

Advanced quantitative text analysis (2022W)

,

Day 3

Today

9:45-10:40	Input: Dictionaries, Regular Expressions
10:40-10:50	Coffee Break
10:50-11:45	Coding
11:45-12:45	Lunch
12:45-13:40	Input: Supervised Machine Learning
13:40-13:50	Coffee Break
13:50-14:45	Coding

Dictionaries, Regular Expressions

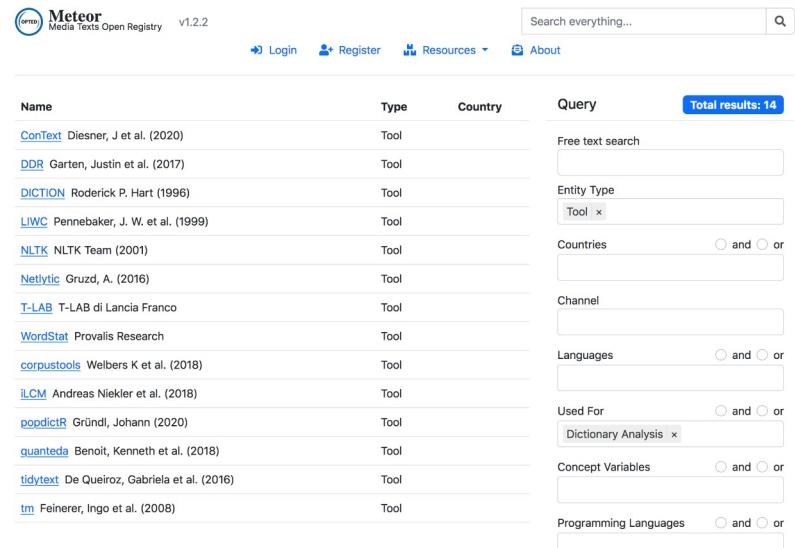
Day 3 Session 1

Dictionary method

Rule-based classification approach when categories are known

- List of words (or phrases) that indicate a category
- Create your own or use/edit existing dictionaries
- Always: validate!

Some existing dictionaries and dictionary application tools



The screenshot shows the Meteor Media Texts Open Registry version 1.2.2. At the top, there is a navigation bar with links for Login, Register, Resources, and About. A search bar is located at the top right. Below the navigation, a table lists various resources categorized by Type (Tool), Country, and Query. To the right of the table, there are several search filters: Free text search, Entity Type (with a 'Tool' option selected), Countries (with radio buttons for 'and' and 'or'), Channel, Languages (with radio buttons for 'and' and 'or'), Used For (with a 'Dictionary Analysis' option selected), Concept Variables (with radio buttons for 'and' and 'or'), and Programming Languages (with radio buttons for 'and' and 'or').

Name	Type	Country	Query
ConText Diesner, J et al. (2020)	Tool		Free text search
DDR Garten, Justin et al. (2017)	Tool		Entity Type
DICTION Roderick P. Hart (1996)	Tool		Countries
LIWC Pennebaker, J. W. et al. (1999)	Tool		Channel
NLTK NLTK Team (2001)	Tool		Languages
Netlytic Gruzd, A. (2016)	Tool		Used For
T-LAB T-LAB di Lancia Franco	Tool		Dictionary Analysis
WordStat Provalis Research	Tool		Concept Variables
corpusTools Welbers K et al. (2018)	Tool		Programming Languages
ilCM Andreas Niekler et al. (2018)	Tool		
popdictR Gründl, Johann (2020)	Tool		
quanteda Benoit, Kenneth et al. (2018)	Tool		
tidytext De Queiroz, Gabriela et al. (2016)	Tool		
tm Feinerer, Ingo et al. (2008)	Tool		

https://meteor.opted.eu/query?dgraph.type=Tool&used_for=0x1ade5

Dictionary use cases

- Selecting text data according to some defined category ('what is relevant data')
- Classifying documents into known categories

Example: Search string for migration news

Table 3 Boolean search strings used for retrieval of migration-related news articles

Country	Language	Search string
Spain	Spanish	asilo* OR inmigran* OR refugiad* OR migrante* OR migratori* OR "sin papeles" OR "campo de desplazados" OR patera* OR emigra* OR "libre circulación" OR "fuga de cerebros"
UK	English	asyl* OR immigrant* OR immigrat* OR migrant* OR migrat* OR refugee* OR foreigner* OR "undocumented worker*" OR "guest worker*" OR "foreign worker*" OR emigrat* OR "freedom of movement" OR "free movement"
Germany	German	asyl* OR immigrant* OR immigriert* OR immigrat* OR migrant* OR migrat* OR flüchtling* OR ausländer* OR zuwander* OR zugewander* OR einwander* OR eingewander* OR gastarbeiter* OR "ausländische arbeitnehmer*" OR emigr* OR auswander* OR ausgewander* OR personenfreizügigkeit* OR arbeitnehmerfreizügigkeit* OR "freier personenverkehr"

Example: Classify actor salience

Objective: salience of women and men migrant's in the news across time and outlet

- Dictionary approach to measure mentions of women and men migrants in German news articles



Example: Classify actor salience

Keyword list more complex than initially thought

Recall: C7 Precisions: 01

```

349
350
351
352 person_f_migrant_endIn_regex = c("[Aa]sylantin", "[Aa]sylbewerberin", "[Aa]
353     "[Zz]uwanderin", "[Ee]inwanderin", "[Gg]a
354     "[Aa]usländische\\w{0,2}\\s[Bb]ürgerin",
355
356
357 f_relation_regex = c("[Mm]eine\\w{0,2}\\s[FF]rau(en)\\s|\\W", "[Mm]ein\\
358     "[Dd]eine\\w{0,2}\\s[FF]rau(en)\\s|\\W", "[Dd]ein\\w{
359     "[Ss]eine\\w{0,2}\\s[FF]rau(en)\\s|\\W", "[Ss]ein\\w{
360     "[Ii]hr\\w{0,2}\\s[FF]rau(en)\\s|\\W", "[Ii]hr\\w{0,2}
361     "[Uu]ns\\w{0,2}\\s[FF]rau(en)\\s|\\W", "[Uu]ns(er)re
362     "[Ee]u(er)re)\\w{0,2}\\s[FF]rau(en)\\s|\\W", "[Ee]u(e
363
364
365
366 f_relation_regex <- paste(apos_closed, f_relation_regex, "", sep = "") #add
367 df_f_relation_regex <- data.frame(f_relation_regex) #as. dataframe

```

Table 2. Dictionary subcategories for the measurement of the concepts: migrant women's and migrant men's salience.

