

# Επεξεργασία Φυσικής Γλώσσας

# Text Classification

Δημήτρης Πασχαλίδης

Τμήμα Επιστήμης Υπολογιστών

Πανεπιστήμιο Κύπρου



# Classification

- Μια αντιστοίχιση  $h$  από δεδομένα εισόδου  $x$  (από instance space  $X$ ) σε μια ετικέτα  $y$  (από categorical space  $Y$ ).
- $X$  = σύνολο όλων documents
- $Y = \{y_1, y_2, \dots, y_k\}$
- $x$  = ένα μονό document
- $y$  = ετικέτα εξόδου για  $x$



# Movie Reviews

□  $Y = \{\text{Positive, Negative, Neutral}\}$

■  $x_1 =$  “This film blew me away, exciting, fast paced, surprisingly gritty, and genuinely had an awesome story...”

■  $y_1 =$  **Positive**

■  $x_2 =$  “... This boring 3 hour movie had no proper plot whatsoever but just a bunch of random scenes and stupid dialogues ...”


■  $y_2 =$  **Negative**



# Opinion Mining

Positive

Negative

 **Big Bird** ✓  
@BigBird

I got the COVID-19 vaccine today! My wing is feeling a little sore, but it'll give my body an extra protective boost that keeps me and others healthy.

 **President Biden** ✓ @POTUS  
United States government official  
Replying to @BigBird and @EricaRHill  
Good on ya, @BigBird. Getting vaccinated is the best way to keep your whole neighborhood safe.

 **Ted Cruz** ✓ @tedcruz  
Replying to @BigBird and @EricaRHill  
Government propaganda...for your 5 year old!

 **Roseanna Renaud** @sdprairiewoman  
Replying to @BigBird and @EricaRHill  
Thank you Big Bird for caring about the health of our youth. Pay no attention to the bully party.

 **Lisa Boothe** ✓ @LisaMarieBoothe  
Replying to @BigBird and @EricaRHill  
Brainwashing children who are not at risk from COVID. Twisted.



# Συγγραφέας: Αρσενικός or Θηλυκός ?

1. By 1925 present-day Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony of Cochin-China; the central area with its imperial capital at Hue was the protectorate of Annam...
2. Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of her greatest assets...

S. Argamon, M. Koppel, J. Fine, A. R. Shimori, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," *Text*, vol. 23, p. 321-346



# Είναι αυτό SPAM ?

**Subject: Important notice!**

**From: Stanford University <newsforum@stanford.edu>**

**Date: October 28, 2011 12:34:16 PM PDT**

**To: undisclosed-recipients;;**

---

Greats News!

You can now access the latest news by using the link below to login to Stanford University News Forum.


<http://www.123contactform.com/contact-form-StanfordNew1-236335.html>

Click on the above link to login for more information about this new exciting forum. You can also copy the above link to your browser bar and login for more information about the new services.

© Stanford University. All Rights Reserved.



# Θέμα άρθρων

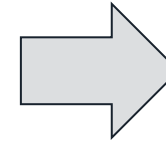


Guardian morning briefing

## Thursday briefing: Russia launches attack on Ukraine

Explosions and airstrikes reported including attack on fighter plane base ... dismay at UN security council ... President Zelenskiy vows 'we will defend ourselves'

by [Warren Murray](#)



- Politics
- Sports
- Lifestyle
- Gossip
- Health
- ...

# Text Classification Πρόβλημα

Δοσμένου ενός κειμένου  $w = (w_1, w_2, \dots, w_T) \in V^*$ ,  
πρόβλεψη ετικέτας  $y \in Y$





# Text Classification Εφαρμογές

Task	X	Y
Language Identification	Text	{ English, Greek, Mandarin, ... }
SPAM Classification	Email	{ SPAM, not SPAM }
Authorship Attribution	Text	{ J.K. Rowling, George R. R. Martin, ... }
Genre Classification	Novel	{ Detective, Romance, Fantasy, ... }
Sentiment Classification	Text	{ Positive, Negative, Neutral, Mixed }
Hate-speech Classification	Text	{ Toxic, Sexist, Hate, Racism, ... }
Veracity Classification	Text	{ Fake, Real, Reliable, Unreliable }

# Classification Μέθοδοι:

## Rule-based Classification

- ❑ Κανόνες που βασίζονται σε μοτίβα, συνδυασμούς λέξεων ή άλλα χαρακτηριστικά.
  - SPAM / UNRELIABLE:
    - ❑ Black-listed πηγές e.g. [www.realnews.124.info](http://www.realnews.124.info)
  - Εάν οι κανόνες βελτιωθούν από τους ειδικούς → υψηλές βαθμολογίες ακρίβειας.

