# CS-C3240 – Machine Learning D

**Round 3: From features to classification**

## Stephan Sigg

Department of Communications and Networking
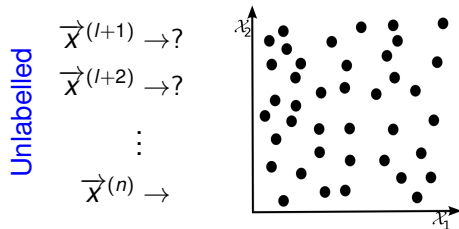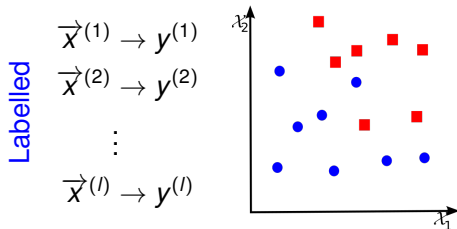Aalto University, School of Electrical Engineering
stephan.sigg@aalto.fi

Version 2.3, January 23, 2022
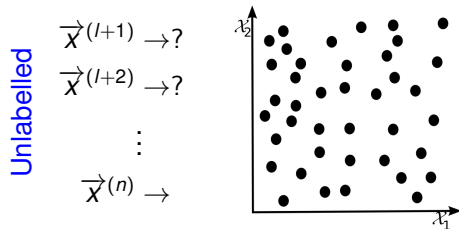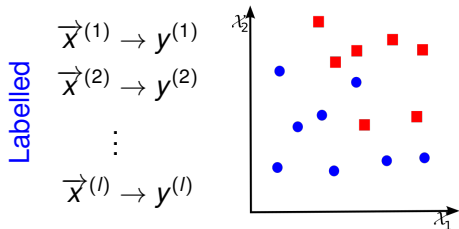
# Overview

Shortage of labelled data

# Shortage of labelled data

Given a set $\mathcal{Z}$ of data points $\overrightarrow{z}_1, \ldots, \overrightarrow{z}_n$, features $\mathcal{X}_1, \ldots, \mathcal{X}_m$ and labels $y_1, \ldots, y_o$, assume partially labelled sets $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ and $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{l+1, \ldots, n\}$ as well as $\overrightarrow{x}^{(k)} = x_1^{(k)}, \ldots, x_m^{(k)}, k \in \{1, \ldots, n\}$

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
January 23, 2022
3 / 8

# Shortage of labelled data

Given a set $\mathcal{Z}$ of data points $\overrightarrow{z}_1, \ldots, \overrightarrow{z}_n$, features $\mathcal{X}_1, \ldots, \mathcal{X}_m$ and labels $y_1, \ldots, y_o$, assume partially labelled sets $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ and $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{l+1, \ldots, n\}$ as well as $\overrightarrow{x}^{(k)} = x_1^{(k)}, \ldots, x_m^{(k)}, k \in \{1, \ldots, n\}$



## Problem:

- Unlabelled training data is often easy to obtain
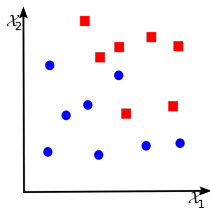- **However:** labelling the data requires significant manual work

**Aalto University**
School of Electrical
Engineering

**mbient**
**Intelligence**

**Stephan Sigg**
January 23, 2022
3 / 8

# Shortage of labelled data

**Automated labeling through semi-supervised learning**

Increase amount of labelled data via semi-supervised learning

1. Start with labelled data $\qquad\qquad\qquad\qquad \vec{z}_i = \langle \vec{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 23, 2022
4 / 8

# Shortage of labelled data

**Automated labeling through semi-supervised learning**

Increase amount of labelled data via semi-supervised learning

1. Start with labelled data $\qquad \overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$

2. Train the classifier on the labelled data $\quad \hat{h}\left(\overrightarrow{\hat{w}}, \cdot\right) = \min_{i=1,..l} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$

