

Language Processing for Different Domains and Genres: Machine Learning Introduction

Caroline Sporleder, Ines Rehbein

Universität des Saarlandes

Wintersemester 2009/10

5.11.2009

Goal:

develop computer programs that automatically improve with experience by learning from representative input (and output) data

Motivation:

- for many problems, the best way of computing the correct output from the input is not known.
- manually determining input output rules by (informed) trial and error is time consuming and typically results in low coverage (but high precision)

Example:

predicting the plural form of a German noun

Example: German Plural Formation

Nine possibilities:

- 1 no ending, no umlaut: *das Zimmer* - *die Zimmer* (rooms)
- 2 no ending, but umlaut: *der Faden* - *die Fäden* (thread)
- 3 -e: *der Hund* - *die Hunde* (dogs)
- 4 -e plus umlaut: *der Stuhl* - *die Stühle* (chairs)
- 5 -er: *das Kind* - *die Kinder* (children)
- 6 -er plus umlaut: *das Lamm* - *die Lämmer* (lambs)
- 7 -n: *die Straße* - *die Straßen* (streets)
- 8 -en: *die Bank* - *die Banken* (banks)
- 9 -s: *das Trio* - *die Trios* (trios)

Hypothesis 1: ending is determined by noun's grammatical gender

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n)

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*
- *der Hund* (m)

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*
- *der Hund* (m) \Rightarrow *die Hunde*

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*
- *der Hund* (m) \Rightarrow *die Hunde*
- *die Bank* (f)

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*
- *der Hund* (m) \Rightarrow *die Hunde*
- *die Bank* (f) \Rightarrow *die Banken*

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*
- *der Hund* (m) \Rightarrow *die Hunde*
- *die Bank* (f) \Rightarrow *die Banken*
- *das Zimmer* (n)

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*
- *der Hund* (m) \Rightarrow *die Hunde*
- *die Bank* (f) \Rightarrow *die Banken*
- *das Zimmer* (n) \Rightarrow *die Zimmerer*

Hypothesis 1: ending is determined by noun's grammatical gender

Rule 1:

- masculine \Rightarrow -e
- neuter \Rightarrow -er
- feminine \Rightarrow -en

Applying Rule 1:

- *das Kind* (n) \Rightarrow *die Kinder*
- *der Hund* (m) \Rightarrow *die Hunde*
- *die Bank* (f) \Rightarrow *die Banken*
- *das Zimmer* (n) \Rightarrow **die Zimmerer* (*die Zimmer*)

Hypothesis 2: morpho-phonological form also influences ending

Hypothesis 2: morpho-phonological form also influences ending

Rule 2: don't add ending if noun already ends in -e, -en, or -er

Hypothesis 2: morpho-phonological form also influences ending

Rule 2: don't add ending if noun already ends in -e, -en, or -er

Applying Rules 1 and 2:

Hypothesis 2: morpho-phonological form also influences ending

Rule 2: don't add ending if noun already ends in -e, -en, or -er

Applying Rules 1 and 2:

- *das Zimmer* (n)

Hypothesis 2: morpho-phonological form also influences ending

Rule 2: don't add ending if noun already ends in -e, -en, or -er

Applying Rules 1 and 2:

- *das Zimmer* (n) \Rightarrow die Zimmer

Hypothesis 2: morpho-phonological form also influences ending

Rule 2: don't add ending if noun already ends in -e, -en, or -er

Applying Rules 1 and 2:

- *das Zimmer* (n) \Rightarrow die Zimmer
- *die Ampel* (f)

Hypothesis 2: morpho-phonological form also influences ending

Rule 2: don't add ending if noun already ends in -e, -en, or -er

Applying Rules 1 and 2:

- *das Zimmer* (n) \Rightarrow *die Zimmer*
- *die Ampel* (f) \Rightarrow *die Ampelen*

Hypothesis 2: morpho-phonological form also influences ending

Rule 2: don't add ending if noun already ends in -e, -en, or -er

Applying Rules 1 and 2:

- *das Zimmer* (n) ⇒ die Zimmer
- *die Ampel* (f) ⇒ **die Ampelen* (*die Ampeln*)

Rule 3: -e and -en become \emptyset and -n if the last syllable of the singular contains a schwa.

