

## **CLASSIFICATION EVALUATION**

Sprachverarbeitung (Vorlesung)

Janis Pagel\*

8 May 2025

\*Based on slides by Nils Reiter

#### Recap

- Machine learning: Let the machine figure out which properties are relevant when
- Feature-based ML: Humans define domain-specific features
- Neural ML: Machine also figures out which features to use
- Train and test data
- Tabular data as input for machine learning systems
- File formats: CSV/TSV
- Why machine learning?
  - Development in NLP/CL over last 30 years
  - Language phenomena in the wild are complex and context-dependent
  - Rule-based systems difficult to develop and maintain



01

# EVALUATION OF MACHINE LEARNING SYSTEMS

#### Introduction

- Starting next weeks: Different machine learning strategies
  - Predictive methods: Given a text, predict some properties of it
- Today: Evaluation
- Goal, in general: Predict (linguistic) categories of text
  - Examples: Parts of speech, syntactic relations, semantic roles, word senses, ...



- For today, we consider the actual ML stuff as a black box
- How exactly do we evaluate? How do we measure how good predictions are?



- For today, we consider the actual ML stuff as a black box
- How exactly do we evaluate? How do we measure how good predictions are?

- Task: Assign a polarity (positive/neutral/negative) to a linguistic expression
- Linguistic expression: sentences, phrases, documents
  - In this example: Documents



- For today, we consider the actual ML stuff as a black box
- How exactly do we evaluate? How do we measure how good predictions are?

- Task: Assign a polarity (positive/neutral/negative) to a linguistic expression
- Linguistic expression: sentences, phrases, documents
  - In this example: Documents
- · Classification task: Instances are sorted into previously known categories



- For today, we consider the actual ML stuff as a black box
- How exactly do we evaluate? How do we measure how good predictions are?

- Task: Assign a polarity (positive/neutral/negative) to a linguistic expression
- Linguistic expression: sentences, phrases, documents
  - In this example: Documents
- · Classification task: Instances are sorted into previously known categories
- Data set: 100 documents that have labels
  - I.e., we know the result to expect



## Annotation Time!

Awesome movie!

Boring as hell

Great start, boring afterwards. Very good acting.



## Experiments





- Goal: Predict the quality on new data
- The program cannot have seen the data, so that it's a realistic test





- Comparison of system output with gold standard
  - "Intrinsic evaluation"
- Two sets of predictions for the items
  - One set from the gold standard
  - One set from the system
- Two aspects to talk about
  - Evaluation metric (how we quantify the performance)
  - Metric interpretation (what we think the metric tells us)



- Comparison of system output with gold standard
  - "Intrinsic evaluation"
- Two sets of predictions for the items
  - One set from the gold standard
  - One set from the system
- Two aspects to talk about
  - Evaluation metric (how we quantify the performance)
  - Metric interpretation (what we think the metric tells us)

- Gold standard: [1, 0, -1, -1]
- System output: [1, -1, 1, 0]
- (positive: 1, neutral: 0, negative: -1)



#### **Extrinsic Evaluation**

- In some cases, reference data for a task doesn't exist or can't be created
- Extrinsic evaluation: Evaluate a downstream application
- Compare performance of downstream application
  - Without your component
  - With your component
- Assumptions
  - Your component helps performance of the downstream application
  - We know how to evaluate the downstream task

![](_page_13_Picture_9.jpeg)

#### **Extrinsic Evaluation**

- In some cases, reference data for a task doesn't exist or can't be created
- Extrinsic evaluation: Evaluate a downstream application
- Compare performance of downstream application
  - Without your component
  - With your component
- Assumptions
  - Your component helps performance of the downstream application
  - We know how to evaluate the downstream task

![](_page_14_Figure_9.jpeg)

![](_page_14_Picture_10.jpeg)

#### Evaluation of Machine Learning Systems

#### Evaluation Metric, Part 1

- Metric Interpretation
- Evaluation Metric, Part 2
- Metric Averages
- Dataset Organization

![](_page_15_Picture_6.jpeg)

#### Accuracy and Error Rate

- Accuracy
  - Percentage of correctly classified instances
  - Example above

• 
$$A = \frac{1}{4} = 0.25 = 25\%$$

• "the higher the better"

![](_page_16_Picture_7.jpeg)

#### Accuracy and Error Rate

- Accuracy
  - Percentage of correctly classified instances
  - Example above

• 
$$A = \frac{1}{4} = 0.25 = 25\%$$

- "the higher the better"
- Error Rate
  - Percentage of *incorrectly* classified instances
  - Example above

• 
$$E = \frac{3}{4} = 0.75 = 75\%$$

• "the lower the better"

![](_page_17_Picture_12.jpeg)