Subcategory name	Description	Measured concepts (examples)	
		Migrant women's salience	Migrant men's salience
1. Impersonal designation	General person nominations for migrants	(immigrant, asylum seeker, etc.) with female or gender-neutral word endings (-In, -in, etc.)	(immigrant, asylum seeker, etc.) with generic masculine form or gender-neutral word endings (-In, -in, etc.)
2. Origin: Single words	Nominations that refer to the territorial origin of a person ^a	(French, African, Syrian, etc.) with female or gender-neutral word endings (-In, -in, etc.)	(French, African, Syrian, etc.) with generic masculine form or gender-neutral word endings (-In, -in, etc.)
3. Origin: Combinations of different word groups ^b	General gender-related person nominations + from + general territorial denominations ^c	(woman, girl, mother, sister) + from + (France, Africa, Syria, etc.)	(man, boy, father, brother) + from + (France, Africa, Syria, etc.)
a	General territorial denominations (adjectives) ^c + General gender-related person nominations	(French, African, Syrian, etc.) + (woman, girl, mother, sister)	(French, African, Syrian) + (man, boy, father, brother)
b			
4. In relation	Relational expressions: possessive pronouns + General gender-related person nominations	(my, your, her, his, our, yours, their) + (woman, girl, mother, sister)	(my, your, her, his, our, yours, their) + (man, boy, father, brother)
5. As phrase	Phrases	e.g., "women and children"	e.g., "men and children"
6. Other expressions		e.g., "women from abroad"	e.g., "men from abroad"

As general notes, the dictionary includes the singular and plural version of all words, word endings (e.g., for prepositions) consider the different cases used in the German language.

^aDownloaded from the CLDR (Unicode Common Locale Data Repository) <http://cldr.unicode.org/>, which holds standard name translations of countries and regions (version v33.1).

^bMeasured at the sentence level.

^cManually compiled by a native speaker for all CLDR territorial denominations, which the German language allows (e.g., no separate word for many smaller islands, e.g., Isle of Man, Curaçao). Assisted by the preeminent German language dictionary *duden.de*.

Validate, validate, validate

Data validation

We ended our class yesterday by sharing resources and experiences about data collection.

What possibilities do we have to evaluate the quality of our data?

Data validation

Motivation: source and data selection determines results and conclusions

Are the selected data **sources** and selected **data** points representative for my target concept or discourse?

- Are they relevant?
- Are they representative?

Side note: Data source validation

Validation techniques:

- Rely on expert opinion
- Rely on data source selection of similar research

Relevance of search string validation

- Sampling based on search strings popular (Stryker et al. 2016) and recommended (Barberá et al., 2021)
- Reviews of search string validation procedures
 - out of 83 content analyses, 39% stated the search terms they used, and only 6% discussed their validity (Stryker et al. 2016)
 - out of 105 content analysis studies, 73.3% stated the search terms they used, only 12.4% reported validity metrics (Mahl et al., 2022)
- Careless application of non-validated search terms may lead to noisy inferences (Mahl et al., 2022)

Search string validation guidelines

- (1) formalize the validation process—that is, carefully protocol every step and decision taken prior to and during data collection
- (2) define the universe of relevant data beforehand, which includes an explicit definition of relevance criteria that have to be met
- (3) triangulate inductive and deductive approaches, such as frequency or keyness analyses, to identify the most indicative terms during search term mining, including a consideration of semantic changes of concepts over time as well as varying national or cultural contexts
- (4) carefully validate search terms identified in search term mining by using well-established performance metrics such as precision, recall, and F1-score

Key validation approach

How close is an automated measurement to a more trusted measurement:

Human understanding of text?

Dictionary validation with manually created baseline

Steps

- Code a subset manually and compare manual with automated classification decisions (via recall, precision, F1)
 - Iterative dictionary improvement
 - Ideally: manual coding and dictionary development is performed by different persons
-

Creation of a manual baseline

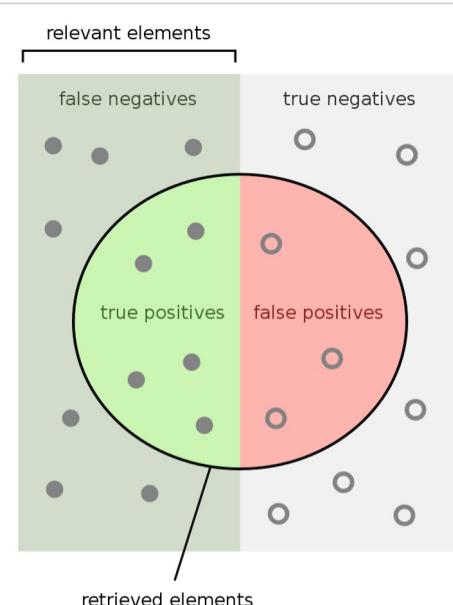
Decisions

- Codebook creation
- Who:
 - Expert coders: Coder recruitment and training sessions
 - Crowdcoders: test questions, majority choice
- Quality assessment: e.g., Inter-coder reliability of involved coders, majority vote
- Documents selected for baseline should be representative for target discourse (e.g., random selection or artificial week)

Recall, precision, F1?

Metrics frequently used to express the validity of a search string & more generally also of automated classification methods

- Precision (P)
- Recall \circledast (R)
- $F1 = 2^* (P * R) / (R + R)$



$$\text{Precision} = \frac{\text{How many retrieved items are relevant?}}{\text{How many relevant items are retrieved?}}$$
$$\text{Recall} = \frac{\text{How many relevant items are retrieved?}}{\text{How many relevant items are there?}}$$

Scharkow, 2012, 133–36

Tools for manual coding

- Google sheets
- AmCAT: <https://vu.amcat.nl/accounts/login/?next=/>
- AnnoTinder: https://github.com/ccs-amsterdam/CCS_annotator

Example: Baseline creation templates

For search strings

Create sampling plan (goal: representative for universe of texts)

Database	Date	Outlet	Number of all articles published that day
APA	Mon 8.10.2018	Standard	136
APA	Tue 9.07.2019	Standard	87
APA	Wed 12.02.2020	Standard	89
APA	Thr 15.04.2021	Standard	98
APA	Fri 27.05.2022	Standard	94

Note: Ideally repeat this procedure for each outlet included; cover the full range of time period investigated

Collect articles

Article id	Date	Text	Manually perceived as relevant (1=yes, 0 = No)	Perceived as relevant by search string (1=yes, 0 = No)
1	Mon 8.10.2018	Kern ist an sich selbst gescheitert. Die SPÖ braucht jetzt mehr Gerechtigkeit und weniger Gockelhaftigkeit ...	1	1
2	Mon 8.10.2018	Impressum und Offenlegung: Herausgeber: Oscar Bronner...	0	0
3	Mon 8.10.2018	Einseitiger Vorschlag. Zu viele Waffen in der Hand der Bürger sind gefährlich. Ein Blick in die USA zeigt, warum. Im Kern geht es...	0	1
...

Code manually
Search with search string

Calculate recall and precision

Dictionary

Pros

- Often needed to select data (search strings)
- High reliability and control
- High transparency and reproducibility

Cons

- Difficulty increases with the latency of the construct

Regular Expressions (regex)

Regular Expressions (regex)

- How to pronounce?