Rule 3: -e and -en become \emptyset and -n if the last syllable of the singular contains a schwa.

Applying Rules 1, 2 and 3:

Rule 3: -e and -en become \emptyset and -n if the last syllable of the singular contains a schwa.

Applying Rules 1, 2 and 3:

- *die Ampel* (f)

Rule 3: -e and -en become \emptyset and -n if the last syllable of the singular contains a schwa.

Applying Rules 1, 2 and 3:

- *die Ampel* (f) \Rightarrow *die Ampeln*

Rule 3: -e and -en become \emptyset and -n if the last syllable of the singular contains a schwa.

Applying Rules 1, 2 and 3:

- *die Ampel* (f) \Rightarrow *die Ampeln*
- *der Nachbar* (m)

Rule 3: -e and -en become \emptyset and -n if the last syllable of the singular contains a schwa.

Applying Rules 1, 2 and 3:

- *die Ampel* (f) \Rightarrow *die Ampeln*
- *der Nachbar* (m) \Rightarrow *die Nachbare*

Rule 3: -e and -en become \emptyset and -n if the last syllable of the singular contains a schwa.

Applying Rules 1, 2 and 3:

- *die Ampel* (f) \Rightarrow *die Ampeln*
- *der Nachbar* (m) \Rightarrow **die Nachbare* (*die Nachbarn*)

Machine Learning German Plural Formation

Learn from input-output pairs:

- Zimmer, Zimmer
- Faden, Fäden
- Hund, Hunde
- Stuhl, Stühle
- Kind, Kinder
- Lamm, Lämmer
- Straße, Straßen
- Bank, Banken
- Trio, Trios
- Ampel, Ampeln
- Nachbar, Nachbarn
- Maus, Mäuse

⇒ **Input** is typically represented as a **feature vector**.

Machine Learning German Plural Formation (2)

What information needs to be represented for the task to be learnable (i.e., which **features** need to be modelled)?

Machine Learning German Plural Formation (2)

What information needs to be represented for the task to be learnable (i.e., which **features** need to be modelled)?

- 1 Zimmer: <n>, Zimmer
- 2 Faden: <m>, Fäden
- 3 Hund: <m>, Hunde
- 4 Stuhl: < m>, Stühle
- 5 Kind: <n>, Kinder
- 6 Lamm: <n>, Lämmer
- 7 Straße: <f>, Straßen
- 8 Bank: <f>, Banken
- 9 Trio: <n>, Trios
- 10 Ampel: <f>, Ampeln
- 11 Nachbar: <m>, Nachbarn
- 12 Maus: <f>, Mäuse

Machine Learning German Plural Formation (2)

What information needs to be represented for the task to be learnable (i.e., which **features** need to be modelled)?

- 1 Zimmer: $\langle n, \text{er} \rangle$, Zimmer
- 2 Faden: $\langle m, \text{en} \rangle$, Fäden
- 3 Hund: $\langle m, \emptyset \rangle$, Hunde
- 4 Stuhl: $\langle m, \emptyset \rangle$, Stühle
- 5 Kind: $\langle n, \emptyset \rangle$, Kinder
- 6 Lamm: $\langle n, \emptyset \rangle$, Lämmer
- 7 Straße: $\langle f, \text{e} \rangle$, Straßen
- 8 Bank: $\langle f, \emptyset \rangle$, Banken
- 9 Trio: $\langle n, \emptyset \rangle$, Trios
- 10 Ampel: $\langle f, \emptyset \rangle$, Ampeln
- 11 Nachbar: $\langle m, \emptyset \rangle$, Nachbarn
- 12 Maus: $\langle f, \emptyset \rangle$, Mäuse

Machine Learning German Plural Formation (2)

What information needs to be represented for the task to be learnable (i.e., which **features** need to be modelled)?

