□ In [71]: pos_docs['terrible']

□ Out[71]: 26

□ In [72]: neg_docs['terrible']

□ Out[72]: 85

□ In [73]: neg_docs['terrific']

□ Out[73]: 19

□ In [74]: pos_docs['terrific']

□ Out[74]: 75

$3 \times \text{pitt},$

$2 \times \text{terrible},$

$0 \times \text{terrific}$

42

'pos'

- $\hat{P}(\text{pos}) = \frac{959}{1900}$
- $\hat{P}(w = \text{pitt}|\text{pos}) = \frac{31}{798\ 742}$
- $\hat{P}(w = \text{terrible}|\text{pos}) = \frac{26}{798\ 742}$
- $\hat{P}(\text{pos} | 3 \times \text{pitt}, 2 \times \text{terrible}, 0 \times \text{terrific}) =$
 $k' \frac{959}{1900} \times \left(\frac{31}{798\ 742}\right)^3 \times \left(\frac{26}{798\ 742}\right)^2$
 $= k' 3.12 \times 10^{-23}$

'neg'

- $\hat{P}(\text{neg}) = \frac{941}{1900}$
- $\hat{P}(w = \text{pitt}|\text{neg}) = \frac{25}{705\ 726}$
- $\hat{P}(w = \text{terrible}|\text{neg}) = \frac{104}{705\ 726}$
- $\hat{P}(\text{pos} | 3 \times \text{pitt}, 2 \times \text{terrible}, 0 \times \text{terrific}) =$
 $k' \frac{941}{1900} \times \left(\frac{25}{705\ 726}\right)^3 \times \left(\frac{104}{705\ 726}\right)^2$
 $= k' 4.78 \times 10^{-22}$

2. NB – Bernoulli model for text classification

43

- How are words turned into features?
- A vocabulary of words, W
- Each word w_i makes a feature f_i
- The possible values for f_i is True and False (1 and 0)
- $f_i = 1$ in a document if and only if it contains w_i .

Bernoulli NB: Decision

$$\arg \max_{s_m \in S} P(s_m | \langle f_1 = v_1, f_2 = v_2, \dots, f_n = v_n \rangle) \approx \arg \max_{s_m \in S} P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m)$$

- f_i refers to a word in the vocabulary
- v_i is 1 or 0 depending on whether the word occurs in the text or not
- n is the number of words in the vocabulary

$$\arg \max_{s_m \in S} P(s_m) \prod_{i=1}^n P(f_i = v_i | s_m) = \arg \max_{s_m \in S} P(s_m) \prod_{i=1}^n P(v_i | s_m)$$

Example: movie reviews NLTK (Bernoulli)

45

□ 'pitt' occurs in 15 'pos' and 6 'neg' reviews

$$\square \hat{P}(pitt = True|pos) = \frac{15}{959} \quad \hat{P}(pitt = True|neg) = \frac{6}{941}$$

$$\square \hat{P}(pitt = False|pos) = \frac{944}{959} \quad \hat{P}(pitt = False|neg) = \frac{935}{941}$$

pitt = True, terrible = True, terrific = False

46

'pos'

- $\hat{P}(pos) = \frac{959}{1900}$
- $\hat{P}(pitt = True|pos) = \frac{15}{959}$
- $\hat{P}(terrible = True|pos) = \frac{26}{959}$
- $\hat{P}(terrific = False|pos) = \frac{959-75}{959}$
- $\hat{P}(pos | pitt = True, terrible = True, terrific = False) = k \frac{959}{1900} \times \frac{15}{959} \times \frac{26}{959} \times \frac{884}{959} = k0.00020$

'neg'

- $\hat{P}(neg) = \frac{941}{1900}$
- $\hat{P}(pitt = True|neg) = \frac{6}{941}$
- $\hat{P}(terrible = True|neg) = \frac{85}{941}$
- $\hat{P}(terrific = False|neg) = \frac{941-19}{941}$
- $\hat{P}(neg | pitt = True, terrible = True, terrific = False) = k \frac{941}{1900} \times \frac{6}{941} \times \frac{85}{941} \times \frac{922}{941} = k0.00028$

$(k = 1/P(pitt = True, terrible = True, terrific = False))$, the same for both classes)

The two models

47

- Multinomial model
 - ▣ Jurafsky and Martin, 3.ed, sec. 4, Sentiment analysis
 - ▣ Related to n-gram models
- Bernoulli
 - ▣ NLTK book, Sec. 6.1, 6.2, 6.5
 - Including the Movie review example
 - ▣ Jurafsky and Martin, 2.ed, sec. 20.2, WSD
- Both
 - ▣ Manning, Raghavan, Schütze, *Introduction to Information Retrieval*, Sec. 13.0-13.3