Regular Expressions (regex)

- How to pronounce?
- /gɪf/
- /dʒɪf/

Regular Expressions (regex)

- How to pronounce?
- /gɪf/
- /dʒɪf/
- /ɹε.ɡεks/
- /ɹε.dʒεks/

Regular Expressions (regex)

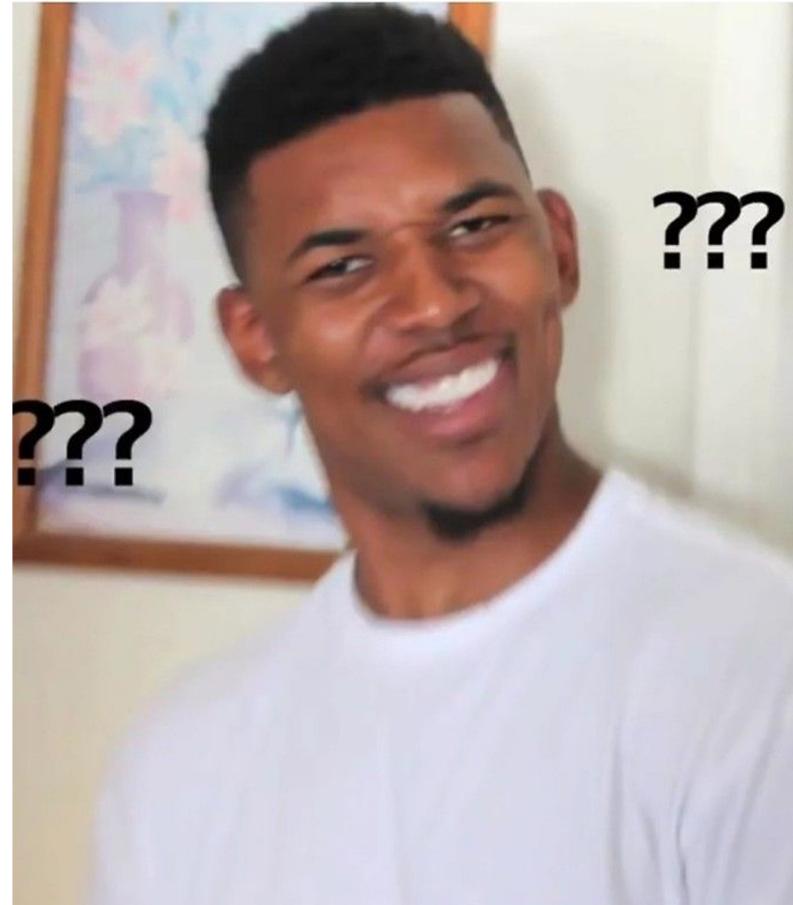
- How to pronounce?
- */gʌf/*
- */dʒɪf/*
- */θθ.gʌks/*
- */rɛ.dʒɛks/*

Regular Expressions (regex)

- How to pronounce?
- /gɪf/
- **/dʒɛks/**
- /rɛ.gɛks/
- **/ɛks.dʒɛks/**

Regular Expressions (regex)

- How to pronounce?
- /gɪf/
- **/dʒɛks/**
- /rɛ.gɛks/
- **/ɛks.dʒɛks/**



regex

- Formal language to specify search strings

regex

- Formal language to specify search strings
 - Insanely difficult
-

regex

- Formal language to specify search strings
- INSANELY difficult

regex

- Formal language to specify search strings
 - ***INSANELY*** difficult
-

regex

- Formal language to specify search strings
- ***INSANELY*** difficult
- Nobody can remember anything

regex

- Formal language to specify search strings
- ***INSANELY*** difficult
- Nobody can remember anything
- Different ***flavours***

regex

- Formal language to specify search strings
- ***INSANELY*** difficult
- Nobody can remember anything
- Different ***flavours***

- “Some people, when confronted with a problem, think “I know, I’ll use regular expressions.” Now they have two problems.” Jamie Zawinski

*Regular Expressions
for Perl, Ruby, PHP,
Python, C, Java, and .NET*

2nd Edition

Regular Expression

Pocket Reference



O'REILLY®

Tony Stubblebine

*Regular Expressions
for Perl, Ruby, PHP,
Python, C, Java, and .NET*

2nd Edition

Regular Expression

Pocket Reference



O'REILLY®

Tony Stubblebine

128 pages!!!

- Disjunctions

Grimmer / Jurafsky
Cheat-sheet

RE	Match	Example Patterns Matched
[mM]oney	Money or money	“Money”
[abc]	‘a’, ‘b’, or ‘c’	“Investing in <u>Iran</u> ”
[1234567890]	any digit	“is <u>dangerous</u> <u>business</u> ”
[\\.]	A period	“sitting on \$ <u>7.5</u> billion dollars”
		“ <u>2005</u> and <u>2006</u> , more than ”
		“\$ <u>150</u> million dollars”
		“ ‘Run!', he screamed.”

- Ranges

RE	Match	Example Patterns Matched
[A-Z]	an upper case letter	“Rep. <u>Anthony</u> <u>Weiner</u> ”
(D-Brooklyn & Queens)		“(D-Brooklyn & <u>Queens</u>)”
[a-z]	a lower case letter	“ACORN’s <u>s</u> ”
[0-9]	a single digit	“(9th CD) ”

- Negations

RE	Match	Example Patterns Matched
[^A-Z]	not an upper case letter	“ACORN’s”
[^Ss]	neither ‘S’ nor ‘s’	“ <u>ACORN</u> ’s”
[^\.]	not a period	“‘Run!’, he screamed.”

- Optional Characters: ?, *, +

RE	Match	Example Patterns Matched
colou?r	Words with u 0 or 1 times	“color” or “colour ”
oo*h!	Words with o 0 or more times	“oh!” or “ooh!” or “oooh!”
o+h!	Words with o 1 or more times	“oh!” or “ooh!” or “oooooh!” or

- Start of the line anchor ^, end of the line anchor \$

Grimmer / Jurafsky
Cheat-sheet

RE	Match	Example Patterns Matched
^ [A-Z]	Upper case start of line	“ <u>Palo Alto</u> ” “the town of <u>Palo Alto</u> ” “ <u>the town of Palo Alto</u> ” “ <u>Palo Alto</u> ” “ <u>Palo Alto</u> ” “ <u>the town of Palo Alto</u> ”
^ [^A-Z]	Not upper case start of line	
^. .	Start of line	
. \$	Identify character that ends a line	“Wait!_” “This is the end.”

- “Or” | statements, Useful short hand

RE	Match	Example Patterns Matched
yours mine	Matches “yours” or “mine”	“it’s either <u>yours</u> or <u>mine</u> ”
\ d	Any digit	“ <u>1-Mississippi</u> ”
\ D	Any non-digit	“ <u>1-Mississippi</u> ”
\ s	Any whitespace character	“ <u>1,_2</u> ”
\ S	Any non-whitespace character	“ <u>1,_2</u> ”
\ w	Any alpha-numeric	“ <u>1-Mississippi</u> ”
\ W	Any non-alpha numeric	“ <u>1-Mississippi</u> ”

How difficult to regex an email

How difficult to regex an email

Rather

How difficult to regex an email

```
(?:[a-zA-Z0-9!#$%&'*+/=?^_-`{|}~-]+(?:\.[a-zA-Z0-9!#$%&'*+/=?^_-`{|}~-]+)*|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21\x23-\x5b\x5d-\x7f]|\\"[\x01-\x09\x0b\x0c\x0e-\x7f])*")@(?:([a-zA-Z0-9](?:[a-zA-Z0-9-]*[a-zA-Z0-9])?\.)+[a-zA-Z0-9](?:[a-zA-Z0-9-]*[a-zA-Z0-9])?|[((?:((?:2(5[0-5]|[0-4][0-9])|1[0-9][0-9]|1[1-9]?[0-9]))\.){3}(?:((?:2(5[0-5]|[0-4][0-9])|1[0-9][0-9]|1[1-9]?[0-9]))|[a-zA-Z0-9-]*[a-zA-Z0-9]:(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f]|\\"[\x01-\x09\x0b\x0c\x0e-\x7f])+)\])
```


