# CS-C3240 – Machine Learning D

**Anomaly detection, Recommender Systems, Online learning**

Stephan Sigg

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

Version 1.0, January 19, 2022

# Learning goals

- Collaborative Filtering
- Recommender Systems
- Batch gradient descent
- Online learning

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
2 / 33

# Outline

Recommender systems

Stochastic Classification

Online learning

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
3 / 33

Elektronik

USB Datenkabel für ...
EUR 15,00
Warum empfehlen?

LG E960 E 960 Nexus ...
EUR 71,49
Warum empfehlen?

GiYa Technology ...
EUR 5,99
Warum empfehlen?

Original Google Nexus ...
EUR 88,88
Warum empfehlen?

Werkzeug Reparatur ...
EUR 15,99
Warum empfehlen?

Schraubenzieher ...
EUR 7,15
Warum empfehlen?

› Alle Empfehlungen in Elektronik anzeigen

Computer & Zubehör

Logitech C270 USB HD ...
EUR 22,37
Warum empfehlen?

EDIMAX EW-7811UN ...
EUR 7,99
Warum empfehlen?

Präzisier Werkzeugset ...
EUR 9,29
Warum empfehlen?

Drahtloses Qi ...
EUR 14,99
Warum empfehlen?

Aelor® [i-Pack] ...
EUR 9,99
Warum empfehlen?

Salcar® TV Stick ...
EUR 13,18
Warum empfehlen?

› Alle Empfehlungen in Computer & Zubehör anzeigen

Musik

Neuerscheinung
Modern Blues
– The Waterboys
EUR 14,99
Warum empfehlen?

Neuerscheinung
Feelings aus der Asche
– Olli Schulz
EUR 14,99
Warum empfehlen?

Bring mich nach Hause
– Wir Sind Helden
Noch 7 Stück auf Lager.
Warum empfehlen?

Ein Leichtes Schwert
– Judith Holofernes
EUR 12,39
Noch 11 Stück auf Lager.
Warum empfehlen?

Tausend Wirre Worte ...
EUR 6,99
Noch 18 Stück auf Lager.
Warum empfehlen?

La Dura (Re-Release)
– Sportfreunde Stiller
EUR 12,12
Noch 1 Stück auf Lager.
Warum empfehlen?

› Alle Empfehlungen in Musik anzeigen

Bücher

Citrus 01
Sakurasta
EUR 6,95
Warum empfehlen?

Bedeutung und ...
James M. Heisig
EUR 23,36
Warum empfehlen?

Wir sind Helden: ...
EUR 9,96
Warum empfehlen?

Die Kana lernen und ...
James W. Heisig
EUR 14,90
Warum empfehlen?

Vokabelkarten ...
Hefei Huang
EUR 14,90
Warum empfehlen?

Kanji und Kana: Die ...
Wolfgang Hadamitzky
EUR 29,80
Warum empfehlen?

Aalto University
School of Electrical
Engineering

Ambient Intelligence

Stephan Sigg
January 19, 2022
4 / 33

# Recommender systems

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
5 / 33

# Recommender systems

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
5 / 33

# Recommender systems

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
5 / 33

# Recommender systems

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
5 / 33

# Recommender systems

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
5 / 33

# Recommender systems

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
January 19, 2022
5 / 33

# Recommender systems

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
5 / 33

# Recommender systems

## Task of Recommender systems

Given these ratings for a number of products, predict likely user-ratings for products that have not yet been rated

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
6 / 33

|  | 🤵 | 🤠 | 😌 | 🌞 | tools | accessories |
|---|---|---|---|---|---|---|
| 📱 | ⭐⭐⭐⭐☆ | ⭐⭐⭐⭐⭐ | ❓ | ☆☆☆☆☆ | 0.9 | 0.0 |
| 🛠 | ⭐⭐⭐⭐☆ | ⭐⭐⭐⭐⭐ | ❓ | ⭐☆☆☆☆ | 1.0 | 0.3 |
| 🔋 | ⭐⭐⭐⭐☆ | ❓ | ⭐⭐☆☆☆ | ❓ | 0.4 | 0.8 |
| 📷 | ❓ | ☆☆☆☆☆ | ⭐⭐⭐⭐⭐ | ⭐⭐⭐⭐⭐ | 0.2 | 0.9 |
| 📲 | ⭐☆☆☆☆ | ⭐☆☆☆☆ | ⭐⭐⭐⭐⭐ | ❓ | 0.0 | 1.0 |