Comparison

48

Multinomial

- Counts how many times a term is present
- Considers
 - ▣ only the present terms
 - ▣ ignores absent terms
- Tends to be the better of the two for longer texts

Bernoulli

- Registers whether a term is present or not
- Considers both
 - ▣ The present terms
 - ▣ The absent terms
- Compatible on shorter snippets

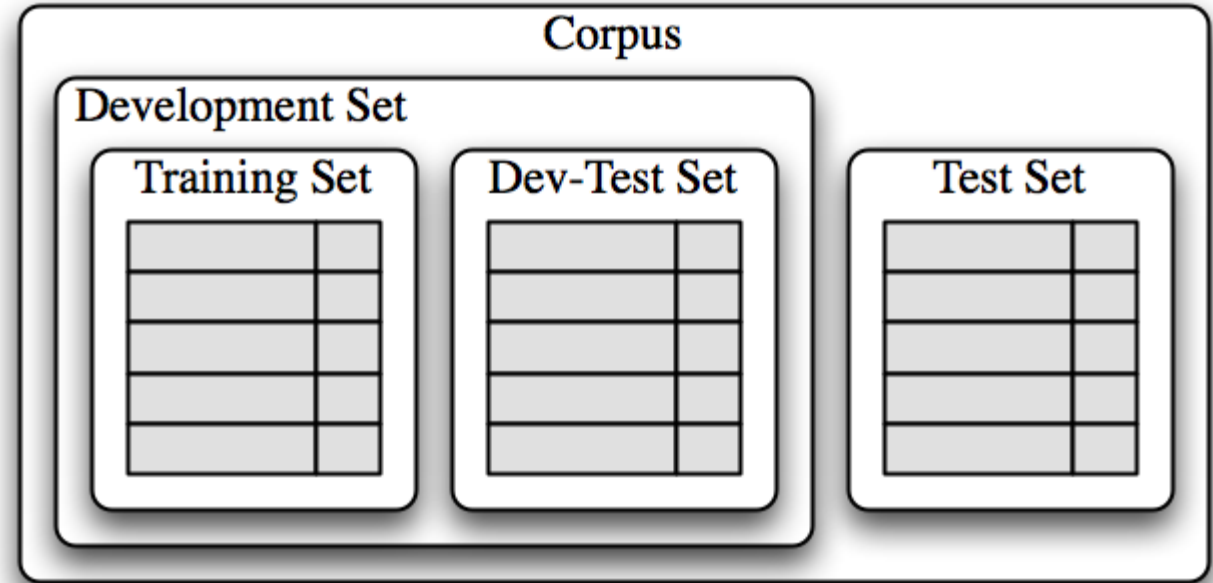
49

Set-up for experiments

Set-up for experiments

50

- Before you start: split into development set and test set.
- Hide the test set
- Split development set into Training and Development-Test set
- Use training set for training a learner



- Use Dev(-Test) for repeated evaluation in the test phase
- **Finally test on the test set!**

Procedure

51

1. Train classifier on training set
2. Test it on dev-test set
3. Compare to earlier runs, is this better?
4. Error analysis: What are the mistakes (on dev-test set)
5. Make changes to the classifier
6. Repeat from 1

=====

- When you have run empty on ideas, test on test set. Stop!

Cross-validation

52

- Small test sets → Large variation in results
- N-fold cross-validation:
 - ▣ Split the development set into n equally sized bins
 - (e.g. $n = 10$)
 - ▣ Conduct n many experiments:
 - In experiment m , use part m as test set and the $n-1$ other parts as training set.
 - ▣ This yields n many results:
 - We can consider the mean of the results
 - We can consider the variation between the results.
 - Statistics!