- 1 Zimmer: $\langle n, er, a\text{-schwa} \rangle$, Zimmer
- 2 Faden: $\langle m, en, e\text{-schwa} \rangle$, Fäden
- 3 Hund: $\langle m, \emptyset, no\text{ schwa} \rangle$, Hunde
- 4 Stuhl: $\langle m, \emptyset, no\text{ schwa} \rangle$, Stühle
- 5 Kind: $\langle n, \emptyset, no\text{ schwa} \rangle$, Kinder
- 6 Lamm: $\langle n, \emptyset, no\text{ schwa} \rangle$, Lämmer
- 7 Straße: $\langle f, e, e\text{-schwa} \rangle$, Straßen
- 8 Bank: $\langle f, \emptyset, no\text{ schwa} \rangle$, Banken
- 9 Trio: $\langle n, \emptyset, no\text{ schwa} \rangle$, Trios
- 10 Ampel: $\langle f, \emptyset, e\text{-schwa} \rangle$, Ampeln
- 11 Nachbar: $\langle m, \emptyset, no\text{ schwa} \rangle$, Nachbarn
- 12 Maus: $\langle f, \emptyset, no\text{ schwa} \rangle$, Mäuse

instance: one input-output pair, where the input is a feature vector

Hund, m, \emptyset , no schwa, Hunde

label or class to be predicted: the output value

Hunde

label can be nominal, numeric, binary etc.

features: types of information encoded in the input

singular form, gender, ending of singular, schwa information

feature values: the values the features assume (for a given instance)

Hund, m, \emptyset , no schwa

values can be nominal, numeric, binary etc.

training set: set of manually labelled instances from which target function (mapping from input to output) is learnt

test set: set of instances to which the trained machine learner is applied, test set labels are used to compute performance of classifier (labels are not known to the classifier)

development set: set of instances used to choose the best parameter settings (typically those that optimise performance on the development set)

Supervised Machine Learning:

system learns target function (mapping from input to output) from a labelled training set (decision tree learners, Naive Bayes, k-NN etc.)

Unsupervised Machine Learning:

no labelled training set, system searches for best model to account for unlabelled test set (clustering etc.)

Semi-Supervised Machine Learning:

combination of supervised and unsupervised; training set consists of a small labelled seed set and a large unlabelled set

Machine Learning Spaces

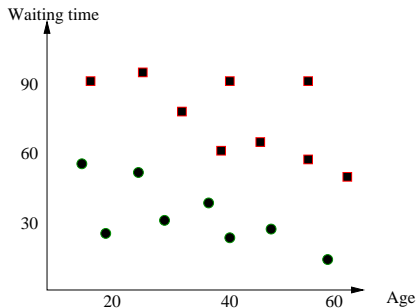
Instances are points (vectors) in n -dimensional space, where n is the number of features

Example: predict customer satisfaction (*happy* vs. *not happy*) of the clients of a call centre from (i) their age and (ii) the number of minutes they were kept in the waiting loop.

Machine Learning Spaces

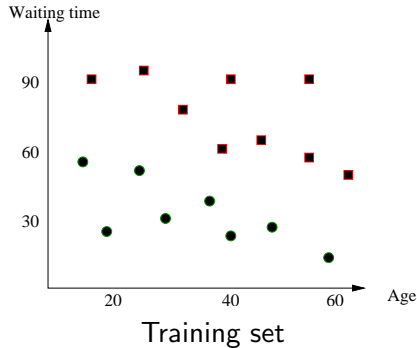
Instances are points (vectors) in n -dimensional space, where n is the number of features

Example: predict customer satisfaction (*happy* vs. *not happy*) of the clients of a call centre from (i) their age and (ii) the number of minutes they were kept in the waiting loop.



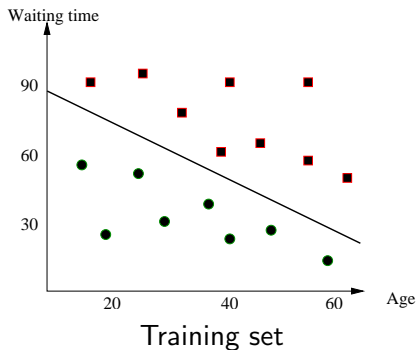
Supervised Machine Learning (1)

Classifier uses training set to try to find the correct decision boundary between the classes.



Supervised Machine Learning (1)

Classifier uses training set to try to find the correct decision boundary between the classes.