**Aalto University**
School of Electrical
Engineering

**A**mbient
**I**ntelligence

**Stephan Sigg**
January 19, 2022
7 / 33

| tools | | accessories |
|---|---|---|
| 0.9 | $x^{(1)}$ | 0.0 |
| 1.0 | $x^{(2)}$ | 0.3 |
| 0.4 | $x^{(3)}$ | 0.8 |
| 0.2 | $x^{(4)}$ | 0.9 |
| 0.0 | $x^{(5)}$ | 1.0 |

Represent items in terms of weighted feature vectors

$$x^{(4)} = \begin{bmatrix} 0.2 \\ 0.9 \end{bmatrix} \quad w^{(1)} = \begin{bmatrix} \mathbf{5} \\ \mathbf{1} \end{bmatrix}$$

$$h(x) \Rightarrow \left( w^{(1)} \right)^T x^{(4)} = 0.2 \cdot 5 + 0.9 \cdot 1 = 1,9$$

**Aalto University**
School of Electrical
Engineering

**A**mbient
**I**ntelligence

**Stephan Sigg**
January 19, 2022
7 / 33

| tools | | accessories |
|---|---|---|
| 0.9 | $x^{(1)}$ | 0.0 |
| 1.0 | $x^{(2)}$ | 0.3 |
| 0.4 | $x^{(3)}$ | 0.8 |
| 0.2 | $x^{(4)}$ | 0.9 |
| 0.0 | $x^{(5)}$ | 1.0 |

Learn weights from provided ratings for single user $j$ (Logistic regression):

$$\min_{w^{(j)}} \frac{1}{2} \sum_{i:y^{(i,j)} \neq ?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^{m} \left( w_k^{(j)} \right)^2$$

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
8 / 33

| | tools | | accessories |
|---|---|---|---|
| | 0.9 | $x^{(1)}$ | 0.0 |
| | 1.0 | $x^{(2)}$ | 0.3 |
| | 0.4 | $x^{(3)}$ | 0.8 |
| | 0.2 | $x^{(4)}$ | 0.9 |
| | 0.0 | $x^{(5)}$ | 1.0 |

Learn weights from provided ratings for all users $1, \ldots, N$:

$$\min_{w^{(1)}, \ldots, w^{(N)}} \frac{1}{2} \sum_{j=1}^{N} \sum_{i: y^{(i,j)} \neq ?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{N} \sum_{k=1}^{m} \left( w_k^{(j)} \right)^2$$

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
9 / 33

# Recommender systems

$$\min_{w^{(1)},\ldots,w^{(N)}} \frac{1}{2} \sum_{j=1}^{N} \sum_{i:y^{(i,j)}\neq?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{N} \sum_{k=1}^{m} \left( w_k^{(j)} \right)^2$$

**Aalto University**
School of Electrical
Engineering

**Imbient Intelligence**

**Stephan Sigg**
January 19, 2022
10 / 33

# Recommender systems

**Optimisation algorithm**

$$\min_{w^{(1)},\ldots,w^{(N)}} \frac{1}{2} \sum_{j=1}^{N} \sum_{i:y^{(i,j)}\neq ?} \left( \left( w^{(j)} \right)^{T} x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{N} \sum_{k=1}^{m} \left( w_k^{(j)} \right)^2$$