Helpers

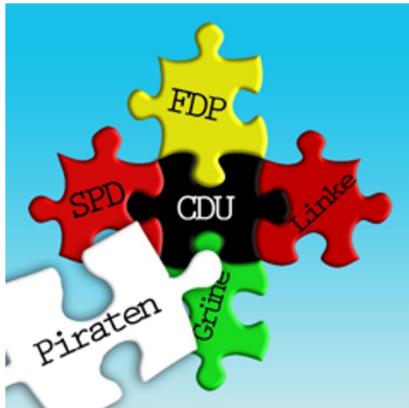
<https://regex101.com>

Recommended R tutorial (Wickham & Grolemund, 2017):
<https://r4ds.had.co.nz/strings.html>

To test regular expressions quickly:
https://spannbaueradam.shinyapps.io/r_regex_tester/

Exercise in R

- Dictionary creation: Visibility of political actors in news headlines
- Calculate recall, precision, F1, Krippendorff's alpha



Supervised classification

Day 3 Session 2

Types of machine learning

1. Supervised
 - o An outcome variable is defined
 - o Focus is on prediction
 2. Unsupervised
 - o No outcome variable has been defined
 - o Focus is on patterns
-

Types of machine learning

1. Supervised
 - An outcome variable is defined
 - Focus is on prediction
 2. Unsupervised
 - No outcome variable has been defined
 - Focus is on patterns
-

Supervised machine learning

Objective:

Classification of documents into pre existing categories

Steps

1. Create a labeled data set
2. Classify documents with supervised learning algorithm
3. Check performance
4. Using the measures

Create a labeled data set

How:

- Human coders annotate parts of the corpus (see slides in session 1 of today)
- Found data (e.g., self-reported profession in users' profile)

Considerations:

- Sampling should be representative for the corpus (e.g., Random, Stratified sample e.g., across time and source)
- Quality of human coding matters (Assess the intercoder reliability)
- Number of documents

Create a labeled data set

Number of documents

- the higher the number of categories and the lower the reliability of the coders, the higher the number of documents (Barberá et al., 2021)
- increase the sizes of manually coded validation dataset as large as possible, preferably to more than $N = 1,300$ (i.e., more than 1% of all data to be examined) assuming acceptable reliability (equal to or higher than .7) (Song et al., 2021)

Table 2. Simulation input parameters.

Factors	Input Parameters
N of human coders	2 (minimum), 5 (intermediate), & 10 (large manual coding)
Intercoder reliability	0.5 (low), 0.7 (acceptable), & 0.9 (high levels of reliability)
N of validation data	600 (0.5%), 1300 (1%), 6500 (5%), & 13000 (10%) of total data
Sampling variability	Random sample vs. nonrandom (biased) subset for validation
Coding per entry	Sole coding vs. duplicated coding for each entry

Song et al., 2021, p. 555

Create a labeled data set

Split labeled data in training data and test data

Training data

- The subset that is used to learn the model parameters

Test data

- Another subset used to evaluate the model's predictive quality
 - Not used for learning!
-

Classify documents with supervised learning

Classifier learns the mapping between features and the labels in the training set

- We define a model $f(Y)=g(X)$
- And apply a learning algorithm to establish which features in X (features extracted from the training documents) matter to recover Y (i.e, the labels of the training documents)
- We fit the model

Classify documents with supervised learning

Considerations:

- Feature representation (Bag of words representation or embeddings)
- Feature selection (remove irrelevant features)
- Classifier selection
 - E.g., Naive Bayes, SVM, KNN, or ensemble methods

Check performance

The fitted model (the trained classifier) is applied to a held-out test set (which is a part of the labeled set but was not used for training the model).

Considerations:

- Danger of overfitting (focus on features that work well with training set but do not generalize)
 - Solutions: cross-validation
- Performance metric (i.e., recall, precision)

Check performance

k -fold cross-validation

- We randomly split the data into k sets (“folds”) of roughly equal size
- Each set is hold out once as test set, while training on the remaining sets
- The problem of a lucky split is reduced



Check performance

Performance metrics

Confusion matrix

		Actual label
Classification (algorithm)	Negative	Positive
Negative	True negative	False negative
Positive	False negative	True positive

$$Accuracy = \frac{True\ Negative + True\ Positive}{True\ Negative + True\ Positive + False\ Negative + False\ Positive}$$

$$Precision_{positive} = \frac{True\ Positive}{True\ Positive + False\ Positive}$$

$$Recall_{positive} = \frac{True\ Positive}{True\ Positive + False\ Negatives}$$

Check performance

Check convergent validity (Adcock & Collier, 2001)

- Compare the measures with other established measures, e.g.,
 - Use trained model to classify texts (open-ended answers relationship uncertainty as described by participants of an online survey) and compare it with the self-assessment in response to a closed question, see Pilny et al. 2019

Using the measures

The classifier is applied to all documents in the corpus

Considerations

- Dealing with the measurement error which occurs through error in the human coding, errors in the classifier
(assess: if error is systematic or not, if consequential for the analysis)

Dictionary vs. supervised machine learning

Category: sentiment

Result: machine learning significantly outperformed dictionary coding



Original Article

The Validity of Sentiment Analysis: Comparing Manual Annotation, Crowd-Coding, Dictionary Approaches, and Machine Learning Algorithms

Wouter van Atteveldt Mariken A. C. G. van der Velden & Mark Boukes

Pages 121-140 | Published online: 28 Jan 2021

Download citation <https://doi.org/10.1080/19312458.2020.1869198>

Full Article

Figures & data

References

Citations

Metrics

Licensing

Reprints & Permissions

View PDF

View EPUB

ABSTRACT

Dictionary vs. supervised machine learning

- Dictionaries can be applied directly to a new corpus (but validate!)
- Supervised machine learning requires (potentially larger amounts) labeled data
- If the training sample is large enough supervised learning will outperform dictionaries

Additional considerations

- Hyperparameter selection
 - Via systematic comparison of different hyperparameters per algorithm
- Random undersampling (Galar et al., 2011)
 - Method to deal with unbalanced classes: use the max. number of positive instances per class and randomly sample the same number of instances of the negative class