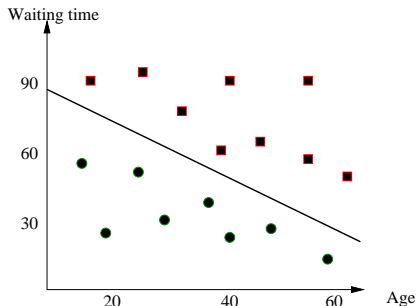


Separability

Whether and how two sets can be separated depends on

- the machine learning algorithm (linear or non-linear)
- the dimensionality of the space (the number of features) and the suitability of the features

Data sets which can be separated by a line (2-d), plane (3-d), or hyperplane (higher dimensional space) are called **linearly separable**.



All machine learners generalise (opposite rote learning)

If they didn't generalise, they wouldn't be able to

All machine learners generalise (opposite **rote learning**)

If they didn't generalise, they wouldn't be able to

- label unseen instances
- ignore outliers (e.g., mislabelled instances in the training set)
⇒ overfitting

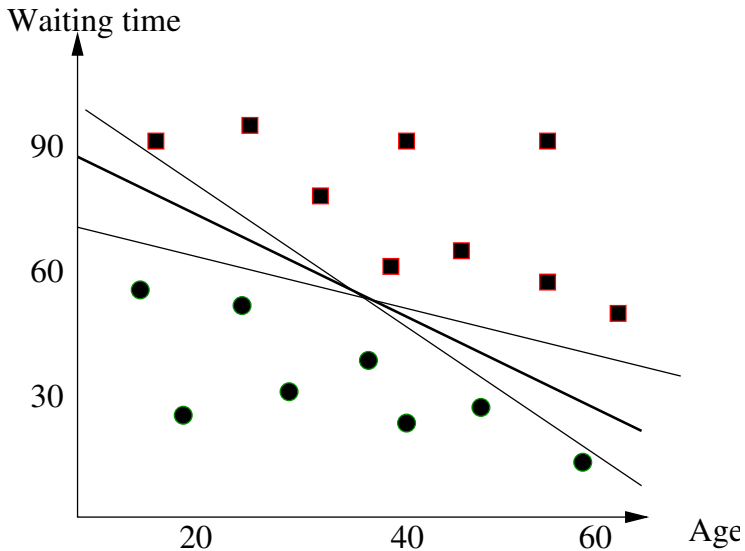
All machine learners generalise (opposite **rote learning**)

If they didn't generalise, they wouldn't be able to

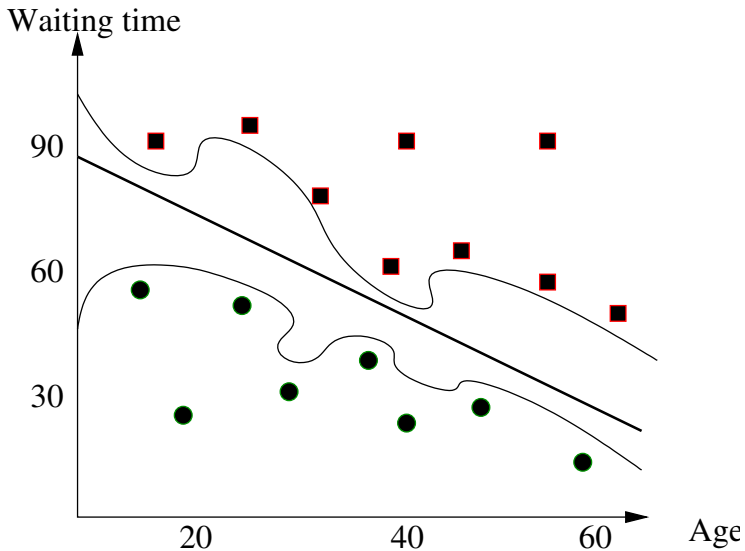
- label unseen instances
- ignore outliers (e.g., mislabelled instances in the training set)
⇒ overfitting

However: different machine learners generalise in different ways
(**inductive bias**)

Generalisation (2)



Generalisation (2)