Gradient descent update:

$$w_k^{(j)} = w_k^{(j)} - \delta \left( \sum_{i=y^{(i,j)}\neq ?} \left( \left( w^{(j)} \right)^{T} x^{(i)} - y^{(i,j)} \right) x_k^{(i)} + \lambda w_k^{(j)} \right)$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
11 / 33

# Recommender systems

**Optimisation algorithm**

$$\min_{w^{(1)},\ldots,w^{(N)}} \frac{1}{2} \sum_{j=1}^{N} \sum_{i:\, y^{(i,j)} \neq ?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{N} \sum_{k=1}^{m} \left( w_k^{(j)} \right)^2$$

Gradient descent update:

$$w_k^{(j)} = w_k^{(j)} - \delta \left( \underbrace{\sum_{i=y^{(i,j)} \neq ?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right) x_k^{(i)} + \lambda w_k^{(j)}}_{\text{partial derivative}} \right)$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
12 / 33

# Recommender systems

## Collaborative filtering



| | tools | | accessories |
|---|---|---|---|
| | 0.9 | $x^{(1)}$ | 0.0 |
| | 1.0 | $x^{(2)}$ | 0.3 |
| | 0.4 | $x^{(3)}$ | 0.8 |
| | 0.2 | $x^{(4)}$ | 0.9 |
| | 0.0 | $x^{(5)}$ | 1.0 |

We are able to calculate the weights given the feature vectors

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
13 / 33

# Recommender systems

**Collaborative filtering**



We are able to calculate the weights given the feature vectors
→ But how do we obtain these feature vectors?

Aalto University
School of Electrical
Engineering

Stephan Sigg
January 19, 2022
13 / 33

# Recommender systems

**Collaborative filtering**

Users might tell their preferences

e.g. more interested in tools or accessories



$$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix} \quad w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$$

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
January 19, 2022
14 / 33

# Recommender systems

**Collaborative filtering**

Users might tell their preferences

e.g. more interested in tools or accessories



$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$  $w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$  $w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$  $w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$

$$\left(w^{(1)}\right)^T x^{(1)} \approx 4; \quad \left(w^{(2)}\right)^T x^{(1)} \approx 5;$$

$$\left(w^{(3)}\right)^T x^{(1)} \approx ?; \quad \left(w^{(4)}\right)^T x^{(1)} \approx 0$$

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
January 19, 2022
15 / 33

# Recommender systems

**Collaborative filtering**

Users might tell their preferences

e.g. more interested in tools or accessories



$$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix} \quad w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$$

$$\left(w^{(1)}\right)^T x^{(1)} \approx 4; \quad \left(w^{(2)}\right)^T x^{(1)} \approx 5;$$

$$\left(w^{(3)}\right)^T x^{(1)} \approx ?; \quad \left(w^{(4)}\right)^T x^{(1)} \approx 0$$

From these weights we can estimate the feature values

**Aalto University**
School of Electrical
Engineering

**mbient
ntelligence**

**Stephan Sigg**
January 19, 2022
15 / 33

# Recommender systems

## Optimisation algorithm

Given the weights/preferences $w^{(1)}, \ldots, w^{(N)}$, we are able to infer a feature $x^{(i)}$

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j:y^{(i,j)} \neq ?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^{m} \left( x_k^{(i)} \right)^2$$

**Aalto University**
**School of Electrical**
**Engineering**

**Stephan Sigg**
**January 19, 2022**
**16 / 33**

# Recommender systems

Optimisation algorithm

Given the weights/preferences $w^{(1)}, \ldots, w^{(N)}$, we are able to infer a feature $x^{(i)}$

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j:y^{(i,j)} \neq ?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^{m} \left( x_k^{(j)} \right)^2$$

Given the weights/preferences $w^{(1)}, \ldots, w^{(N)}$, we are able to infer $x^{(1)}, \ldots, x^{(n)}$

$$\min_{x^{(1)}, \ldots, x^{(n)}} \frac{1}{2} \sum_{i=1}^{n} \sum_{j:y^{(i,j)} \neq ?} \left( \left( w^{(j)} \right)^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{i=1}^{n} \sum_{k=1}^{m} \left( x_k^{(i)} \right)^2$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
17 / 33