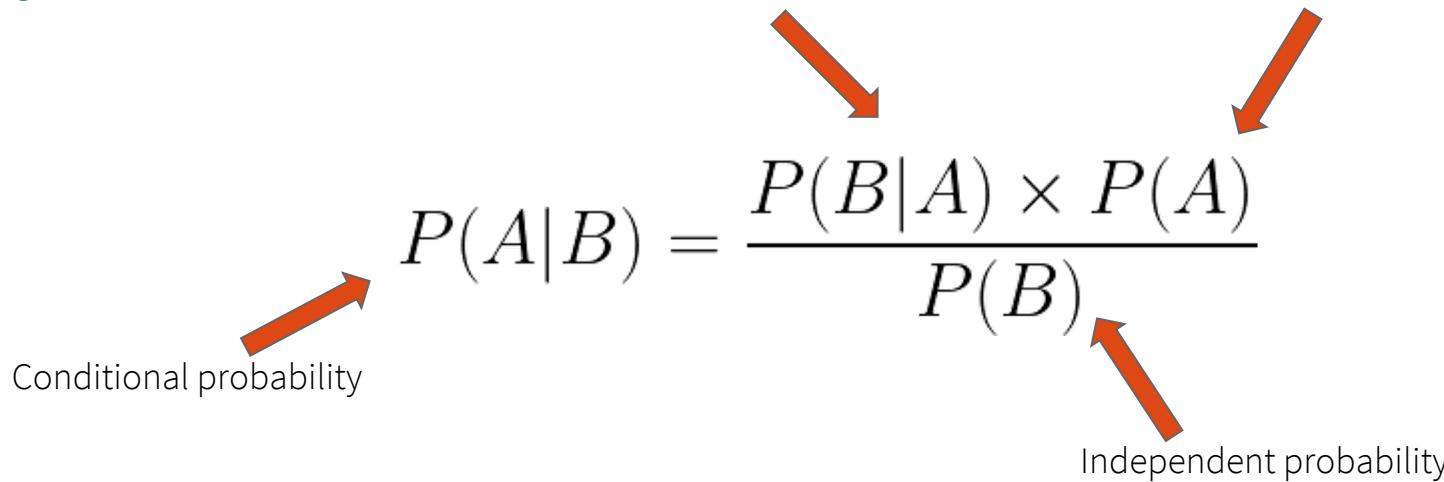
Naïve Bayes Classifier

- Probabilistic classifier
 - Simple
 - Fast
 - Good Accuracy
-

Bayes Theorem

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

Bayes Theorem

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


The diagram illustrates the components of Bayes' Theorem. It features the formula $P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$. Four orange arrows point from labels to specific terms in the formula: one arrow points from 'Conditional probability' to the term $P(B|A)$; another arrow points from 'Independent probability' to the term $P(A)$; a third arrow points from 'Conditional probability' to the term $P(A|B)$; and a fourth arrow points from 'Independent probability' to the term $P(B)$.

Conditional probability

Independent probability

Conditional probability

Independent probability

Bayes Theorem

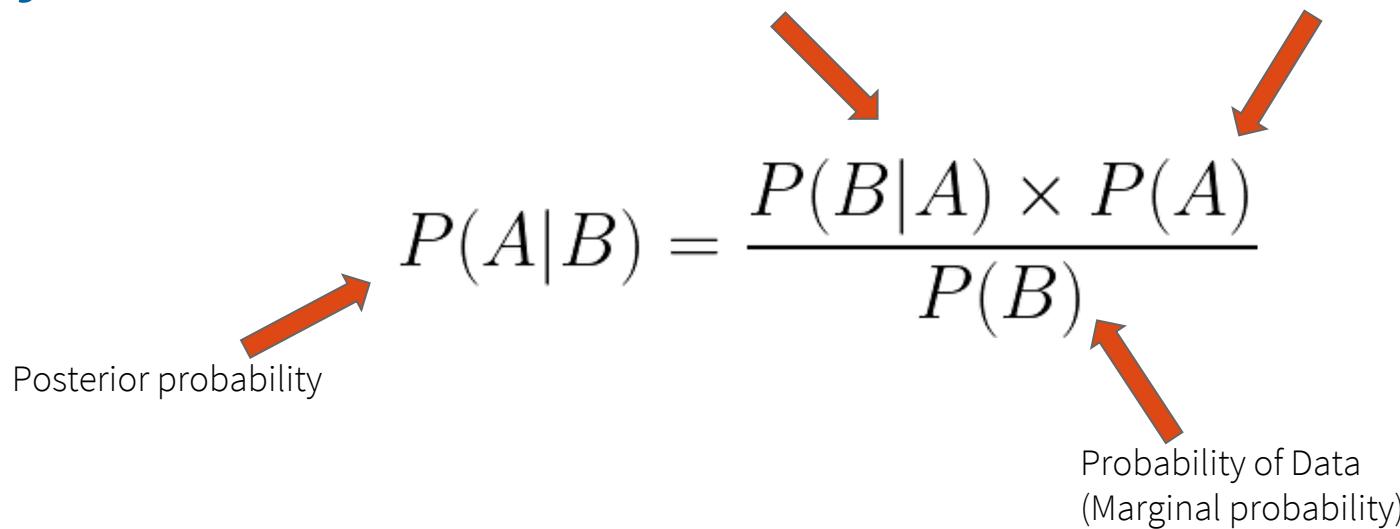
$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

Posterior probability

Likelihood

Prior probability

Probability of Data
(Marginal probability)



The diagram illustrates Bayes Theorem with a central equation $P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$. Four red arrows point from text labels to specific terms in the equation:

- An arrow points from "Posterior probability" to the term $P(A|B)$.
- An arrow points from "Likelihood" to the term $P(B|A)$.
- An arrow points from "Prior probability" to the term $P(A)$.
- An arrow points from "Probability of Data (Marginal probability)" to the term $P(B)$.

Bayes Theorem

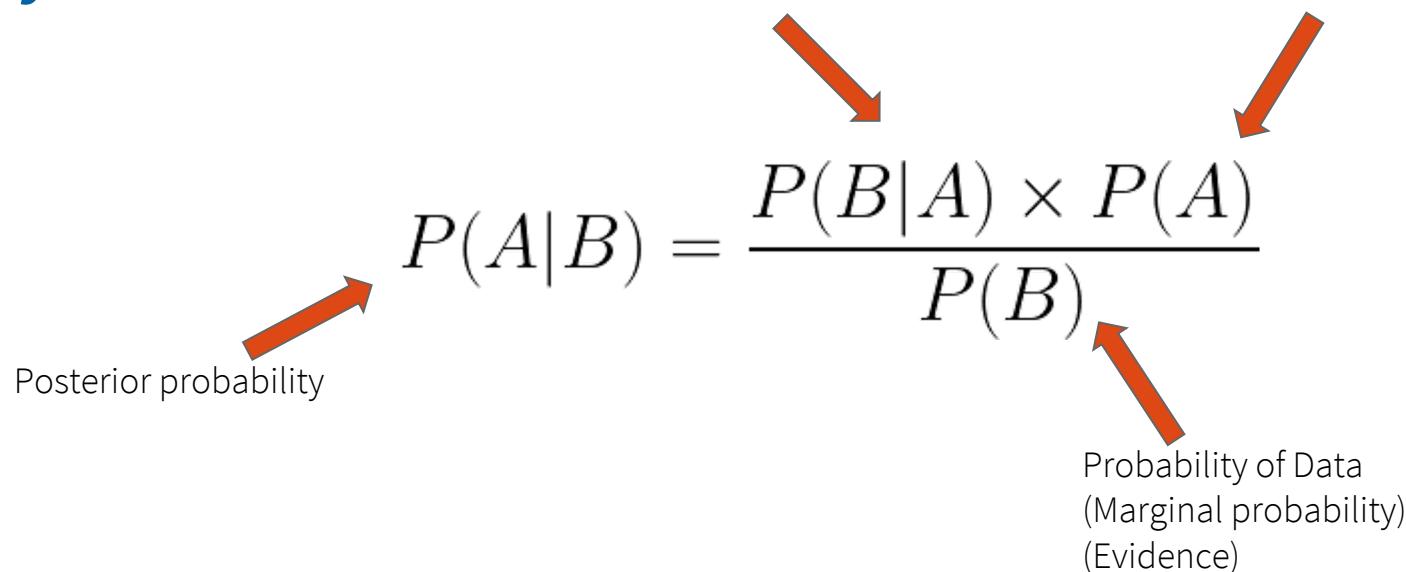
$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

Posterior probability

Likelihood

Prior probability

Probability of Data
(Marginal probability)
(Evidence)



The diagram illustrates Bayes Theorem with a central equation $P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$. Four red arrows point from text labels to specific components of the equation:

- An arrow points from "Posterior probability" to the term $P(A|B)$.
- An arrow points from "Likelihood" to the term $P(B|A)$.
- An arrow points from "Prior probability" to the term $P(A)$.
- An arrow points from "Probability of Data (Marginal probability) (Evidence)" to the term $P(B)$.