**Δημιουργία / Διατήρηση / Φιλτράρισμα κανόνων**  
→ **Time Consuming & Resource Demanding**



# Classification Μέθοδοι : Supervised Machine Learning

□ Δεδομένα εισόδου:

- ένα έγγραφο  $d$
- ένα σταθερό σύνολο κλάσεων / ετικετών  $C = \{ c_1, c_2, \dots, c_k \}$
- ένα εκπαιδευτικό σύνολο annotated έγγραφα

$$m = \{ (d_1, c_1, \dots, (d_m, c_m) ) \}$$

□ Έξοδος:

- ένας μαθημένος classifier  $y : d \rightarrow c$

## Examples of Classifiers:

- Naive Bayes
- Logistic Regression
- Support Vector Machines
- k-Nearest Neighbors

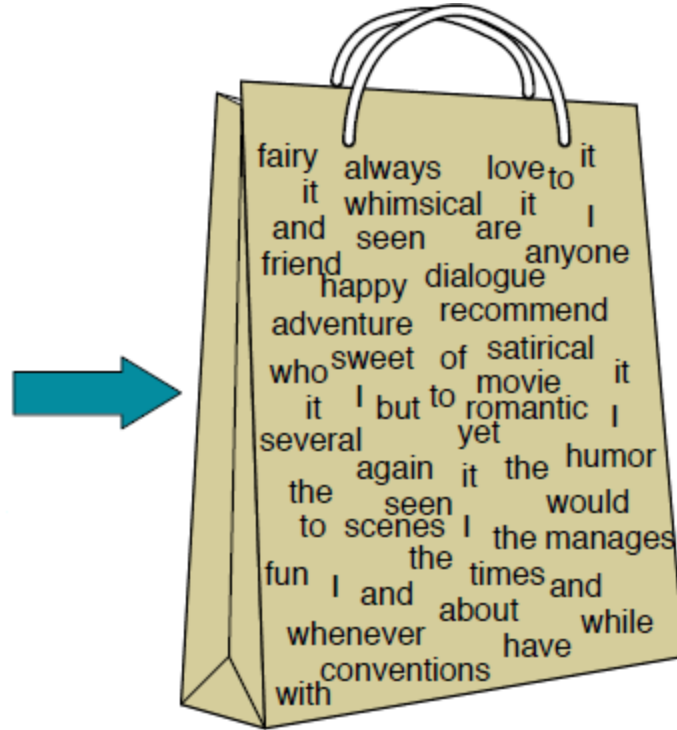
...

# Text Vectorization

- ❑ Οι υπολογιστές δεν κατανοούν το κείμενο, τη σύνταξη και τη σχέση μεταξύ των λέξεων → **Χρειάζεστε έναν τρόπο αναπαράστασης λέξεων με αριθμούς.**
- ❑ Text vectorization: Η διαδικασία μετατροπής κειμένου σε αριθμητική αναπαράσταση.
  - Δημοφιλής βασική αναπαράσταση κειμένου:  
**Bag of Words (BoW)**

## Bag-of-Words Representation

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

- Ο bag-of-words είναι μια αναπαράσταση σταθερού μήκους, η οποία αποτελείται από ένα διάνυσμα καταμέτρησης λέξεων:

$\mathbf{s}$  = “*It was the best of times, it was the words of times*”



$\mathbf{x}$  = [ *aardvark*, ..., *best*, ..., *it*, ..., *of*, ..., *zyther* ]

- Το μήκος του  $\mathbf{x}$  ισούται με το μέγεθος του λεξιλογίου
- Για κάθε  $\mathbf{x}$ , μπορεί να υπάρχουν πολλά  $w$ , ανάλογα με τη σειρά των λέξεων

# Linear Classification επί BoW

□ Εστω  $\psi(\mathbf{x}, \mathbf{y})$  χαρακτηρίζει την συμβατότητα του bow  $\mathbf{x}$  και ετικέτας  $\mathbf{y}$ , τότε  $\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \psi(\mathbf{x}, \mathbf{y})$

□ Σε ένα linear classifier, Αυτή η συνάρτηση

βαθμολόγησης έχει τη μορφή 
$$\psi(\mathbf{x}, \mathbf{y}) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_{j=1} \theta_j \cdot f_j(\mathbf{x}, \mathbf{y})$$

- όπου  $\boldsymbol{\theta}$  είναι ένα διάνυσμα **weights** / **coefficients** και  $\mathbf{f}$  είναι ένα **feature function**.