#### Accuracy and Error Rate

- Accuracy
  - Percentage of correctly classified instances
  - Example above

• 
$$A = \frac{1}{4} = 0.25 = 25\%$$

- "the higher the better"
- Error Rate
  - Percentage of *incorrectly* classified instances
  - Example above

• 
$$E = \frac{3}{4} = 0.75 = 75\%$$

- "the lower the better"
- A + E = 1, E = 1 A and A = 1 E

![](_page_18_Picture_13.jpeg)

![](_page_19_Picture_3.jpeg)

• G = [1, 0, 1], S = [0, 0, 1]  
• A = ?, E = ?  
• A = 
$$\frac{2}{3}$$
 = 0.66, E =  $\frac{1}{3}$  = 0.33

![](_page_20_Picture_3.jpeg)

![](_page_21_Picture_3.jpeg)

• G = [1, 0, 1], S = [0, 0, 1]  
• A = ?, E = ?  
• A = 
$$\frac{2}{3}$$
 = 0.66, E =  $\frac{1}{3}$  = 0.33  
• G = ['f", "m", "u", "m", "f"], S = ['m", "f", "u", "m", "f"]  
• A = ?, E = ?  
• A =  $\frac{3}{5}$  = 0.6, E =  $\frac{2}{5}$  = 0.4

![](_page_22_Picture_3.jpeg)

Examples

• G = [1, 0, 1], S = [0, 0, 1] • A = ?, E = ? • A =  $\frac{2}{3}$  = 0.66, E =  $\frac{1}{3}$  = 0.33 • G = ["f", "m", "u", "m", "f"], S = ["m", "f", "u", "m", "f"] • A = ?, E = ? • A =  $\frac{3}{5}$  = 0.6, E =  $\frac{2}{5}$  = 0.4

(We don't need the original data for evaluation, we are just comparing gold standard classes with system output. We don't even need to know what the classes represent.)

![](_page_23_Picture_4.jpeg)

#### Evaluation of Machine Learning Systems

- Evaluation Metric, Part 1
- Metric Interpretation
- Evaluation Metric, Part 2
- Metric Averages
- Dataset Organization

![](_page_24_Picture_6.jpeg)

How good are 60% accuracy?

![](_page_25_Picture_1.jpeg)

- Something to compare with
- Justification for investing research time
- Predecessor system
  - E.g., the one from last year
- Competing system
  - E.g., the one from Düsseldorf University
- Very simple system
  - E.g., a single feature decides everything
- Dummy system
  - E.g., if we make random decisions
  - Most common baseline
- · It's allowed to specify multiple baselines

![](_page_26_Picture_13.jpeg)

UNIVERSITÄT ZU KÖLN

<ul> <li>Something to compare with</li> </ul>	System	Accuracy
<ul> <li>Justification for investing research time</li> </ul>	Model 1	56
Predecessor system	Model 2	53
• E.g., the one from last year	Model 3	58
Competing system	Baseline 1	33
<ul> <li>E.g., the one from Düsseldorf University</li> </ul>	Baseline 2	45
<ul> <li>Very simple system</li> <li>E.g., a single feature decides everything</li> </ul>	Table: Results table in publicatio	
Dummy system		
<ul> <li>E.g., if we make random decisions</li> </ul>		

- Most common baseline
- It's allowed to specify multiple baselines

A simple solution to the problem

- How well can the task be solved without investing (a lot of) time and work?
- What is a simple solution, and how well does it solve the problem?

![](_page_28_Picture_4.jpeg)

A simple solution to the problem

- How well can the task be solved without investing (a lot of) time and work?
- What is a simple solution, and how well does it solve the problem?
- Baselines are used for comparison in experiments
- 'Real' algorithms should be able to beat the baseline, i.e., achieve higher accuracy
- · Baselines have obvious shortcomings, are not expected to work every time
  - Although, sometimes they work surprisingly well

![](_page_29_Picture_8.jpeg)

#### **Group Exercises**

What are reasonable baselines for these tasks?

- Detecting nouns in German texts
- Detecting sentence boundaries
- Detecting fake news
- Detecting the gender of dramatic characters (18-19th century)
- Predict the pos tag of the word after a determiner
- Given a corpus consisting of 'the Universal Declaration of Human Rights', 'Lord of the Rings' and the minutes of the European Parliament. Predict the origin of a random sentence.