**Aalto University**
School of Electrical
Engineering

**mbient**
**Intelligence**

**Stephan Sigg**
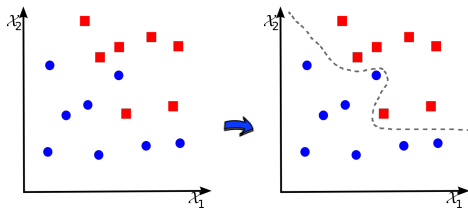January 23, 2022
4 / 8

# Shortage of labelled data

**Automated labeling through semi-supervised learning**

Increase amount of labelled data via semi-supervised learning

1. Start with labelled data $\quad\quad \overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$

2. Train the classifier on the labelled data $\quad \hat{h}\left(\overrightarrow{\hat{w}}, \cdot\right) = \min_{i=1,..l} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 23, 2022
4 / 8

# Shortage of labelled data

**Automated labeling through semi-supervised learning**

Increase amount of labelled data via semi-supervised learning

1. Start with labelled data $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \dots, l\}$

2. Train the classifier on the labelled data $\hat{h}\left(\overrightarrow{w}, \cdot\right) = \min_{i=1,\dots l} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$

3. Use $\hat{h}\left(\overrightarrow{w}, \cdot\right)$ to learn labels for $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{1, \dots, l\}$ $\hat{y}^{(j)} = \hat{h}\left(\overrightarrow{w}, \overrightarrow{x}^{(j)}\right)$

**Aalto University**
School of Electrical
Engineering

**Imbient**
**Intelligence**

**Stephan Sigg**
January 23, 2022
4 / 8

# Shortage of labelled data

**Automated labeling through semi-supervised learning**

Increase amount of labelled data via semi-supervised learning
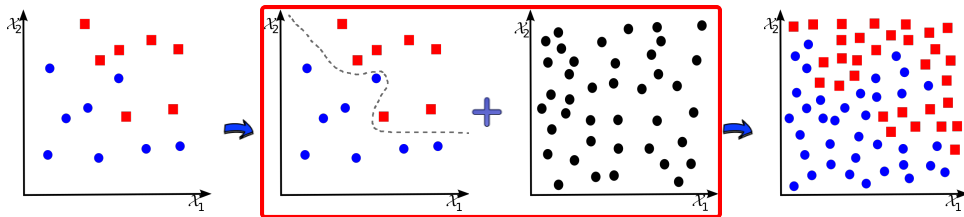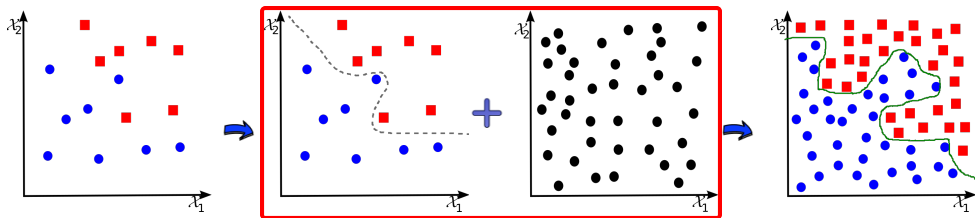
1. Start with labelled data $\quad \overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$

2. Train the classifier on the labelled data $\quad \hat{h}\left(\overrightarrow{w}, \cdot\right) = \min_{i=1,..l} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$

3. Use $\hat{h}\left(\overrightarrow{w}, \cdot\right)$ to learn labels for $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{1, \ldots, l\}$ $\quad \hat{y}^{(j)} = \hat{h}\left(\overrightarrow{w}, \overrightarrow{x}^{(j)}\right)$

4. Train new classifier $\hat{h}'$ on $\langle \overrightarrow{x}^{(1)}, y^{(1)} \rangle, \ldots, \langle \overrightarrow{x}^{(l+1)}, \hat{y}^{(l+1)} \rangle$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