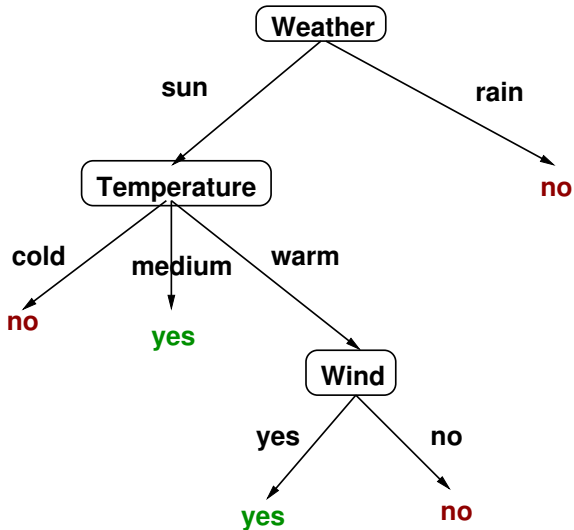


Example: Decision Tree Learning (1)

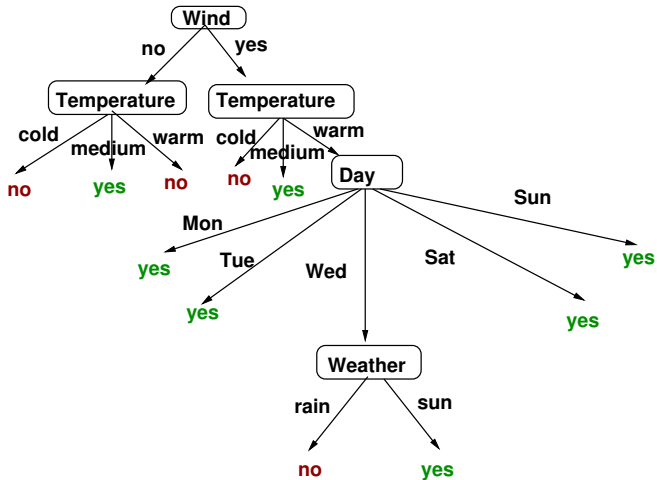
Training Data for “Play Tennis” Task

day	weather	temperature	wind	play_tennis?
Tues	sun	warm	no	no
Sun	rain	cold	no	no
Mon	sun	medium	no	yes
Wed	rain	warm	yes	no
Sat	sun	warm	yes	yes
Wed	sun	warm	yes	yes
Mon	sun	warm	yes	yes
Sun	sun	warm	yes	yes

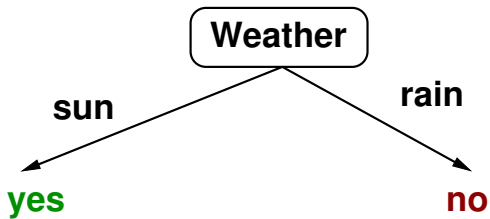
Example: Alternative Decision Trees (2)



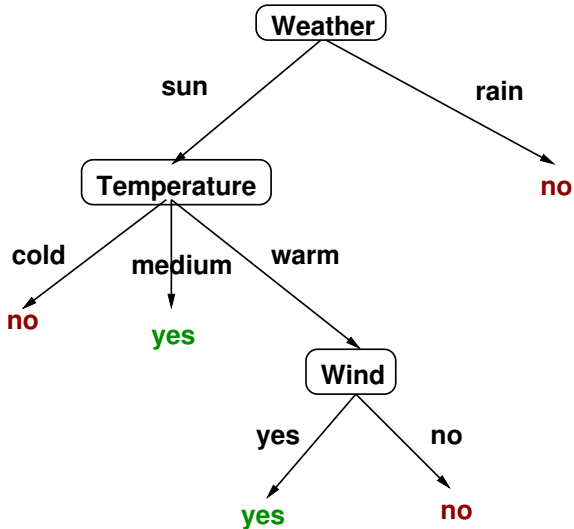
Example: Alternative Decision Trees (2)



Example: Alternative Decision Trees (2)



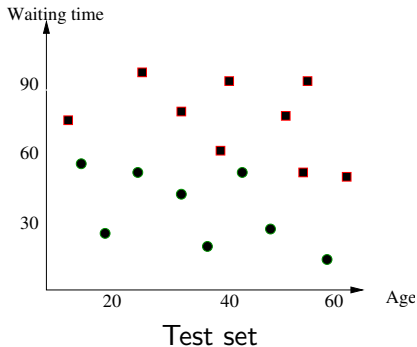
Example: Alternative Decision Trees (2)



Inductive bias: the simplest tree that fits the data (Occam's razor)