Bayes Theorem

$$P(A|B) \propto P(B|A) \times P(A)$$

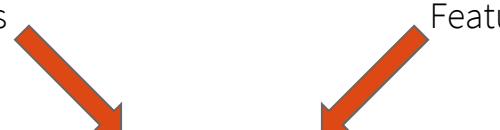
Naïve Bayes Classifier

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

Naïve Bayes Classifier

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

Class Features



Naïve Bayes Classifier

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

Class Features



Features are assumed to be independent. Hence, “**Naïve**”

Naïve Bayes Classifier

$$P(C_k|\mathbf{x}) = \frac{P(C_k) \times P(\mathbf{x}|C_k)}{P(\mathbf{x})}$$

Naïve Bayes Classifier

$$P(C_k|\mathbf{x}) \propto P(C_k) \times P(\mathbf{x}|C_k)$$

Naïve Bayes Classifier

$$\begin{aligned} p(C_k \mid x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &\propto p(C_k) p(x_1 \mid C_k) p(x_2 \mid C_k) p(x_3 \mid C_k) \dots \\ &\propto p(C_k) \prod_{i=1}^n p(x_i \mid C_k), \end{aligned}$$

Decision Rule

$$\hat{y} = \operatorname{argmax} p(C_k) \prod_{i=1}^n p(x_i|C_k)$$



Implemented in many stats/ML packages



Support Vector Machine

- Comes from computer science
 - Very good
 - Rather difficult math
-
- Considered one of the best off-the-shelf classification algorithms

Hyperplane

- $n-1$ dimensional plane that separates the n -dimensional space

Hyperplane

- n-1 dimensional plane that separates the n-dimensional space
- 2-dimensional hyperplane:

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$$

Hyperplane

- n-1 dimensional plane that separates the n-dimensional space
- 2-dimensional hyperplane:
- Line equation

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$$

Hyperplane

- n-1 dimensional plane that separates the n-dimensional space
- 2-dimensional hyperplane:
- Line equation

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 = 0$$

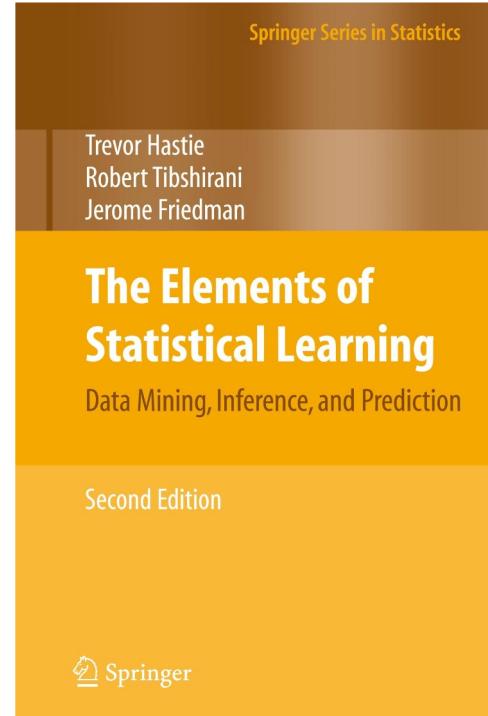
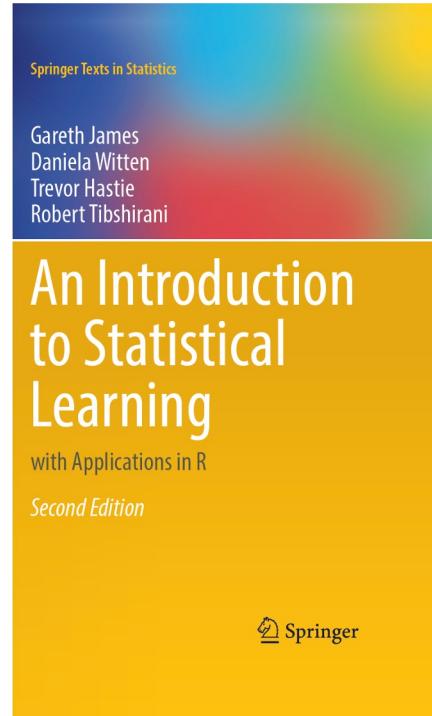
$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p = 0$$

