



## Data Science: Applied Text Mining (S42)

26 – 29 July 2021

Course location:

Course Director: Dr. Ayoub Bagheri

E-mail: ms.summerschool@uu.nl

| Saturday 24 & Sunday 25 July 2021 |             |   |
|-----------------------------------|-------------|---|
| Time                              | Activity    | Description   |
| 12.00 – 18.00                     | Key pick up | <i>You will find the exact key pick up location in the pre-departure information, which becomes available after you have paid the course fee.</i> |

| Day    | Time          | Type               | Description   | Location |
|--------|---------------|--------------------|---|----------|
| Monday | 09:00 – 10:30 | Lecture            | Introduction to text mining<br>Preprocessing text<br>Feature extraction   |          |
|        | 10:45 – 11:45 | Computer Lab       | Python in Google Colaboratory: <ul style="list-style-type: none"> <li>- Cleaning text</li> <li>- NLTK, Gensim, spaCy</li> </ul>                         |          |
|        | 11:45 – 12:30 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab   |          |
|        | 14:00 – 15:15 | Lecture            | Text classification: <ul style="list-style-type: none"> <li>- Binary classification</li> <li>- Multi-class classification</li> </ul> Sentiment analysis |          |
|        | 15:30 – 16:30 | Computer Lab       | Document-term matrix<br>Sentiment classification<br>News classification<br>Classification evaluation  |          |
|        | 16:30 – 17:00 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab   |          |



| Day            | Time          | Type               | Description   | Location |
|----------------|---------------|--------------------|---|----------|
| <b>Tuesday</b> | 09:00 – 10:30 | Lecture            | Feature selection in text: <ul style="list-style-type: none"> <li>- Dimensionality reduction</li> <li>- Principal component analysis</li> </ul>                   |          |
|                | 10:45 – 11:45 | Computer Lab       | Pointwise mutual information<br>PCA, t-SNE visualisation  |          |
|                | 11:45 – 12:30 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab   |          |
|                | 14:00 – 15:30 | Lecture            | Text clustering: <ul style="list-style-type: none"> <li>- K-Means, MBK, DBScan</li> <li>- Non-negative matrix factorization</li> <li>- Topic modelling</li> </ul> |          |
|                | 15:45 – 16:30 | Computer Lab       | Text clustering<br>Building an LDA model<br>Visualise the topics-words distributions<br>Optimal number of clusters<br>Clustering evaluation                       |          |
|                | 16:30 – 17:00 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab   |          |

| Day              | Time          | Type               | Description  | Location |
|------------------|---------------|--------------------|--|----------|
| <b>Wednesday</b> | 09:00 – 10:30 | Lecture            | Word embedding: <ul style="list-style-type: none"> <li>- Distributional hypothesis</li> <li>- Singular value decomposition</li> <li>- CBOW vs Skip-gram</li> </ul> |          |
|                  | 10:45 – 11:45 | Computer Lab       | Embedding layer<br>Retrieve Pre-trained embeddings<br>Visualise the embeddings<br>FastText (using subword information)   |          |
|                  | 11:45 – 12:30 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab  |          |
|                  | 14:00 – 15:30 | Lecture            | Convolutional neural networks<br>Recurrent neural networks <ul style="list-style-type: none"> <li>- LSTM, GRU</li> <li>- BERT</li> </ul>                           |          |
|                  | 15:45 – 16:30 | Computer Lab       | Multi-class text classification  |          |
|                  | 16:30 – 17:00 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab  |          |



| Day      | Time          | Type               | Description   | Location |
|----------|---------------|--------------------|---|----------|
| Thursday | 09:00 – 10:30 | Lecture            | Text summarization<br>Advanced deep learning with text: <ul style="list-style-type: none"><li>- Attention models</li><li>- Text generation with deep learning</li></ul> |          |
|          | 10:45 – 11:45 | Computer Lab       | Text summarization<br>Simple attention model for text   |          |
|          | 11:45 – 12:30 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab   |          |
|          | 14:00 – 15:30 | Lecture            | Bias and fairness in text mining<br>Responsible text mining   |          |
|          | 15:45 – 16:30 | Computer Lab       | A complete pipeline<br>Advanced topics  |          |
|          | 16:30 – 17:00 | Plenary Discussion | Students and teachers discuss and present their solutions to the computer lab   |          |

For information about the Social Programme,  
please visit the [Utrecht Summer School website!](#)