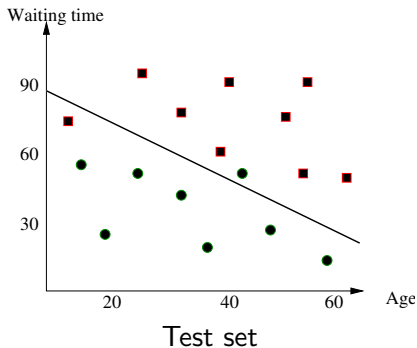
Applying the Classifier (1)

Assumption: the **decision boundary** learnt from the training set is also a good decision boundary for any test set **because the test set is drawn from the same distribution as the training set** (i.e., the two sets are not fundamentally different)



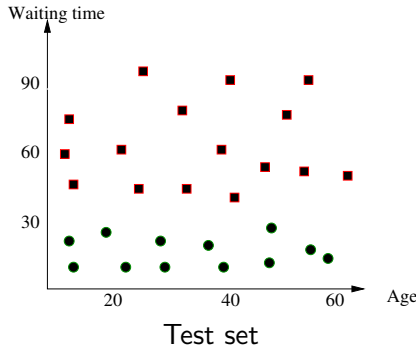
Applying the Classifier (1)

Assumption: the **decision boundary** learnt from the training set is also a good decision boundary for any test set **because the test set is drawn from the same distribution as the training set** (i.e., the two sets are not fundamentally different)



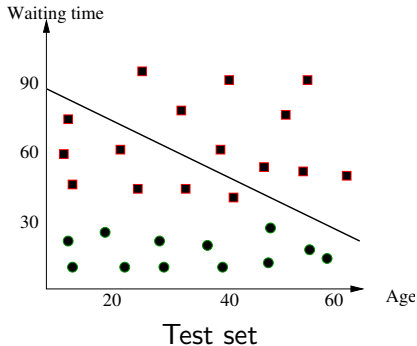
Applying the Classifier (2)

What happens if the test data are differently distributed?



Applying the Classifier (2)

What happens if the test data are differently distributed?



Co-Training

- two classifiers, representing **different views** of the data
- train both on a small labelled seed set
- apply both to large unlabelled set
- classifier A trains classifier B and vice versa

Example: Named Entity Recognition

Peter saw a documentary on **United Cakes Ltd.**

It was said they were planning to buy shares in **Henderson.**

Co-Training

- two classifiers, representing **different views** of the data
- train both on a small labelled seed set
- apply both to large unlabelled set
- classifier A trains classifier B and vice versa

Example: Named Entity Recognition

Peter saw a documentary on **United Cakes Ltd.**

It was said they were planning to buy shares in **Henderson.**

Co-Training

- two classifiers, representing **different views** of the data
- train both on a small labelled seed set
- apply both to large unlabelled set
- classifier A trains classifier B and vice versa

Example: Named Entity Recognition

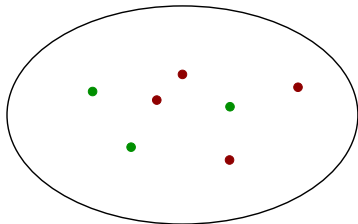
Peter saw a documentary on **United Cakes Ltd.**

