

Computational Communication Science 2

Week 8 - Lecture

»A Deep Dive into Supervised Machine Learning«

Marthe Möller
Anne Kroon

A.M.Moller@uva.nl, @MartheMoller
a.c.kroon@uva.nl, @annekroon

May, 2022

Digital Society Minor, University of Amsterdam

Today

What is SML?

The principles behind SML

SML step by step

Some classical ML models

Validating models

SML: Strenghts and Challenges

What is SML?

What is SML?

Select all images with cats



Reset

Submit

Yu, J., Ma, X., & Han, T. (2016). Four-Dimensional Usability Investigation of Image CAPTCHA. *arXiv preprint arXiv:1612.01067*.

What is SML?



Read more about this project in: Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., & LeCun, Y.

(2014). OverFeat: Integrated recognition, localization and detection using convolutional networks. *arXiv:1312.6229*

[cs]. Retrieved December 23, 2021, from <http://arxiv.org/abs/1312.6229>

What is SML?

Machine Learning: A process whereby a machine learns how to predict a variable.

What is SML?

Supervised Machine Learning (SML): “A form of machine learning, where we aim to predict a variable that, for a least part of our data is known.”

“The goal of Supervised Machine Learning: estimate a model based on some data, and then use the model to predict the expected outcome for some new cases, for which we do not know the outcome yet.”

Van Atteveldt, W., Trilling, D., & Calderon, C. A. (2022). *Computational analysis of communication*.

Wiley-Blackwell

What is SML?

Machine Learning has a lot of similarities to regression analysis!

The principles behind SML

The principles behind SML

$$y = \text{constant} + b_1 * x_1 + b_2 * x_2$$

x_1 = bark? (0 = no, 1 = yes)

x_2 = tail? (0 = no, 1 = yes)

y = Is this a dog? (0 = definitely no, 1 = definitely yes)

The principles behind SML

$$y = \text{constant} + b_1 * x_1 + b_2 * x_2$$

$$y = 0 + 0.8 * x_1 + 0.2 * x_2$$

$$y = 0 + 0.8 * 1 + 0.2 * 0$$

The principles behind SML

$$y = 0 + 0.8 * 1 + 0.2 * 0$$

$$0.8 = 0 + 0.8 * 1 + 0.2 * 0$$

The principles behind SML

$$0.8 = 0 + 0.8 * 1 + 0.2 * 0$$

Classification: a predictive modeling problem where a class label is predicted for a given example of input data.

The principles behind SML

Machine Learning Lingo	Statistics Lingo
Feature	Independent variable
Label	Dependent variable
Labeled dataset	Dataset with both independent and dependent variables
To train a model	To estimate
Classifier	Model to predict nominal outcomes
To annotate	To (manually) code

Adapted from: Van Atteveldt, Trilling, & Arcilla (2021)

The principles behind SML

Machine Learning: using a (regression) formula to predict a label.

Traditional usage of formulas in CS: to explain

Usage of formulas in ML: to predict

Zooming out

We talked about:

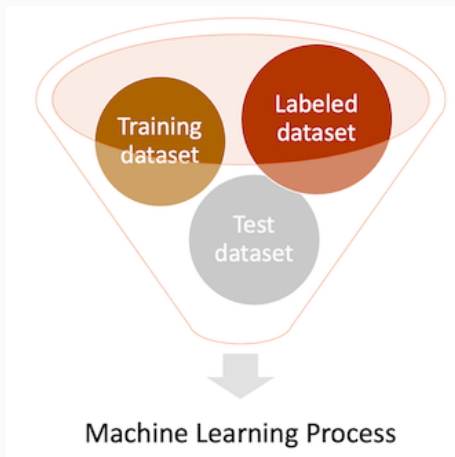
- The principles behind SML

Next, we will talk about:

- The steps of SML

SML step by step

SML step by step



SML step by step

Train
classifier



Test
classifier



Evaluate
classifier

SML step by step



Van Zoonen, W., & Van der Meer, T. G. (2016). Social media research: The application of supervised machine learning in organizational communication research.. *Computers in Human Behavior*, 63, 132–141.