![](_page_30_Picture_9.jpeg)

### **Majority Baseline**

- Select the most frequent category
- Works well in un-even data distributions
  - I.e., if one category is more frequent than the others
- Can be hard to beat
  - E.g. word sense disambiguation

![](_page_31_Picture_6.jpeg)

## **Random Baseline**

- Randomly select a category
- Works well in even distributions
  - I.e., if all categories are equally frequent

![](_page_32_Picture_4.jpeg)

#### Evaluation of Machine Learning Systems

- Evaluation Metric, Part 1
- Metric Interpretation
- Evaluation Metric, Part 2
- Metric Averages
- Dataset Organization

![](_page_33_Picture_6.jpeg)

#### **Per Class Evaluation**

- Accuracy gives us an overall score
- But we want to know more details:
  - Some classes are more important for applications
  - Error analysis!
- We want to evaluate per class (i.e., per polarity)

![](_page_34_Picture_6.jpeg)

Different Kinds of Errors

Polarity	Document
positive neutral negative	Awesome movie! Great start, boring afterwards. Very good acting. Boring as hell

Table: Gold Standard

![](_page_35_Picture_4.jpeg)

Different Kinds of Errors

Polarity	Document
positive neutral negative	Awesome movie! Great start, boring afterwards. Very good acting. Boring as hell

Table: Gold Standard

Variant	Output
GS	1, 0, -1, 1, 1, 0, -1, 1
Model 1 Model 2	1, 0, -1, 1, 1, 0, <mark>1</mark> , 1 1, 0, -1, 1, - <mark>1</mark> , 0, -1, 1

![](_page_36_Picture_5.jpeg)

**Different Kinds of Errors** 

![](_page_37_Picture_2.jpeg)

Figure: Visual representation of errors, focussing on -1 class

![](_page_37_Picture_4.jpeg)

**Different Kinds of Errors** 

![](_page_38_Figure_2.jpeg)

Figure: Visual representation of errors, focussing on -1 class

![](_page_38_Picture_4.jpeg)

**Different Kinds of Errors** 

![](_page_39_Figure_2.jpeg)

Figure: Visual representation of errors, focussing on -1 class

![](_page_39_Picture_4.jpeg)

## **Different Kinds of Errors**

![](_page_40_Figure_1.jpeg)

![](_page_40_Picture_2.jpeg)

#### **Different Kinds of Errors**

![](_page_41_Figure_1.jpeg)

true positive (tp) Correctly classified as target category
true negative (tn) Correctly classified as not target category

![](_page_41_Picture_3.jpeg)

#### **Different Kinds of Errors**

![](_page_42_Figure_1.jpeg)

true positive (tp) Correctly classified as target category
true negative (tn) Correctly classified as not target category
false positive (fp) Incorrectly classified as target category
false negative (fn) Incorrectly classified as not target category

![](_page_42_Picture_3.jpeg)

## Accuracy, revisited

Accuracy: Percentage of correctly classified instances

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

![](_page_43_Picture_3.jpeg)

#### Accuracy, revisited

Accuracy: Percentage of correctly classified instances

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

Error rate: Percentage of incorrectly classified instances

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

![](_page_44_Picture_5.jpeg)

Given the documents that the system marked as -1, how many of those are really -1?

![](_page_45_Picture_2.jpeg)

Given the documents that the system marked as -1, how many of those are really -1?

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

![](_page_46_Picture_3.jpeg)

Given the documents that the system marked as -1, how many of those are really -1?

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

How many of the -1 documents did the system find?

![](_page_47_Picture_4.jpeg)

Given the documents that the system marked as -1, how many of those are really -1?

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

How many of the -1 documents did the system find?

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

![](_page_48_Picture_5.jpeg)

• Enumerator: tp

![](_page_49_Picture_2.jpeg)

- Enumerator: tp
- Precision
  - Denominator: tp + fp
  - Number of things that the system labelled as target category (correct and incorrect)
- Recall
  - Denominator: tp + fn
  - Number of things that the gold standard contained as target category (what the system should have found)

![](_page_50_Picture_8.jpeg)

Importance/Weighting

- Weighting between P and R is application-dependent (and difficult to decide!)
- Guiding question: Which kind of error is more severe?

![](_page_51_Picture_4.jpeg)

Importance/Weighting

- Weighting between P and R is application-dependent (and difficult to decide!)
- Guiding question: Which kind of error is more severe?
- If findings are inspected by humans
  - Precision errors are easy to spot, but recall errors cannot be detected
  - But: humans tend to trust computers

![](_page_52_Picture_7.jpeg)

Importance/Weighting

- Weighting between P and R is application-dependent (and difficult to decide!)
- Guiding question: Which kind of error is more severe?
- If findings are inspected by humans
  - · Precision errors are easy to spot, but recall errors cannot be detected
  - But: humans tend to trust computers
- Severity of consequences

![](_page_53_Picture_8.jpeg)

Importance/Weighting

- Weighting between P and R is application-dependent (and difficult to decide!)
- Guiding question: Which kind of error is more severe?
- If findings are inspected by humans
  - · Precision errors are easy to spot, but recall errors cannot be detected
  - But: humans tend to trust computers
- Severity of consequences