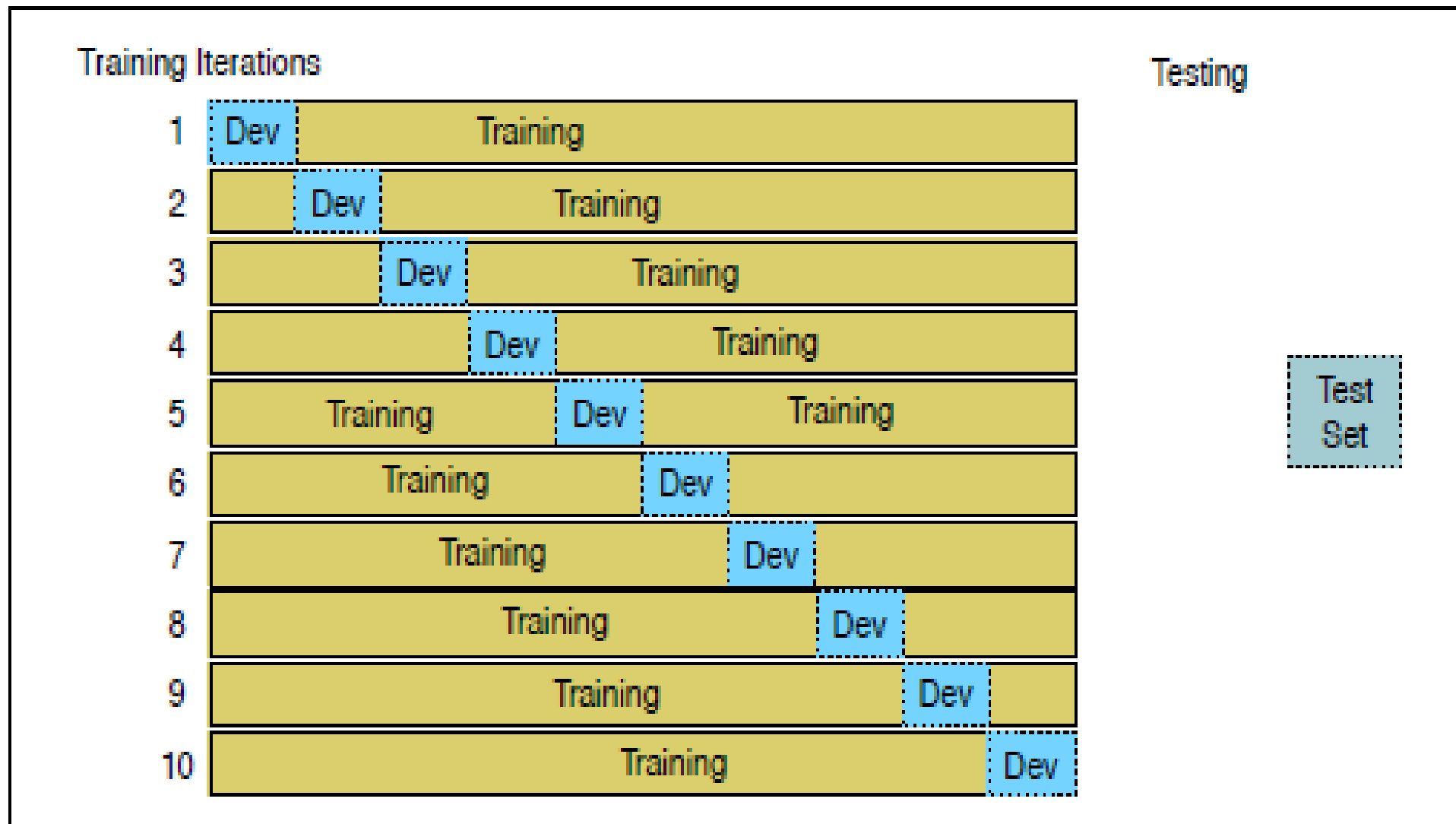


Figure 6.7 10-fold crossvalidation

Evaluation measure: Accuracy

55

- What does accuracy 0.81 tell us?
- Given a test set of 500 sentences:
 - ▣ The classifier will classify 405 correctly
 - ▣ And 95 incorrectly
- A good measure given:
 - ▣ The 2 classes are equally important
 - ▣ The 2 classes are roughly equally sized
 - ▣ Example:
 - Woman/man
 - Movie reviews: pos/neg

But

56

- For some tasks, the classes aren't equally important
 - ▣ Worse to lose an important mail than to receive yet another spam mail
- For some tasks the different classes have different sizes.

Information retrieval (IR)