Learning problem: εύρεση σωστών weights  $\boldsymbol{\theta}$ , υποθέτοντας  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$

# Πιθανοτικό Classification

- ❑ Naive Bayes είναι πιθανοτικός classifier με βάση την ακόλουθη στρατηγική:
  - Ορισμός μοντέλου πιθανότητας  $p(x, y)$
  - Εκτιμήστε τις παραμέτρους του μοντέλου πιθανότητας με τη μέγιστη πιθανότητα στο training dataset.
- ❑ **Bayes Rule** για ένγραφο  $d$  και κλάση / ετικέτα  $c$  :

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



# Naive Bayes Classifier

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

MAP είναι “maximum a posteriori” = Πιο πιθανή κλάση

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

Bayes Rule

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

Dropping the denominator

$P(d)$  καταργείται επειδή είναι περιττή. Για  $d$ , η πιθανότητά του είναι η ίδια για όλα τα  $c$ .

# Naive Bayes Classifier



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

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)$$

Έγγραφο  $d$   
εκπροσωπείτε ως  
features  $x_1 \dots x_n$

$O(|X|^n \times |C|)$  parameters  $\rightarrow$  Εάν έχουμε πολλά χαρακτηριστικά, τότε χρειαζόμαστε μεγάλο αριθμό παραδειγμάτων training.

# Multinomial Naive Bayes

## Independence Assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

- **BoW Assumption:** Η τοποθέτηση λέξεων δεν έχει σημασία.
- **Conditional Independence:** Feature πιθανότητες  $P(x_i | c_j)$  είναι ανεξάρτητες δεδομένης της κατηγορίας / ετικέτας  $c$ .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$



# Multinomial Naive Bayes

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c_j) \prod_{x \in X} P(x | c)$$

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

**Για text classification:**

Θέσεις = όλες οι θέσεις  
λέξεων στο έγγραφο τεστ.

# Προβλήματα με Probability Product

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

- ❑ Ο πολλαπλασιασμός πολλών πιθανοτήτων μπορεί να οδηγήσει σε αριθμητικό underflow.

- $0.0006 \times 0.0007 \times 0.0009 \times 0.01 \times 0.5 \times 0.000008 \dots$

- ❑ Λύση: Χρήση logs  $\rightarrow \log(a \times b) = \log(a) + \log(b)$

➔

$$c_{NB} = \operatorname{argmax}_{c_j \in C} \left[ \log P(c_j) + \sum_{i \in \text{positions}} \log P(x_i | c_j) \right]$$

# Learning του Naive Bayes

□ Learning του Naive Bayes:

- Χρήση του Laplace (add-1) smoothing για την αντιμετώπιση μηδενικών πιθανοτήτων.

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

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

## 1. Απόσπασμα λεξιλογίου V από corpus

## 2. Υπολογισμός P(c<sub>j</sub>) όρων:

$\forall c_j \in C \rightarrow \text{docs}_{c_j} = \text{docs από } c_j$

$$P(c_j) \leftarrow \frac{| \text{docs}_{c_j} |}{| \text{total \# documents} |}$$

## 3. Υπολογισμός P(w<sub>k</sub> | c<sub>j</sub>) όρων :

$\text{text}_{c_j} = \text{συγχώνευση docs}_{c_j}$

$\forall w_k \in V \rightarrow n_k = \text{occur}(w_k, \text{text}_{c_j})$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha | \text{Vocabulary} |}$$

# Sentiment Analysis με Naive Bayes

Document	Category
Just plain boring	Negative
Entirely predictable and lacks energy	Negative
No surprises and very few laughs	Negative
Very powerful	Positive
The most fun film of the summer	Positive
Predictable with no fun	?

Training  
 Test

**1. + / - Πιθανότητες**

$$\hat{P}(c_j) = \frac{N_{c_j}}{N_{total}} \quad P(-) = 3 / 5$$

$$\quad \quad \quad \quad \quad \quad P(+) = 2 / 5$$

# Sentiment Analysis με Naive Bayes

Document	Category
Just plain boring	Negative
Entirely predictable and lacks energy	Negative
No surprises and very few laughs	Negative
Very powerful	Positive
The most fun film of the summer	Positive
Predictable <del>with</del> no fun	?