#### Example (Test performance in a pandemic)

- Individual health: Mistakenly being in quarantine is a severe limitation, and might have economic consequences
- Public health: Find more infections, even if it means a few people are mistakenly put in quarantine

![](_page_54_Picture_11.jpeg)

### **F-Score**

- Sometimes, it is convenient to combine precision and recall into a single number
- F-Score is common way to do that (it's a fancy way of averaging)
  - $\beta$  can be used to weight precision and recall differently
  - $\beta = 1$  means equal weighting,  $\beta = 2$  weighs recall two times as high as precision,  $\beta = 0.5$  weighs recall two times as low as precision
- F-Measure corresponds to the harmonic mean

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

![](_page_55_Picture_7.jpeg)

#### **F-Score**

- Sometimes, it is convenient to combine precision and recall into a single number
- F-Score is common way to do that (it's a fancy way of averaging)
  - $\beta$  can be used to weight precision and recall differently
  - $\beta = 1$  means equal weighting,  $\beta = 2$  weighs recall two times as high as precision,  $\beta = 0.5$  weighs recall two times as low as precision
- F-Measure corresponds to the harmonic mean

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

- Most commonly chosen value for  $\beta$  is 1
- The equation simplifies to:

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

![](_page_56_Picture_10.jpeg)

#### Evaluation of Machine Learning Systems

- Evaluation Metric, Part 1
- Metric Interpretation
- Evaluation Metric, Part 2
- Metric Averages
- Dataset Organization

![](_page_57_Picture_6.jpeg)

## **Results in Scientific Papers**

System	Class	Precision	Recall
	-1	45	75
Model 1	0	54	61
	1	78	12
	Macro Average	59	49
	Micro Average	55	56
	-1	0	0
	0	100	0
Baseline 1	1	0	0
	Macro Average	33	0
	Micro Average	75	0

Table: Example table with results

![](_page_58_Picture_3.jpeg)

#### Micro- and Macro-Average

• Macro-Average: Arithmetic mean

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

• Micro-Average: Weighted arithmetic mean

$$\overline{x} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$$

• Takes into account how frequent categories are

![](_page_59_Picture_6.jpeg)

#### Micro- and Macro-Average

• Macro-Average: Arithmetic mean

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

• Micro-Average: Weighted arithmetic mean

$$\overline{x} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$$

• Takes into account how frequent categories are

Class	Freq. (= $w$ )	Ρ	R
А	7	50	90
В	1	80	10
С	2	90	20
Macro Average		73	40
Micro Average		61	68

![](_page_60_Picture_7.jpeg)

#### Evaluation of Machine Learning Systems

- Evaluation Metric, Part 1
- Metric Interpretation
- Evaluation Metric, Part 2
- Metric Averages
- Dataset Organization

![](_page_61_Picture_6.jpeg)

• What if a chosen test split results in higher scores because items in split are easier than others by chance?

![](_page_62_Picture_2.jpeg)

- What if a chosen test split results in higher scores because items in split are easier than others by chance?
- Solution: Cross validation
  - Make multiple splits of data so that every part of data is tested once
  - Number of cross validation splits are called folds

![](_page_63_Picture_5.jpeg)

- What if a chosen test split results in higher scores because items in split are easier than others by chance?
- Solution: Cross validation
  - Make multiple splits of data so that every part of data is tested once
  - Number of cross validation splits are called folds

![](_page_64_Picture_5.jpeg)

- What if a chosen test split results in higher scores because items in split are easier than others by chance?
- Solution: Cross validation
  - Make multiple splits of data so that every part of data is tested once
  - Number of cross validation splits are called folds

![](_page_65_Figure_5.jpeg)

![](_page_65_Picture_7.jpeg)

- What if a chosen test split results in higher scores because items in split are easier than others by chance?
- Solution: Cross validation
  - Make multiple splits of data so that every part of data is tested once
  - Number of cross validation splits are called folds

![](_page_66_Figure_5.jpeg)

#### Randomness

- Some test options or algorithms involve random numbers
  - E.g., cross validation
- Results could be unrealistically good, by chance

![](_page_67_Picture_4.jpeg)

#### Randomness

- Some test options or algorithms involve random numbers
  - E.g., cross validation
- Results could be unrealistically good, by chance
- Simple solution: Run the experiments repeatedly (e.g., 1000 times)

![](_page_68_Picture_5.jpeg)

02

# **SUMMARY**

## Summary

- Evaluation of ML models is important
  - · Because we don't know in advance what works and what does not
- Two components
  - Comparison to a baseline
    - Previous or dummy model
  - Calculation of precision/recall
    - Precision: How many of those marked as category X by the model are truly category X?
    - Recall: How many of those that are category X has the model marked as X?
  - Training/test split or cross validation

![](_page_70_Picture_10.jpeg)