It was said they were planning to **buy shares in Henderson.**

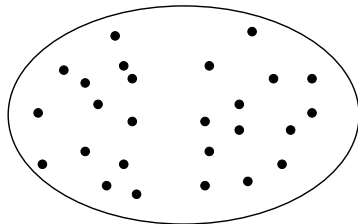
Co-Training: Example

Classifier A

Classifier B

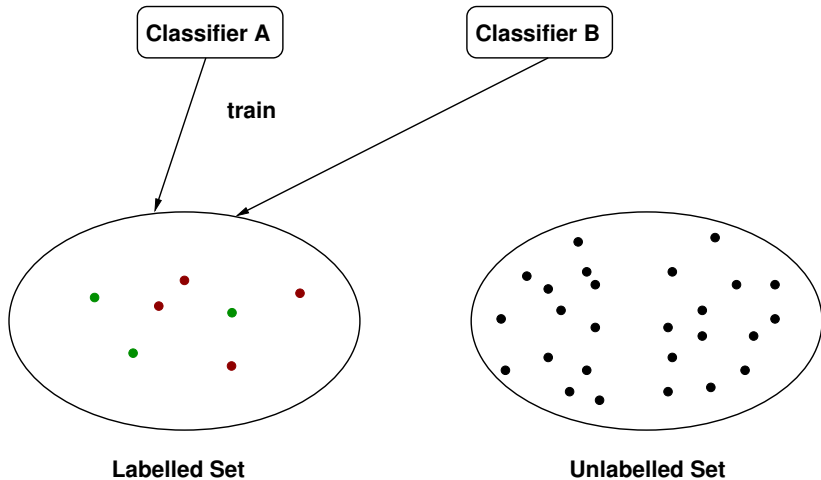


Labeled Set

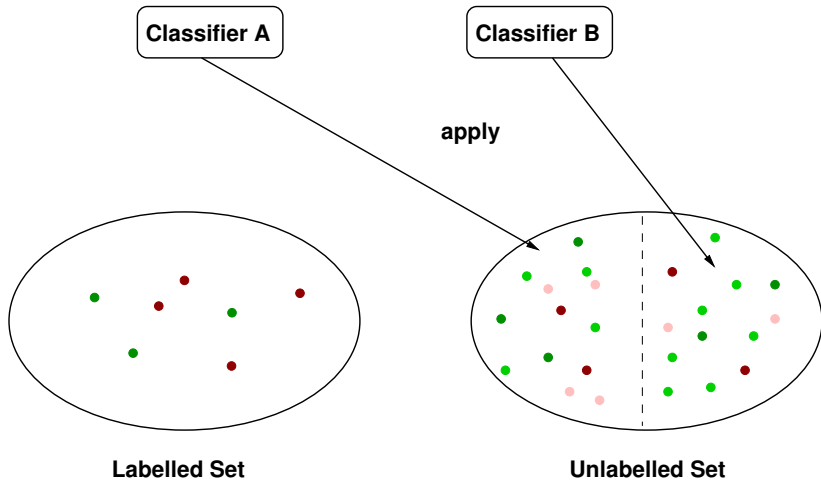


Unlabelled Set

Co-Training: Example



Co-Training: Example

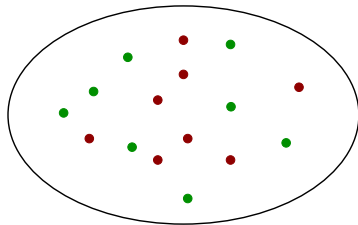


Co-Training: Example

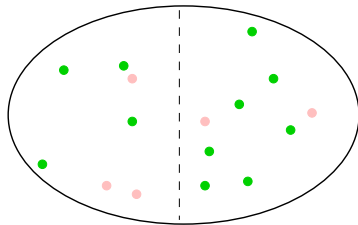
Classifier A

Classifier B

select new training examples

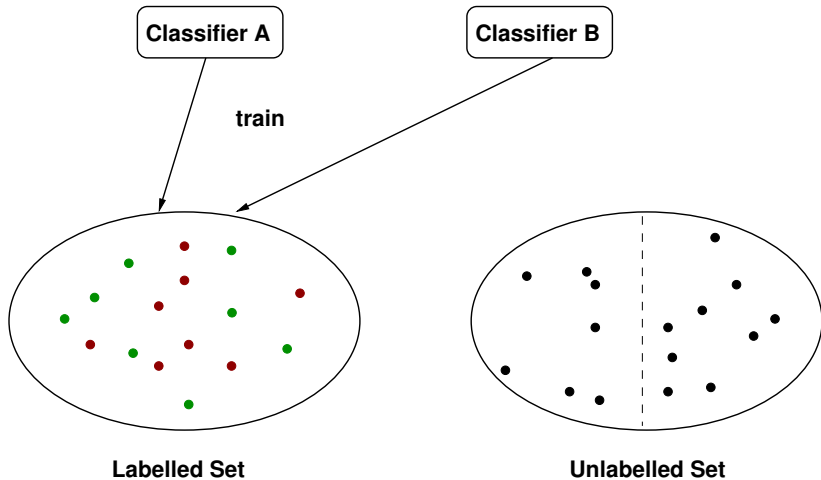


Labeled Set



Unlabelled Set

Co-Training: Example



Self-Training

- similar to co-training but only one classifier
- not clear whether it works