# Recommender systems

## Naive (iterative) Collaborative filtering algorithm

Given $x^{(1)}, \ldots, x^{(n)}$, we are able to estimate $w^{(1)}, \ldots, w^{(N)}$

Given $w^{(1)}, \ldots, w^{(N)}$, we are able to estimate $x^{(1)}, \ldots, x^{(n)}$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
18 / 33

# Recommender systems

## Naive (iterative) Collaborative filtering algorithm

Given $x^{(1)}, \ldots, x^{(n)}$, we are able to estimate $w^{(1)}, \ldots, w^{(N)}$

Given $w^{(1)}, \ldots, w^{(N)}$, we are able to estimate $x^{(1)}, \ldots, x^{(n)}$

## Collaborative filtering (naive)

Init: Randomly initialise the $w^{(i)}$

Repeat:
- Estimate the $x^{(i)}$ from the $w^{(i)}$
- Estimate the $w^{(i)}$ from the $x^{(i)}$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
18 / 33

# Recommender systems

### Naive (iterative) Collaborative filtering algorithm

Given $x^{(1)}, \ldots, x^{(n)}$, we are able to estimate $w^{(1)}, \ldots, w^{(N)}$

Given $w^{(1)}, \ldots, w^{(N)}$, we are able to estimate $x^{(1)}, \ldots, x^{(n)}$

## Collaborative filtering (naive)

Init: Randomly initialise the $w^{(i)}$

Repeat:
- Estimate the $x^{(i)}$ from the $w^{(i)}$
- Estimate the $w^{(i)}$ from the $x^{(i)}$

$\rightarrow$ CF iteratively improves the estimates for $x^{(i)}$ and $w^{(i)}$

$\rightarrow$ Algorithm <u>collaborates</u> with users: by providing some information about their preferences, it computes and improves the features

**Aalto University**
School of Electrical
Engineering

Ambient
Intelligence

**Stephan Sigg**
January 19, 2022
18 / 33

# Recommender systems

## Collaborative filtering

$$\min_{w^{(1)},\dots,w^{(N)}} \quad \frac{1}{2}\sum_{j=1}^{N}\sum_{i:y^{(i,j)}\neq ?}\left(\left(w^{(j)}\right)^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2}\sum_{j=1}^{N}\sum_{k=1}^{m}\left(w_k^{(j)}\right)^2$$

$$\min_{x^{(1)},\dots,x^{(n)}} \quad \frac{1}{2}\sum_{i=1}^{n}\sum_{j:y^{(i,j)}\neq ?}\left(\left(w^{(j)}\right)^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2}\sum_{i=1}^{n}\sum_{k=1}^{m}\left(x_k^{(i)}\right)^2$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
19 / 33

# Recommender systems

**Collaborative filtering**

$$\min_{w^{(1)},\ldots,w^{(N)}} \quad \frac{1}{2}\sum_{j=1}^{N}\sum_{i:y^{(i,j)}\neq?}\left(\left(w^{(j)}\right)^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2}\sum_{j=1}^{N}\sum_{k=1}^{m}\left(w_k^{(j)}\right)^2$$

$$\min_{x^{(1)},\ldots,x^{(n)}} \quad \frac{1}{2}\sum_{i=1}^{n}\sum_{j:y^{(i,j)}\neq?}\left(\left(w^{(j)}\right)^T x^{(i)} - y^{(i,j)}\right)^2 + \frac{\lambda}{2}\sum_{i=1}^{n}\sum_{k=1}^{m}\left(x_k^{(i)}\right)^2$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
20 / 33