<https://doi.org/10.1016/j.chb.2016.05.028>

Zooming out

We talked about:

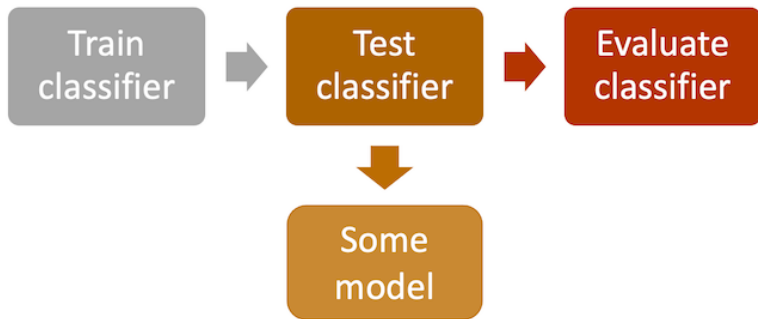
- The principles behind SML
- The steps of SML

Next, we will talk about:

- Some commonly used ML models

Some classical ML models

Regression



Regression

Media literate?

Not at all

Very much

1

2

3

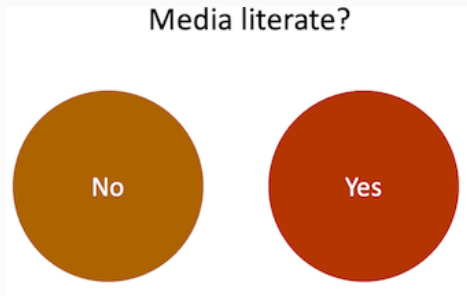
4

5

6

7

Regression



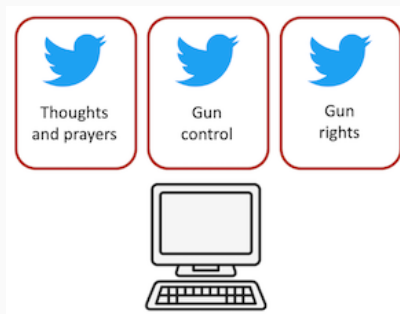
Regression



Zhang, Y., Shah, D., Foley, J., Abhishek, A., Lukito, J., Suk, J., Kim, S. J., Sun, Z., Pevehouse, J., & Garlough, C. (2019). Whose lives matter? mass shootings and social media discourses of sympathy and policy, 2012–2014.

Journal of Computer-Mediated Communication, 24(4), 182–202. <https://doi.org/10.1093/jcmc/zmz009>

Regression

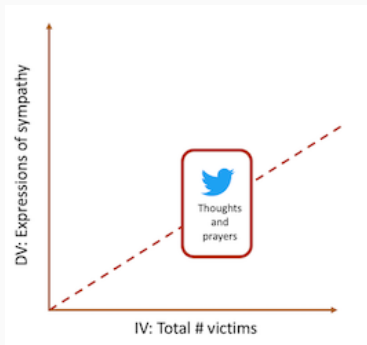


Zhang, Y., Shah, D., Foley, J., Abhishek, A., Lukito, J., Suk, J., Kim, S. J., Sun, Z., Pevehouse, J., & Garlough, C.

(2019). Whose lives matter? mass shootings and social media discourses of sympathy and policy, 2012–2014.

Journal of Computer-Mediated Communication, 24(4), 182–202. <https://doi.org/10.1093/jcmc/zmz009>

Regression

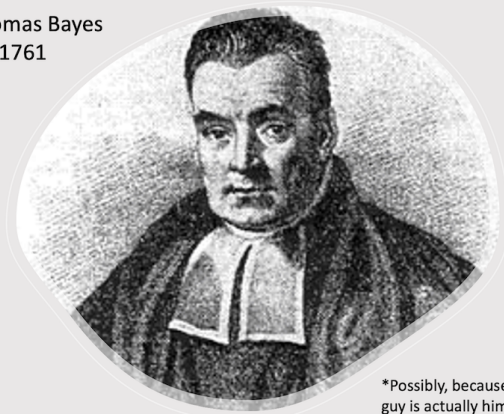


Zhang, Y., Shah, D., Foley, J., Abhishek, A., Lukito, J., Suk, J., Kim, S. J., Sun, Z., Pevehouse, J., & Garlough, C. (2019). Whose lives matter? mass shootings and social media discourses of sympathy and policy, 2012–2014.

Journal of Computer-Mediated Communication, 24(4), 182–202. <https://doi.org/10.1093/jcmc/zmz009>

Naïve Bayes

Possibly* Thomas Bayes
1702 – 1761



*Possibly, because it is unclear if this guy is actually him, but there is no other (claimed) portrait of him.