Training  
Test

## 1. + / - Πιθανότητες

$$\hat{P}(c_j) = \frac{N_{c_j}}{N_{total}} \quad P(-) = 3 / 5$$

$$P(+) = 2 / 5$$

## 2. Clean Texts

Remove “*with*”.



# 3. Υπολογισμός Training Likelihoods

Document	Category
Just plain boring	Negative
Entirely predictable and lacks energy	Negative
No surprises and very few laughs	Negative
Very powerful	Positive
The most fun film of the summer	Positive
Predictable <del>with</del> no fun	?

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

$$P(\text{"predictable"}|-) = \frac{1 + 1}{14 + 20}$$

$$P(\text{"predictable"}|+) = \frac{0 + 1}{9 + 20}$$

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

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



# 4. Test Set Scoring

Document	Category
Just plain boring	Negative
Entirely predictable and lacks energy	Negative
No surprises and very few laughs	Negative
Very powerful	Positive
The most fun film of the summer	Positive
Predictable <b>with</b> no fun	?

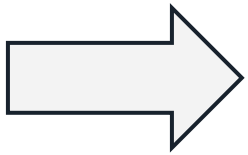
$$P(-)P(S| -) = \frac{3}{5} \times \frac{2 \times 2 \times 1}{34^3} = 6.1 \times 10^{-5}$$

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

Το έγγραφο δοκιμής  
είναι πιο πιθανό να είναι  
**NEGATIVE** class

# Naive Bayes και LM

- ❑ Naive Bayes classifier μπορεί να χρησιμοποιήσει διαφορετικά features:
  - URLs, emails, dictionaries, numerical και network features.
- ❑ Ωστόσο, αν χρησιμοποιούμε μόνο word features, και εξετάστε το whole corpus (not just a subset of words)



Naive Bayes είναι παρόμοιο με Language Modeling από τις προηγούμενες διαλέξεις.

# Naive Bayes και LM

- ❑ Αντιστοίχιση κάθε λέξης:  $P( word | c )$
- ❑ Ανάθεση κάθε πρότασης:  $P( s | c ) = \prod P( word | c )$

Class POSITIVE	
0.1	I
0.1	love
0.01	this
0.05	fun
0.1	film

I	love	this	fun	film
0.1	0.1	0.05	0.01	0.1

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

# Naive Bayes ως LM

□ Ποια τάξη αποδίδει τη μεγαλύτερη πιθανότητα σε  $s$ ?

Class POSITIVE	
0.1	I
0.1	love
0.01	this
0.05	fun
0.1	film

Class NEGATIVE	
0.2	I
0.001	love
0.01	this
0.005	fun
0.1	film

I	love	this	fun	film
0.1	0.1	0.05	0.01	0.1
0.2	0.001	0.01	0.005	0.1

$$P(s | POSITIVE) = 0.00000005$$

$$P(s | NEGATIVE) = 0.00000001$$

# Naive Bayes ως LM

□ Ποια τάξη αποδίδει τη μεγαλύτερη πιθανότητα σε  $s$ ?

Class POSITIVE	
0.1	I
0.1	love
0.01	this
0.05	fun
0.1	film

Class NEGATIVE	
0.2	I
0.001	love
0.01	this
0.005	fun
0.1	film

I	love	this	fun	film
0.1	0.1	0.05	0.01	0.1
0.2	0.001	0.01	0.005	0.1

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

$$P(s | NEGATIVE) = 0.0000001$$

# Άλλοι Classifiers

- Naive Bayes ονομάζεται έτσι λόγω του:
  - Αγνοώντας αφελώς τις εξαρτήσεις μεταξύ των λέξεων και τη σειρά τους, αντιμετωπίζοντας κάθε λέξη ως εξίσου ενημερωτική.
  - Οι feature πιθανότητες  $P(x_i | c_j)$  είναι ανεξάρτητες δεδομένης της κατηγορίας / ετικέτας  $c$ .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$

- Discriminative classifiers αποφύγουν αυτό το ζήτημα, χωρίς να μοντελοποιούν τη "παραγωγική" πιθανότητα  $p(x)$

# Perceptron Classifier

□ Υπόθεση που βασίζεται σε σφάλματα και όχι στην υπόθεση ανεξαρτησίας.