# Recommender systems

**Collaborative filtering**

$$\min_{w^{(1)},\ldots,w^{(N)}} \quad \frac{1}{2}\sum_{j=1}^{N}\sum_{i:y^{(i,j)}\neq?}\left(\left(w^{(j)}\right)^{T}x^{(i)}-y^{(i,j)}\right)^{2}+\frac{\lambda}{2}\sum_{j=1}^{N}\sum_{k=1}^{m}\left(w_{k}^{(j)}\right)^{2}$$

$$\min_{x^{(1)},\ldots,x^{(n)}} \quad \frac{1}{2}\sum_{i=1}^{n}\sum_{j:y^{(i,j)}\neq?}\left(\left(w^{(j)}\right)^{T}x^{(i)}-y^{(i,j)}\right)^{2}+\frac{\lambda}{2}\sum_{i=1}^{n}\sum_{k=1}^{m}\left(x_{k}^{(i)}\right)^{2}$$

Minimize $w^{(i)}$ and $x^{(i)}$ simultaneously:

$$\min_{x^{(1)},\ldots,x^{(n)},w^{(1)},\ldots,w^{(N)}} \quad \frac{1}{2}\sum_{i,j:y^{(i,j)}\neq?}\left(\left(w^{(j)}\right)^{T}x^{(i)}-y^{(i,j)}\right)^{2}$$

$$+\frac{\lambda}{2}\sum_{i=1}^{n}\sum_{k=1}^{m}\left(x_{k}^{(i)}\right)^{2}+\frac{\lambda}{2}\sum_{j=1}^{N}\sum_{k=1}^{m}\left(w_{k}^{(j)}\right)^{2}$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
21 / 33

# Recommender systems

Init: Randomly initialise the $w^{(j)}$ and $x^{(i)}$

Optimisation: Simultaneously minimise the above function for $w^{(j)}$ and $x^{(i)}$

Gradient descent:

$$x_k^{(i)} = x_k^{(i)} - \delta \left( \sum_{j=y^{(i,j)} \neq ?} \left( \left(w^{(j)}\right)^T x^{(i)} - y^{(i,j)} \right) w_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$w_k^{(j)} = w_k^{(j)} - \delta \left( \sum_{i=y^{(i,j)} \neq ?} \left( \left(w^{(j)}\right)^T x^{(i)} - y^{(i,j)} \right) x_k^{(i)} + \lambda w_k^{(j)} \right)$$

Prediction: For a user $i$ with parameters $w^{(j)}$ and an item with learned featues $x$, estimate a rating of $\left(w^{(j)}\right)^T x$

**Aalto University**
School of Electrical
Engineering

Ambient
Intelligence

**Stephan Sigg**
January 19, 2022
22 / 33

# Outline

Recommender systems

Stochastic Classification

Online learning

**Aalto University**
School of Electrical
Engineering

**m**bient
**I**ntelligence

**Stephan Sigg**
January 19, 2022
23 / 33

# Stochastic Classification

Model training becomes slow with increasing data size

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
January 19, 2022
24 / 33

# Stochastic Classification

Model training becomes slow with increasing data size

$\rightarrow$ Because of repeated looping over the complete data set until convergence

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
24 / 33

# Stochastic Classification

Model training becomes slow with increasing data size

$\rightarrow$ Because of repeated looping over the complete data set until convergence

## Solution

- Randomly iterate the update only over a subset of items instead of repeatedly considering the whole data set.

**Aalto University**
School of Electrical
Engineering

**Ambient
Intelligence**

**Stephan Sigg**
January 19, 2022
24 / 33

# Stochastic Classification

**Example: Gradient descent**



$$\text{minimize } L[(\mathcal{X}, \mathcal{Y}), h(\cdot)] \quad = \quad \frac{1}{2n} \sum_{i=1}^{n} \left( w^T x^{(i)} - y^{(i)} \right)^2$$

$$\text{Repeat } \forall j : w_j \quad = \quad w_j - \delta \cdot \frac{\partial L[(\mathcal{X}, \mathcal{Y}), h(\cdot)]}{\partial w_j}$$

$$\rightarrow \forall j : w_j \quad = \quad w_j - \delta \cdot \frac{1}{n} \sum_{i=1}^{n} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