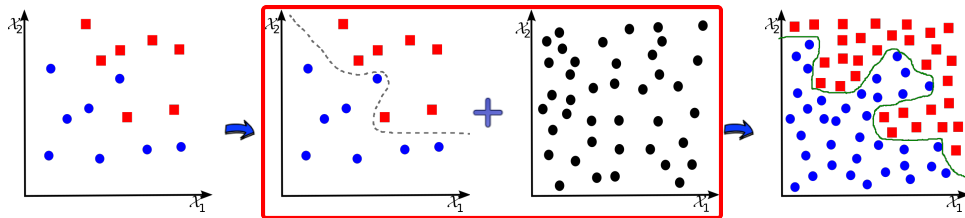**Stephan Sigg**
January 23, 2022
4 / 8

# Shortage of labelled data

**Automated labeling through semi-supervised learning**

Remark:

- No guaranteed success → Empirical validation required
- Introducing weights to samples can reduce dependency on learned labels

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
**January 23, 2022**
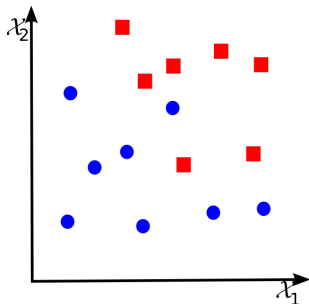5 / 8

# Shortage of labelled data – Automatic labelling

Provided independent feature sub-sets (perspectives) $\{\mathcal{X}\}_s$ with $\bigcup_s \{\mathcal{X}\}_s = \mathcal{X}_1, \ldots, \mathcal{X}_m$ and $\bigcap_s \{\mathcal{X}\}_s = \emptyset$, multiple classification models $h_s(\overrightarrow{w}_s, \overrightarrow{x})$ are trained to these sub-sets using the labelled data $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ to iteratively label unlabelled data $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{l+1, \ldots, n\}$

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
January 23, 2022
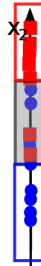6 / 8

# Shortage of labelled data – Automatic labelling

Provided independent feature sub-sets (perspectives) $\{\mathcal{X}\}_s$ with $\bigcup_s\{\mathcal{X}\}_s = \mathcal{X}_1, \ldots, \mathcal{X}_m$ and $\bigcap_s\{\mathcal{X}\}_s = \emptyset$, multiple classification models $h_s(\overrightarrow{w}_s, \overrightarrow{x})$ are trained to these sub-sets using the labelled data $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ to iteratively label unlabelled data $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{l+1, \ldots, n\}$

**1** Train classifier $h_s$ for each $\{\mathcal{X}\}_s$

$$\hat{h}_s(\overrightarrow{w}_s, \cdot) = \min_{i \in \{\mathcal{X}_s\}} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$$

**Aalto University**
School of Electrical
Engineering

Ambient Intelligence

**Stephan Sigg**
January 23, 2022
6 / 8

# Shortage of labelled data – Automatic labelling

Provided independent feature sub-sets (perspectives) $\{\mathcal{X}\}_s$ with $\bigcup_s \{\mathcal{X}\}_s = \mathcal{X}_1, \ldots, \mathcal{X}_m$ and $\bigcap_s \{\mathcal{X}\}_s = \emptyset$, multiple classification models $h_s(\overrightarrow{w}_s, \overrightarrow{x})$ are trained to these sub-sets using the labelled data $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ to iteratively label unlabelled data $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{l+1, \ldots, n\}$

1. Train classifier $h_s$ for each $\{\mathcal{X}\}_s$  $\qquad \hat{h}_s(\overrightarrow{w_s}, \cdot) = \min_{i \in \{\mathcal{X}_s\}} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$

2. Apply $\hat{h}_s(\overrightarrow{w_s}, \cdot)$ to $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$  $\qquad \hat{y}^{(j)} = \hat{h}_s(\overrightarrow{w_s}, \overrightarrow{x}^{(j)})$