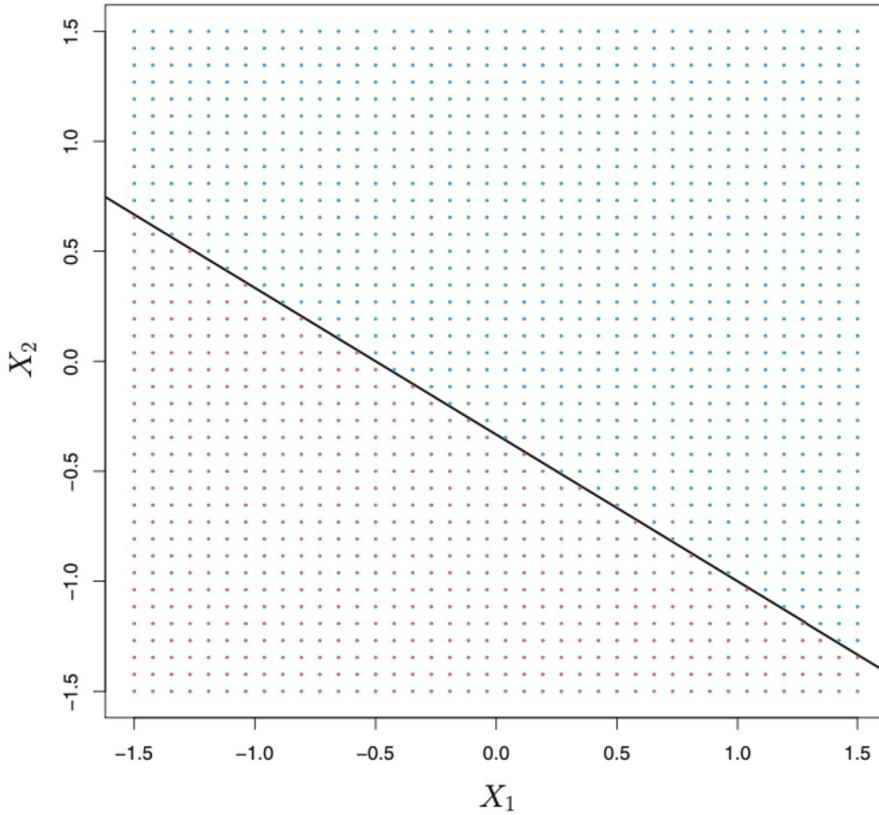
Classification

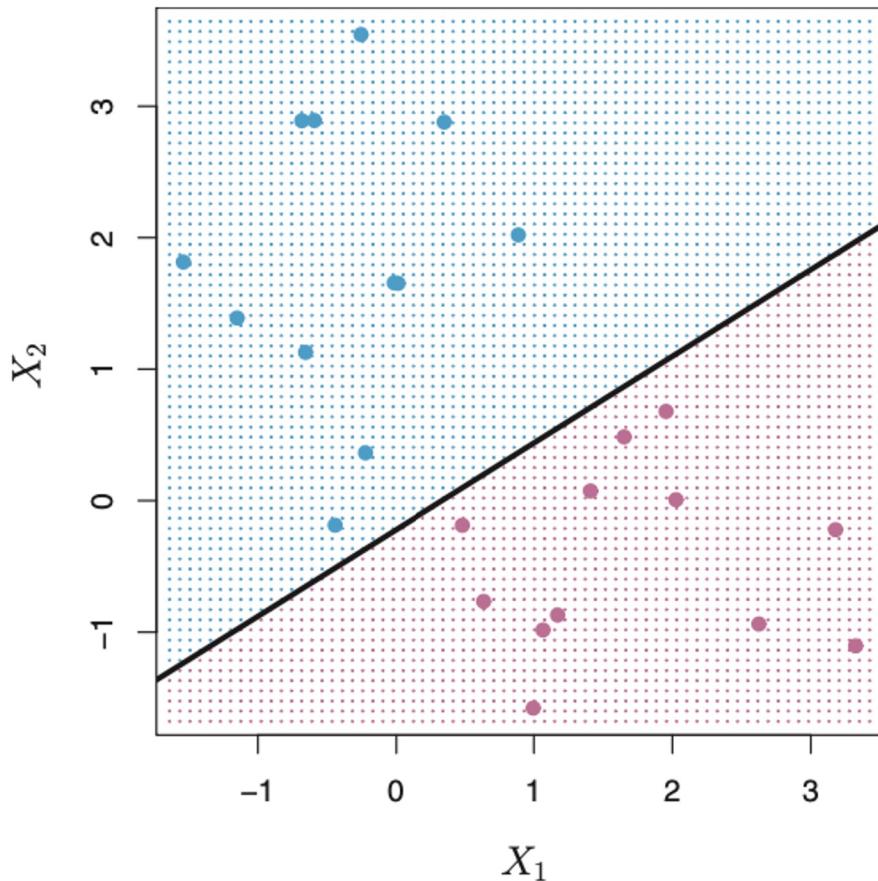
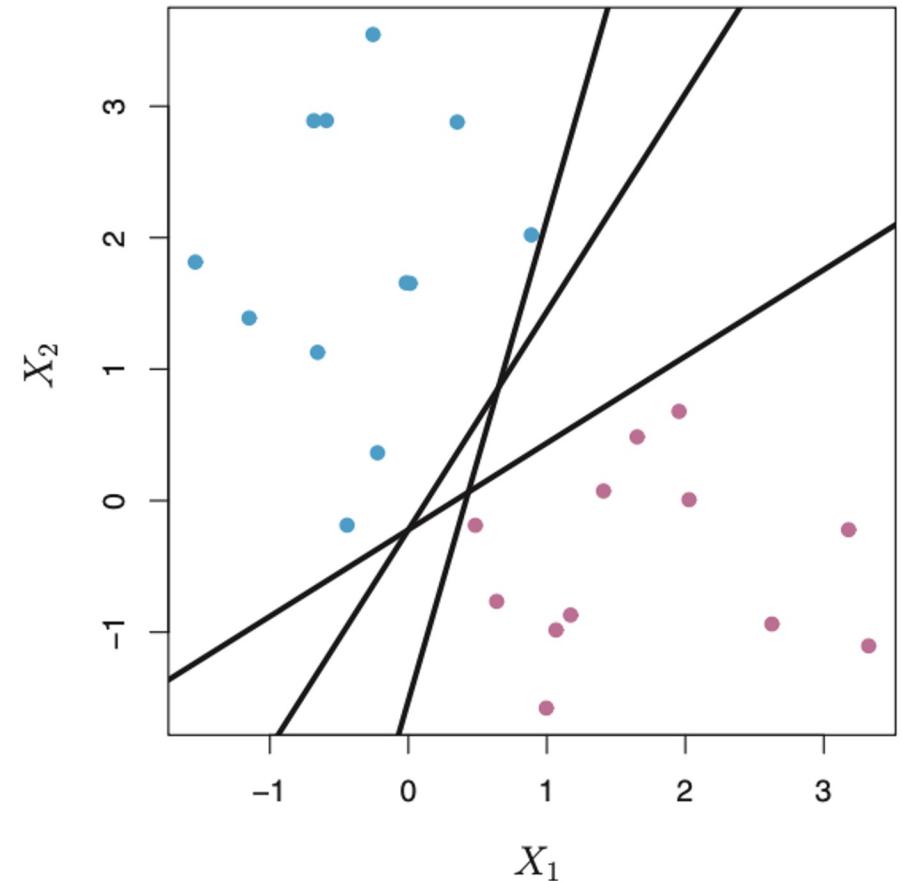
$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p > 0$$

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p < 0$$

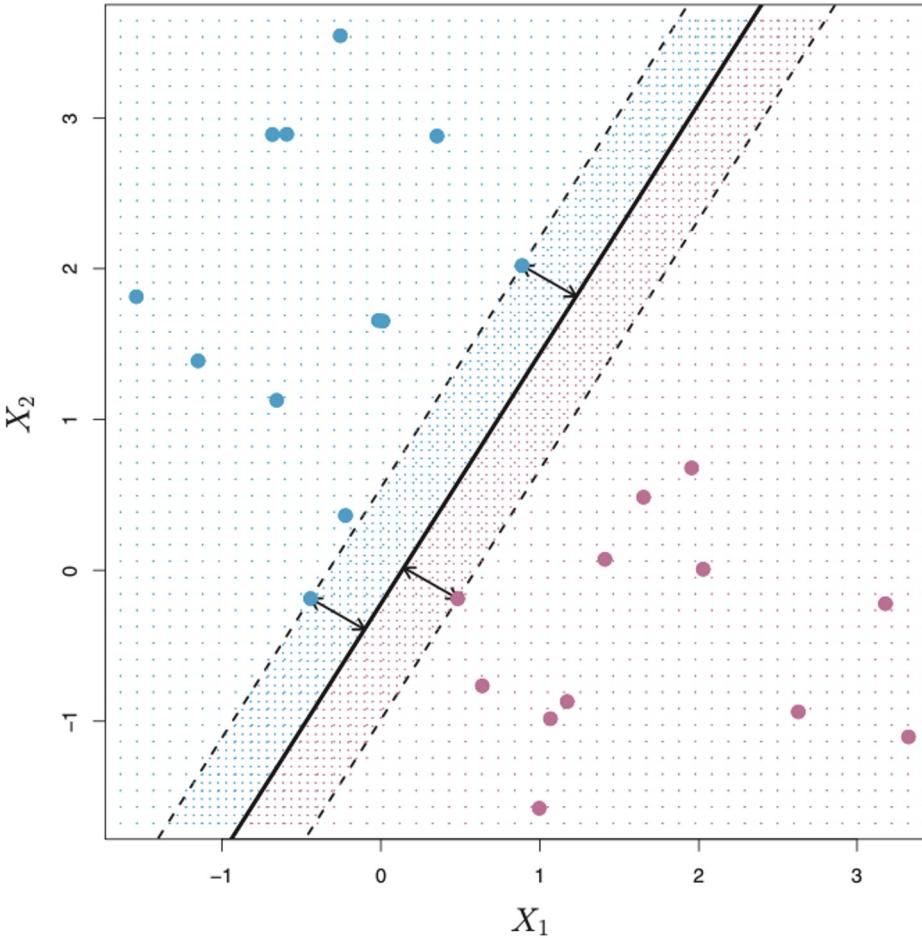
Following images from:







SV Classifier



Support Vector Machine

- Non-linear version of the Support Vector Classifier
- Extension using **Kernels**

$$f(x) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i \langle x, x_i \rangle$$

$$f(x) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i K(x, x_i)$$

$$f(x) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i \langle x, x_i \rangle$$



Kernel function

$$f(x) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i K(x, x_i)$$

$$f(x) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i \langle x, x_i \rangle$$

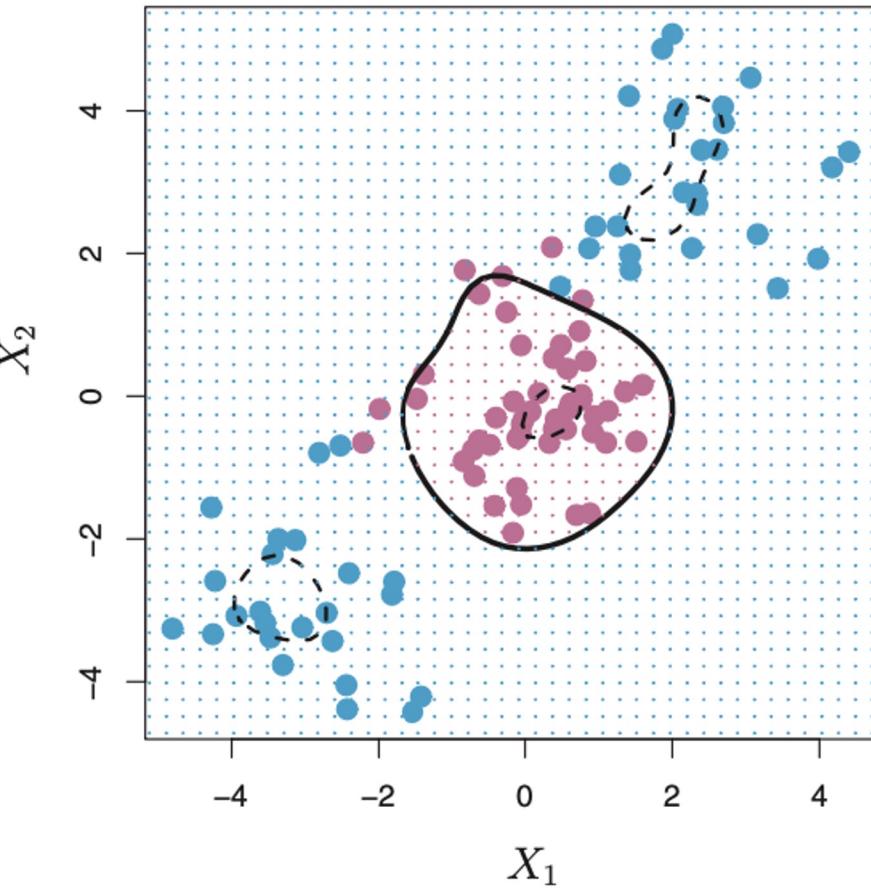
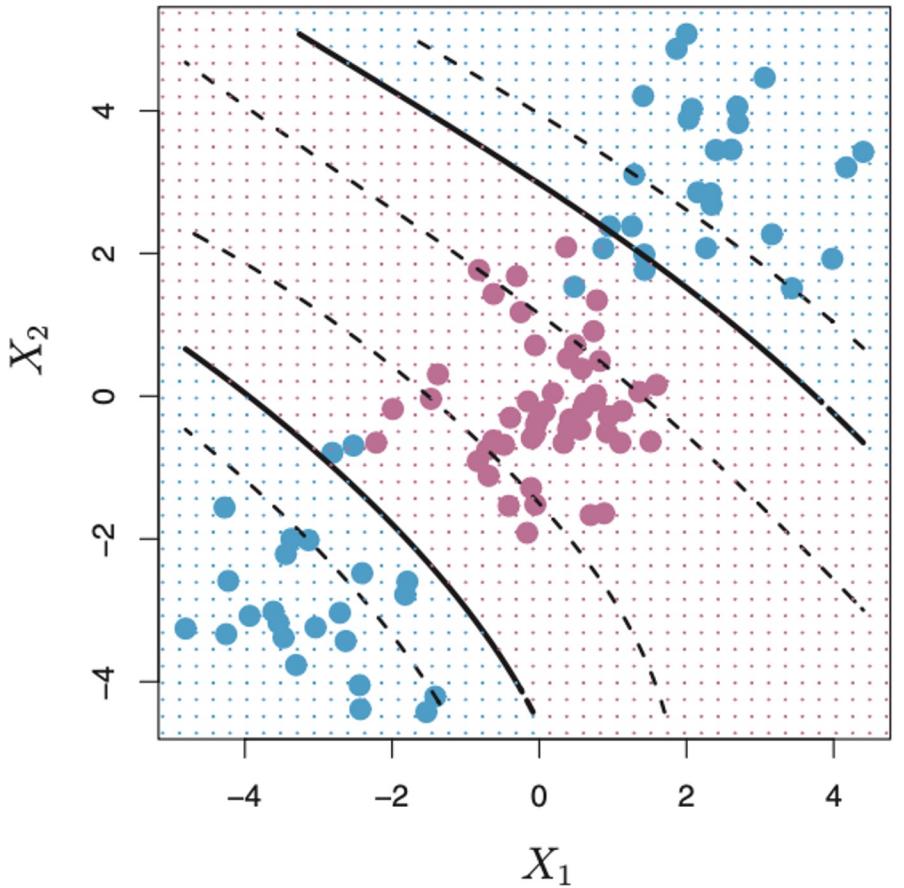
Polynomial Kernel
Non-linear

Kernel function

$$f(x) = \beta_0 + \sum_{i \in \mathcal{S}} \alpha_i K(x, x_i)$$



$$K(x_i, x_{i'}) = (1 + \sum_{j=1}^p x_{ij} x_{i'j})^d.$$



Kernel Trick

- *Actual name*

Kernel Trick

- ***Actual name***
- Attempt to place n-dimensional data into n+1 dimensional space



-5.0

-2.5

0.0

2.5

5.0

x1



-5.0

-2.5

0.0
x1

2.5

5.0