57

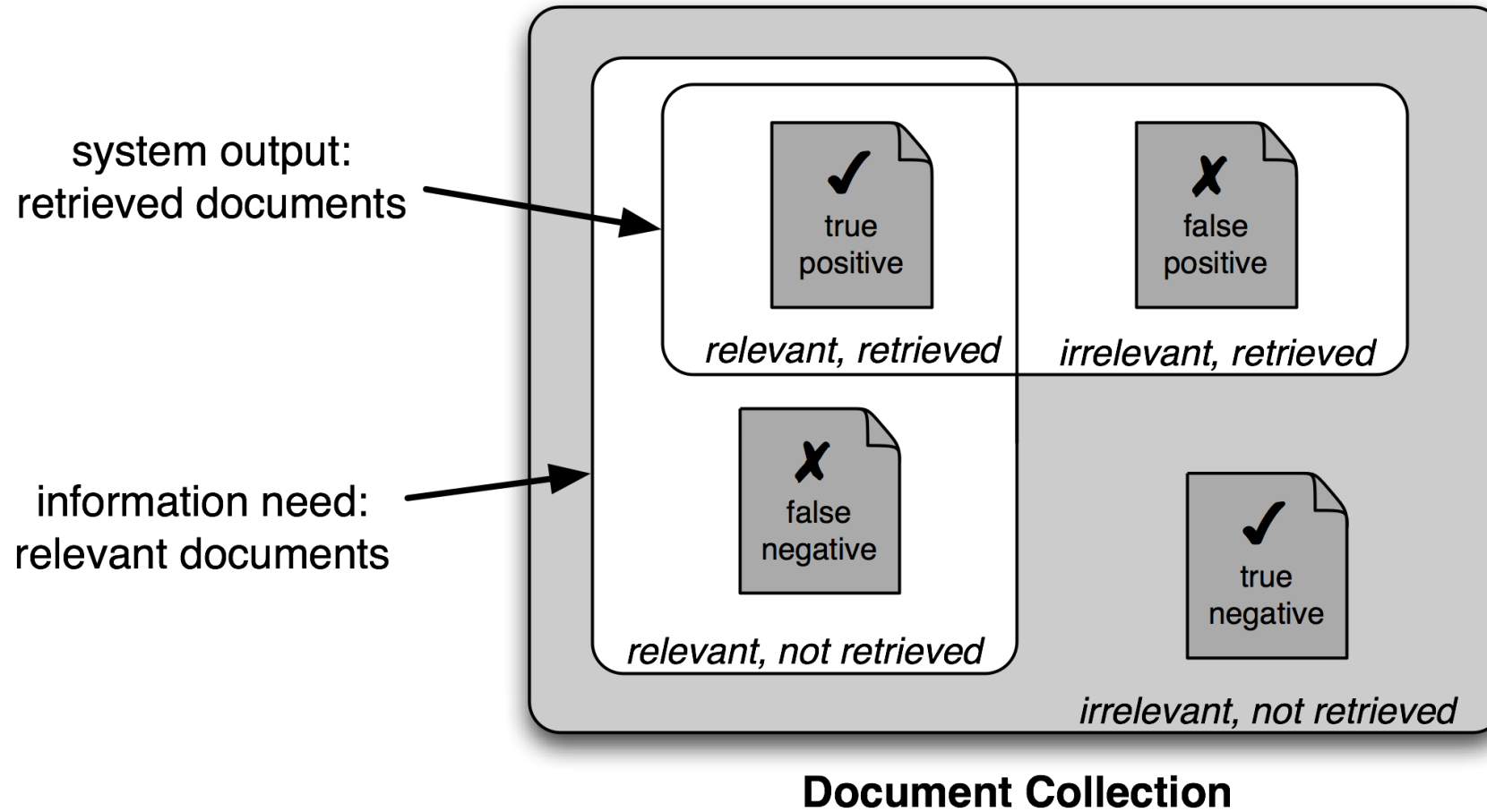
- Traditional IR, e.g. a library
 - ▣ Goal: Find all the documents on a particular topic out of 100 000 documents,
 - Say there are 5
 - ▣ The system delivers 10 documents: all irrelevant
 - What is the accuracy?

- For these tasks, focus on
 - ▣ The relevant documents
 - ▣ The documents returned by the system

- Forget the
 - ▣ Irrelevant documents which are not returned

IR - evaluation

58



Confusion matrix

59

		<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 = $\frac{tp+tn}{tp+fp+tn+fn}$

Figure 6.4 Contingency table

- Beware what the rows and columns are:
 - NLTKs ConfusionMatrix swaps them compared to this table

Evaluation measures

60

		Is in C	
		Yes	NO
Classifier	Yes	tp	fp
	No	fn	tn

- Accuracy: $(tp+tn)/N$
- Precision: $tp/(tp+fp)$
- Recall: $tp/(tp+fn)$

- F-score combines P and R
- $F_1 = \frac{2PR}{P+R} \left(= \frac{1}{\frac{1}{R} + \frac{1}{P}} \right)$
- F_1 called “harmonic mean”
- General form
 - $F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}}$
 - for some $0 < \alpha < 1$

Confusion matrix

		<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	$\text{precision}_s = \frac{200}{3+30+200}$
		$\text{recall}_u = \frac{8}{8+5+3}$	$\text{recall}_n = \frac{60}{10+60+30}$	$\text{recall}_s = \frac{200}{1+50+200}$	

Figure 6.5 Confusion matrix for a three-class categorization task, showing for each pair of classes (c_1, c_2), how many documents from c_1 were (in)correctly assigned to c_2

□ Precision, recall and f-score can be calculated for each class against the rest

Today - Classification

62

- Motivation
- Classification
- Naive Bayes classification
- NB for text classification
 - ▣ The multinomial model
 - ▣ The Bernoulli model
- Experiments: training, test and cross-validation
- Evaluation