**Aalto University**
School of Electrical
Engineering

**mbient**
**Intelligence**

**Stephan Sigg**
January 19, 2022
25 / 33

# Stochastic Classification

**Example: Gradient descent**



$\rightarrow$ For single gradient descent-step, algorithms loops over all samples!

$$\text{minimize } L[(\mathcal{X}, \mathcal{Y}), h(\cdot)] = \frac{1}{2n} \sum_{i=1}^{n} \left( w^T x^{(i)} - y^{(i)} \right)^2$$

$$\text{Repeat } \forall j : w_j = w_j - \delta \cdot \frac{\partial L[(\mathcal{X}, \mathcal{Y}), h(\cdot)]}{\partial w_j}$$

$$\rightarrow \forall j : w_j = w_j - \delta \cdot \frac{1}{n} \sum_{i=1}^{n} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

**Aalto University**
School of Electrical
Engineering

**Stephan Sigg**
January 19, 2022
25 / 33

# Stochastic Classification

**Example: Gradient descent**

$\rightarrow$ For a single gradient descent-step, the algorithms loops over all samples!

$$
\begin{aligned}
\text{minimize } L[(\mathcal{X}, \mathcal{Y}), h(\cdot)] &= \frac{1}{2n} \sum_{i=1}^{n} \left( w^T x^{(i)} - y^{(i)} \right)^2 \\
\text{Repeat } \forall j : w_j &= w_j - \delta \cdot \frac{\partial L[(\mathcal{X}, \mathcal{Y}), h(\cdot)]}{\partial w_j} \\
\rightarrow \forall j : w_j &= w_j - \delta \cdot \frac{1}{n} \sum_{i=1}^{n} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}
\end{aligned}
$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
26 / 33

# Stochastic Classification

**Example: Gradient descent**

$\rightarrow$ For a single gradient descent-step, the algorithms loops over all samples!

To speed up the algorithm, compute gradient descent updates from individual training samples (randomly ordered)
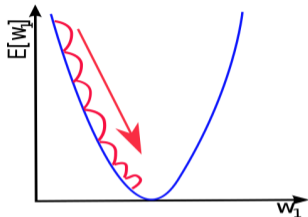
$$\text{minimize } L[(\mathcal{X}, \mathcal{Y}), h(\cdot)] \;=\; \frac{1}{2n} \sum_{i=1}^{n} \left( w^T x^{(i)} - y^{(i)} \right)^2$$

$$\text{Repeat } \forall j : w_j \;=\; w_j - \delta \cdot \frac{\partial L[(\mathcal{X}, \mathcal{Y}), h(\cdot)]}{\partial w_j}$$

$$\rightarrow \forall j : w_j \;=\; w_j - \delta \cdot \frac{1}{n} \sum_{i=1}^{n} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

**Aalto University**
School of Electrical
Engineering

**mbient**
**ntelligence**

**Stephan Sigg**
January 19, 2022
26 / 33

# Stochastic Classification

**Example: Gradient descent**

<u>Standard:</u>

$$\text{minimize } L[(\mathcal{X}, \mathcal{Y}), h(\cdot)] = \frac{1}{2n} \sum_{i=1}^{n} \left( w^T x^{(i)} - y^{(i)} \right)^2$$

$$\text{Repeat } \forall j : w_j = w_j - \delta \cdot \frac{\partial}{\partial w_j} E[w_j]$$

$$\rightarrow \forall j : w_j = w_j - \delta \cdot \frac{1}{n} \sum_{i=1}^{n} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
27 / 33

# Stochastic Classification

**Example: Gradient descent**

<u>Standard:</u>

$$\text{minimize } E[W] \;=\; \frac{1}{2n} \sum_{i=1}^{n} \left( w^T x^{(i)} - y^{(i)} \right)^2$$

$$\text{Repeat } \forall j : w_j \;=\; w_j - \delta \cdot \frac{\partial}{\partial w_j} E[w_j]$$

$$\rightarrow \forall j : w_j \;=\; w_j - \delta \cdot \frac{1}{n} \sum_{i=1}^{n} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

---

<u>Stochastic:</u>

Repeat over all training examples $i$ (random order):

$$\Rightarrow \forall j : \quad w_j = w_j - \delta \cdot \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
28 / 33

# Stochastic Classification

The stochastic implementation will generally move the parameters towards the global minimum (...but not always!)