**Aalto University**
School of Electrical
Engineering

**Ambient**
Intelligence

**Stephan Sigg**
January 23, 2022
6 / 8

# Shortage of labelled data – Automatic labelling

Provided independent feature sub-sets (perspectives) $\{\mathcal{X}\}_s$ with $\bigcup_s \{\mathcal{X}\}_s = \mathcal{X}_1, \ldots, \mathcal{X}_m$ and $\bigcap_s \{\mathcal{X}\}_s = \emptyset$, multiple classification models $h_s(\overrightarrow{w}_s, \overrightarrow{x})$ are trained to these sub-sets using the labelled data $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ to iteratively label unlabelled data $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{l+1, \ldots, n\}$

1. Train classifier $h_s$ for each $\{\mathcal{X}\}_s$ $\qquad \hat{h}_s(\overrightarrow{\hat{w}}_s, \cdot) = \min_{i \in \{\mathcal{X}_s\}} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$

2. Apply $\hat{h}_s(\overrightarrow{\hat{w}}_s, \cdot)$ to $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ $\qquad \hat{y}^{(j)} = \hat{h}_s(\overrightarrow{\hat{w}}_s, \overrightarrow{x}^{(j)})$

3. Add $\langle \overrightarrow{x}, \hat{y}^{(j)} \rangle$ with highest confidence to $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$

**Aalto University**
School of Electrical
Engineering

**Ambient**
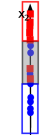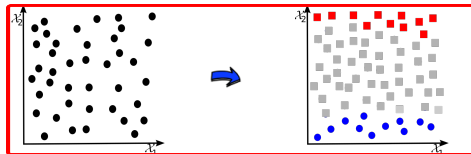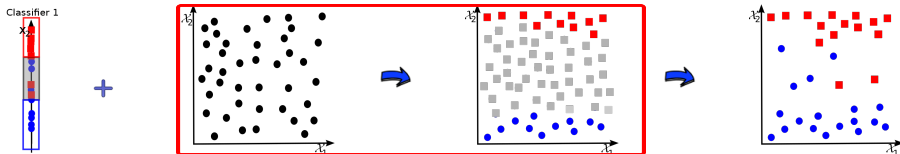**Intelligence**

**Stephan Sigg**
January 23, 2022
6 / 8

# Shortage of labelled data – Automatic labelling

Co-training

Provided independent feature sub-sets (perspectives) $\{\mathcal{X}\}_s$ with $\bigcup_s\{\mathcal{X}\}_s = \mathcal{X}_1, \ldots, \mathcal{X}_m$ and $\bigcap_s\{\mathcal{X}\}_s = \emptyset$, multiple classification models $h_s(\overrightarrow{w}_s, \overrightarrow{x})$ are trained to these sub-sets using the labelled data $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ to iteratively label unlabelled data $\overrightarrow{z}_j = \langle \overrightarrow{x}^{(j)} \rangle, j \in \{l+1, \ldots, n\}$

1. Train classifier $h_s$ for each $\{\mathcal{X}\}_s$ $\qquad \hat{h}_s(\overrightarrow{\hat{w}_s}, \cdot) = \min_{i \in \{\mathcal{X}_s\}} \mathcal{L}\left(h(\overrightarrow{w}, \overrightarrow{x}^{(i)}), y^{(i)}\right)$

2. Apply $\hat{h}_s(\overrightarrow{\hat{w}_s}, \cdot)$ to $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$ $\qquad \hat{y}^{(j)} = \hat{h}_s(\overrightarrow{\hat{w}_s}, \overrightarrow{x}^{(j)})$

3. Add $\langle \overrightarrow{x}, \hat{y}^{(j)} \rangle$ with highest confidence to $\overrightarrow{z}_i = \langle \overrightarrow{x}^{(i)}, y^{(i)} \rangle, i \in \{1, \ldots, l\}$

4. Iterate over over all classifiers $h_s$ until convergence reached

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 23, 2022
6 / 8

# Questions?

Stephan Sigg
`stephan.sigg@aalto.fi`

Si Zuo
`si.zuo@aalto.fi`

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 23, 2022
7 / 8

# Literature

- C.M. Bishop: Pattern recognition and machine learning, Springer, 2007.
- R.O. Duda, P.E. Hart, D.G. Stork: Pattern Classification, Wiley, 2001.

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 23, 2022
8 / 8