□ Perceptron Learning Κανόνας:

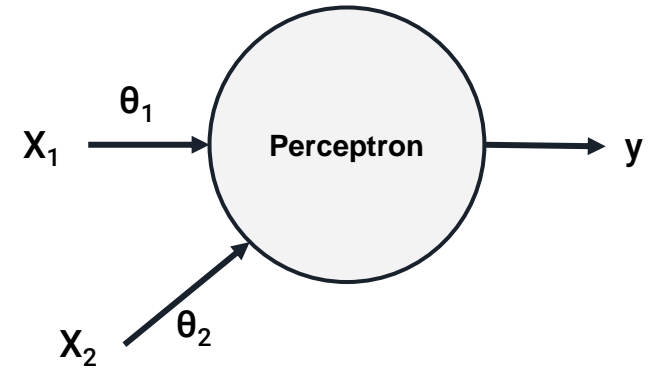
■ Εκτελέστε το τρέχον classifier σε μια παρουσία training data, υπολογίζοντας  $y^{\wedge} = \operatorname{argmax}_y \psi(x^{(i)}, y)$

■ Εάν η πρόβλεψη είναι λανθασμένη:

○ Αυξήστε τα βάρη για το features της πραγματικής ετικέτας.

○ Μειώστε τα βάρη για το features της προβλεπόμενης ετικέτας.  $\theta \leftarrow \theta + f(x^{(i)}, y^{(i)}) - f(x^{(i)}, y^{\wedge})$

■ Επαναλάβετε μέχρι τα training instances είναι σωστά classified, ή εξαντλείται ο χρόνος.





# Perceptron Classifier

---

## Algorithm 3 Perceptron learning algorithm

---

```

1: procedure PERCEPTRON( $\mathbf{x}^{(1:N)}, \mathbf{y}^{(1:N)}$ )
2:    $t \leftarrow 0$ 
3:    $\boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$ 
4:   repeat
5:      $t \leftarrow t + 1$ 
6:     Select an instance  $i$ 
7:      $\hat{y} \leftarrow \operatorname{argmax}_y \boldsymbol{\theta}^{(t-1)} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$ 
8:     if  $\hat{y} \neq y^{(i)}$  then
9:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$ 
10:    else
11:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$ 
12:  until tired
13:  return  $\boldsymbol{\theta}^{(t)}$ 

```

---

- ❑ Objective: Ελαχιστοποίηση loss function στα βάρη.
- ❑ Loss functions θα πρέπει να είναι:
  - Ένας καλός πληρεξούσιος της ακρίβειας των classifiers.
  - Εύκολη βελτιστοποίηση.

# Περισσότερα για Sentiment Analysis

## ❑ Αντιμετώπιση της Άρνησης

I really like this movie.

■ I really **don't** like this movie.

■ **Doesn't** let us get bored.

■ **Don't** miss it!

❑ Η άρνηση αντιστρέφει το νόημα του συναισθήματος.

❑ Απλή μέθοδος baseline: Δημιουργία νέων λέξεων / tokens με προσθήκη **NOT\_** μεταξύ εμφάνισης άρνησης και στίξης.



Didn't like this movie.

Didn't NOT\_like NOT\_this NOT\_movie.

# Περισσότερα για Sentiment Analysis

## ☐ Sentiment Lexicons

- Αν δεν έχετε αρκετά training / labeled data χρησιμοποιήστε προκατασκευασμένες λίστες λέξεων, με το όνομα **Lexicons**.
- **MPQA Subjectivity Cues**: 2718 POSITIVE και 4912 NEGATIVE λέξεις / φράσεις με ένταση (ισχυρές / αδύναμες)
- **General Inquirer**: 1915 POSITIVE και 2291 NEGATIVE, με annotations Ισχυρών vs Αδύναμων, Ενεργητικών vs Παθητικών, Υπερεκτιμημένων vs. Υποτιμημένων και Ηδονής, Πόνου, Αρετής, Κακίας, Κινήτρων, Γνωστικού Προσανατολισμού, etc



# Περισσότερα για Sentiment Analysis

- **Opinion Lexicon:** 2006 POSITIVE και 4789 NEGATIVE λέξεις / φράσεις με ένταση (δυνατές/αδύναμες)
- **AFINN:** AFINN lexicon είναι μια λίστα αγγλικών όρων που βαθμολογούνται για το σθένος μεταξύ -5 (NEGATIVE) και +5 (POSITIVE)
- **Loughran-McDonald:** Αγγλικό sentiment lexicon δημιουργήθηκε για χρήση με οικονομικά έγγραφα → **Financial Sentiment Analysis**

## ☐ Χρησιμοποιώντας Lexicons

- Διάνυσμα των μετρήσεων κάθε φορά που εμφανίζεται μια λέξη από το λεξικό.