**Aalto University**
School of Electrical
Engineering

**Ambient**
Intelligence

**Stephan Sigg**
January 19, 2022
29 / 33

# Stochastic Classification

The stochastic implementation will generally move the parameters towards the global minimum (...but not always!)

> Greatly speeds up the gradient descent steps as it does not loop over all samples in each single iteration

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
29 / 33

# Stochastic Classification

The stochastic implementation will generally move the parameters towards the global minimum (...but not always!)

Greatly speeds up the gradient descent steps as it does not loop over all samples in each single iteration

Tradeoff use $1 \leq k \leq n$ random examples for each gradient descent update $l = 1, 1 + k, 1 + 2k, \ldots$

$$\Rightarrow \forall j : w_j = w_j - \delta \cdot \frac{1}{l} \sum_{l=i}^{i+k} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

**Aalto University**
School of Electrical
Engineering

**I**mbient
**I**ntelligence

**Stephan Sigg**
January 19, 2022
29 / 33

# Stochastic Classification

The stochastic implementation will generally move the parameters towards the global minimum (...but not always!)

Greatly speeds up the gradient descent steps as it does not loop over all samples in each single iteration

Tradeoff use $1 \leq k \leq n$ random examples for each gradient descent update
$I = 1, 1 + k, 1 + 2k, \ldots$

$$\Rightarrow \forall j : w_j = w_j - \delta \cdot \frac{1}{l} \sum_{l=i}^{i+k} \left( w_j x_j^{(i)} - y^{(i)} \right) \cdot x_j^{(i)}$$

$k > 1$ might be faster than $k = 1$ for parallelized code

**Aalto University**
School of Electrical
Engineering

**mbient**
**ntelligence**

**Stephan Sigg**
January 19, 2022
29 / 33

# Outline

Recommender systems

Stochastic Classification

Online learning

Aalto University
School of Electrical
Engineering

Ambient
Intelligence

Stephan Sigg
January 19, 2022
30 / 33

# Online learning

## Online learning

Given a continuous stream of input data, update the parameters of your algorithm on-the-fly

Example: learn user behaviour from website users

**Aalto University**
School of Electrical
Engineering

**Ambient**
**Intelligence**

**Stephan Sigg**
January 19, 2022
31 / 33

# Online learning

## Online learning

Given a continuous stream of input data, update the parameters of your algorithm on-the-fly

Example: learn user behaviour from website users

Similar to stochastic classification: Update the parameters based on individual training examples

$$\Rightarrow \forall j: \quad w_j = w_j - \delta \cdot \left( h(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)}$$

**Aalto University**
School of Electrical
Engineering

**mbient**
**ntelligence**

**Stephan Sigg**
January 19, 2022
31 / 33

# Online learning

## Online learning

Given a continuous stream of input data, update the parameters of your algorithm on-the-fly

Example: learn user behaviour from website users

Similar to stochastic classification: Update the parameters based on individual training examples

$$\Rightarrow \forall j: \quad w_j = w_j - \delta \cdot \left( h(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)}$$

$\Rightarrow$ Able to adapt to changing user behaviour over time

**Aalto University**
School of Electrical
Engineering

**mbient**
**ntelligence**

**Stephan Sigg**
January 19, 2022
31 / 33

# Questions?

Stephan Sigg
`stephan.sigg@aalto.fi`

Si Zuo
`si.zuo@aalto.fi`

**Aalto University**
School of Electrical
Engineering

**Ambient Intelligence**

**Stephan Sigg**
January 19, 2022
32 / 33

# 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 19, 2022
33 / 33