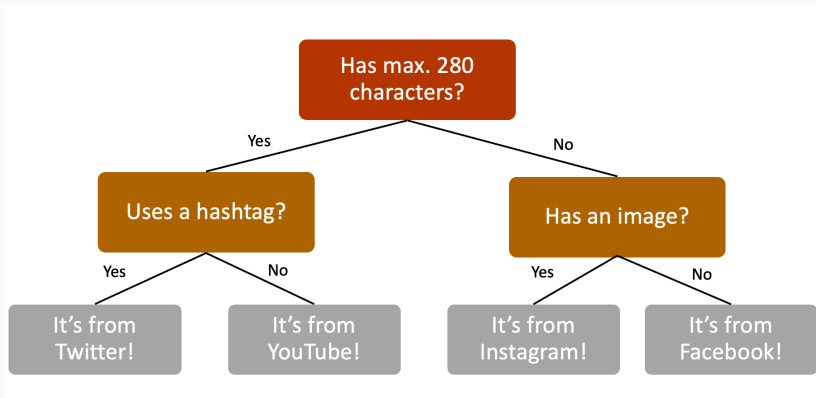
Naïve Bayes

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

Mathematicians' language for: the probability of A is B is the case/present/true.

$$P(\text{label} | \text{features}) = \frac{P(\text{features}|\text{label}) \cdot P(\text{label})}{P(\text{features})}$$

Decision Trees and Random Forests



Decision Trees and Random Forests

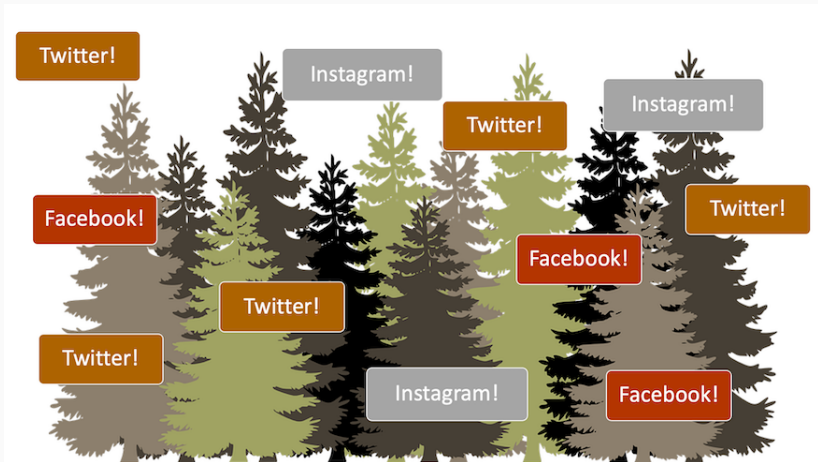
Advantages of decision trees:

- Transparency
- Suitable for non-linear relationships

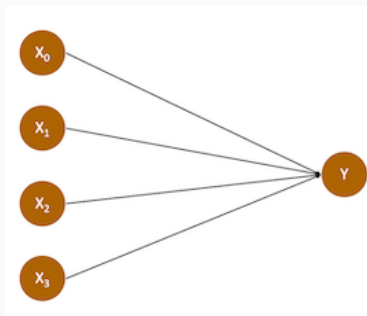
Disadvantages of decision trees:

- Loss of nuance due to yes/no-design
- Cannot correct early mistakes
- Prone to overfitting

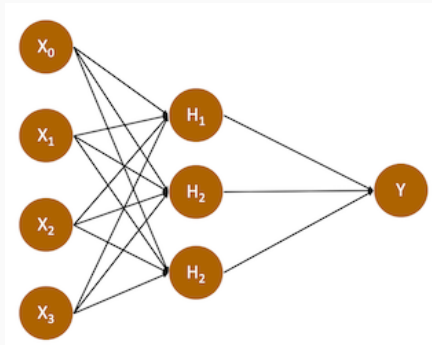
Decision Trees and Random Forests



Neural Networks



Neural Networks



Ha, Y., Park, K., Kim, S. J., Joo, J., & Cha, M. (2021). Automatically detecting image–text mismatch on instagram with deep learning. *Journal of Advertising*, 50(1), 52–62.

<https://doi.org/10.1080/00913367.2020.1843091>

Recap

Many different models available for machine learning.

How do you know what is the best for your case? Try it out and validate!