# Text Classification Αξιολόγηση

	Correct	Not Correct
Selected	TP	FP
Not Selected	FN	TN



Confusion Matrix

- Correct:** Τι γνωρίζουμε ότι είναι σωστό στο dataset.
- Not Correct:** Τι γνωρίζουμε ότι δεν είναι σωστό στο dataset.
- Selected:** Τι επέστρεψε ο classifier (αποτέλεσμα).
- Not Selected:** Τι δεν επέστρεψε ο classifier.
- TP (True Positive):** Τι επέστρεψε ο classifier και είναι ορθό.
- FP (False Positive):** Τι επέστρεψε ο classifier και δεν είναι ορθό.
- FN (False Negative):** Τι δεν επέστρεψε ο classifier αλλά ήταν ορθό.
- TN (True Negative):** Τι δεν επέστρεψε ο classifier και δεν ήταν ορθό.

# Precision, Recall, και F Measure

- ❑ **Precision:** % των στοιχείων που πεστρεψε και είναι σωστά
- ❑ **Recall:** % σωστών στοιχείων που επέστρεψε
- ❑ **F Measure:** Συνδυασμένο μέτρο που αξιολογεί την Precision και Recall tradeoff (Weighted Harmonic Mean)

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

$$F = 2PR$$

Συνήθως χρησιμοποιείται ως ισορροπημένη F1 με  $\beta=1$  (or  $\alpha=1/2$ )

	Correct	Not Correct
Selected	TP	FP
Not Selected	FN	TN

# Multi-Class Αξιολόγηση

□ **Precision:** Ποσοστό εγγράφων της κλάσης  $i$  και αφορούν την τάξη  $i$ :

$$\frac{c_{ii}}{\sum_j c_{ij}}$$

□ **Recall:** Ποσοστό εγγραφών της κλάσης  $i$  που έγιναν classified σωστά:

$$\frac{c_{ii}}{\sum_j c_{ji}}$$

□ **Accuracy:** Ποσοστό εγγραφών που έγιναν classified σωστά:

$$\frac{\sum_i c_{ii}}{\sum_j \sum_i c_{ij}}$$

# Micro- και Macro- Averaging

□ Για multi-class classification, συνδυάστε μετρήσεις απόδοσης ανά κατηγορία με micro- και macro- averaging.

- **Macro-Averaging:** Υπολογισμός απόδοσης για κάθε κλάση και μετά average.
- **Micro-Averaging:** Υπολογισμός contingency table για όλες τις τάξεις και αξιολόγηση

	<u>Class 1</u>		<u>Class 2</u>		<u>Class 3</u>	
	Truth: YES	Truth: NO	Truth: YES	Truth: NO	Truth: YES	Truth: NO
CLF: YES	10	10	90	10	100	20
CLF: NO	10	970	10	890	20	1860

■ **Macro-Averaged Precision:**  $(0.5 + 0.9) / 2 = 0.7$

■ **Micro-Averaged Precision:**  $100 / 120 = 0.83$

■ Micro-Averaged Score is dominated by score on common classes.



# Περίληψη

- Τύποι Classification:
  - Rule-based Classification
  - Supervised Classification
- Εισαγωγή στο Text Vectorization → Είσοδος σε ML models.
  - Παράδειγμα vectorization: Bag-of-Words
- Πιθανοτικό Classification
  - Παράδειγμα Naive Bayes
  - Naive Bayes ως Language Model
- Classification Αξιολόγηση.



# Resources

- Jurafsky, D. and H. Martin Justin, Chapter 4. "Naive Bayes and Sentiment Classification" Speech and Language Processing
- Jurafsky, D. and H. Martin Justin, Chapter 5. "Logistic Regression" Speech and Language Processing
- Python Tutorial on Text Classification using SKLearn:  
[https://scikit-learn.org/stable/tutorial/text\\_analytics/working\\_with\\_text\\_data.html](https://scikit-learn.org/stable/tutorial/text_analytics/working_with_text_data.html)