Zooming out

We talked about:

- The principles behind SML
- The steps of SML
- Some commonly used ML models

Next, we will talk about:

- Validating models

Validating models

Validating Models

Precision quantifies the number of positive class predictions that actually belong to the positive cases.

OR: How much of what we found is actually correct?

Recall quantifies the number of positive class prediction made out of all positive examples in the dataset.

OR: How many of the cases that we wanted to find did we actually find?

Validating Models

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

Validating Models

		Predicted Class	
		Positive	Negative
Actual Class	Positive	150 (TP)	20 (FN)
	Negative	50 (FP)	180 (TN)

Precision is calculated as: $\frac{TP}{TP+FP}$

In our case $\frac{150}{150+50}$ which is 0.75

Recall is calculated as $\frac{TP}{TP+FN}$

In our case $\frac{150}{150+20}$ which is 0.88

Validating Models

Table 2
Relationship classification performance and number of training tweets, random sampling approach.

		100	200	500	1000	2000	3000	4000
Linear support vector machine classifier	AC	0.63	0.65	0.70	0.73	0.80	0.84	0.91
	PC	0.45	0.48	0.59	0.62	0.76	0.80	0.90
	RC	0.38	0.43	0.51	0.59	0.71	0.79	0.86
	AUC	0.41	0.45	0.59	0.61	0.69	0.76	0.85
Naïve Bayes classifier	KA	0.09	0.10	0.39	0.41	0.54	0.65	0.79
	AC	0.63	0.65	0.71	0.75	0.82	0.86	0.91
	PC	0.42	0.46	0.62	0.68	0.81	0.86	0.92
	RC	0.27	0.33	0.47	0.49	0.61	0.69	0.79
Logistic regression classifier	AUC	0.33	0.38	0.60	0.62	0.69	0.77	0.84
	KA	0.08	0.13	0.39	0.40	0.56	0.67	0.78
	AC	0.66	0.67	0.71	0.74	0.79	0.85	0.89
	PC	0.48	0.51	0.63	0.70	0.78	0.89	0.93
	RC	0.04	0.22	0.35	0.39	0.53	0.64	0.73
	AUC	0.08	0.31	0.51	0.55	0.62	0.74	0.82
	KA	0.01	0.09	0.21	0.32	0.48	0.64	0.74

Van Zoonen, W., & Van der Meer, T. G. (2016). Social media research: The application of supervised machine learning in organizational communication research.. *Computers in Human Behavior*, 63, 132–141.

<https://doi.org/10.1016/j.chb.2016.05.028>

Zooming out

We talked about:

- The principles behind SML
- The steps of SML
- Some commonly used ML models
- Validating models

Next, we will talk about:

- Strengths and challenges associated to SML

SML: Strenghts and Challenges

Strengths and Challenges

Strengths:

- Easier to code large datasets
- Enhances replicable research
- Easier to study "natural" human behavior

Disadvantages:

- Resource constraints
- Ethical considerations
- Criticism required (see next slide)

Strengths and Challenges



Belastingdienst koos voor extra controle vooral burgers met lage inkomens

23 november 2021 00:47
Laatste update: 1 dag geleden

761 NU.nl-reacties



Het systeem van de Belastingdienst koos ervoor om de kinderopvangtoeslag vooral bij mensen met een laag inkomen extra te controleren. Dat heeft de fiscus toegegeven in antwoord op vragen van

Trouw en RTL Nieuws.

Fraudejacht bij Toeslagen

Belastingdienst controleerde extra bij lage inkomens in jacht op fraude

22 november 2021 22:59

Trouw

Toeslagen

Belastingdienst ging vooral achter lage inkomens aan

Om toeslagen te controleren op fouten en fraude gebruikte de Belastingdienst een zelflerend algoritme. Dat selecteerde vooral lage inkomens voor controle.

Jan Kleinnijenhuis 22 november 2021

De Belastingdienst heeft jarenlang specifiek burgers met een laag inkomen geselecteerd voor extra controle op

Zooming out I)

We talked about:

- The principles behind SML
- The steps of SML
- Some commonly used ML models
- Validating models
- Strengths and challenges associated to SML

Zooming out II

This week's tutorial:

- Hands-on approach to take a further look into validating models
- Tutorial goals:
 - To get you some experience with the SML process and selecting a model
 - To provide a stepping stone so that you can (independently) advance your machine learning skills