

Werkzeuge für das wissenschaftliche Arbeiten

Python for Machine Learning and Data Science

Magnus Bender
bender@ifis.uni-luebeck.de
Wintersemester 2022/23

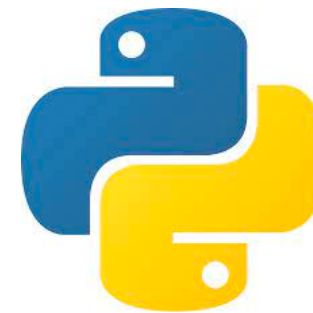
Inhaltsübersicht

1. Programmiersprache Python

a) *Einführung, Erste Schritte*

b) *Grundlagen*

c) *Fortgeschritten*



2. Auszeichnungssprachen

a) **LaTeX, Markdown**

L^AT_EX



3. Benutzeroberflächen und Entwicklungsumgebungen

a) Jupyter Notebooks lokal und in der Cloud (Google Colab)

4. Versionsverwaltung

a) Git, GitHub



5. Wissenschaftliches Rechnen

a) NumPy, SciPy



6. Datenverarbeitung und -visualisierung

a) Pandas, matplotlib, NLTK

7. Machine Learning (scikit-learn)

a) Grundlegende Ansätze (Datensätze, Auswertung)

b) Einfache Verfahren (Clustering, ...)



8. DeepLearning

a) TensorFlow, PyTorch, HuggingFace Transformers



Themen

I. Projektaufgabe 2

- Herangehensweise & Tipps

II. Auszeichnungssprachen

- Inhalt, Struktur & Form
- Semantische Auszeichnung
- Markdown
- LaTeX



Heute

Projektaufgabe 2

„Objektorientierung in Python“

- Klasse „Datensatz“ zur Verwaltung von Daten
 - Jedes Datum hat einen (eindeutigen) Schlüssel (Namen)
 - Oberklasse spezifiziert verschiedene Schnittstellen eines Datensatzes
 - Datensätze können erstellt, vereinigt, iteriert, ... werden
- Daten werden „extern“ in einen Datensatz hinzugefügt
- Datensätze werden „extern“ genutzt.

dataset_usage.py

```
from dataset import DataSetItem
from implementation import DataSet

data = DataSet([
    DataSetItem("Name 1", 11, "Inhalt 1")
])
data += DataSetItem("Name 2", 12, "Inhalt 2")

del data["Name 2"]

for item in data:
    print(item)
```

Hier im Sinne von „außerhalb des zu schreibenden Programms“

Herangehensweise & Tipps

1. Anforderungen an Klasse
2. Spezifikation jeder Schnittstelle exakt umsetzen
3. Schrittweise lösen
 - A. Schnittstellen einzeln nacheinander umsetzen
 - B. Schnittstellen testen

- Inhalte des ersten Teils (erste drei Vorlesungen) beachten



Warum diese Art von Aufgabe?

Eigene Daten, die von einer Bibliothek verarbeitet werden sollen und dafür passenden *ansprechbar* sein müssen.



II.

Auszeichnungssprachen

Aber zuerst:
Fragen zur Aufgabe 2?

Inhalt, St

Inhalt

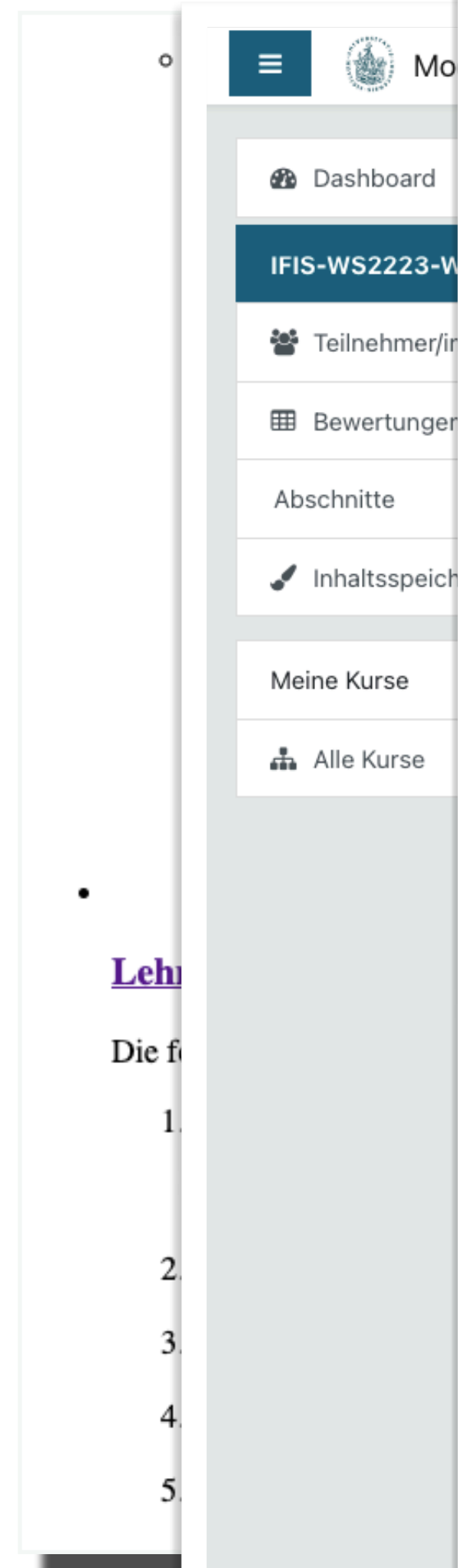
- Bedeutung eines Textes
- Die Wörter (und Sätze)

Struktur

- Aufbau des Dokuments
- Die Absätze, Kapitel und Überschriften

Form

- Aussehen des Dokuments
- Die Darstellung, wie Farben, Schriftart, Markierungen und Kästen



2. Subjective Content Descriptions

Algorithm 1 Training the SCD-word distribution matrix $\delta(\mathcal{D})$

```
1: function BUILDMATRIX( $\mathcal{D}$ ,  $g(\mathcal{D})$ )
2:   Input: Corpus  $\mathcal{D}$ , Set of SCDs  $g(\mathcal{D})$ 
3:   Output: SCD-word distribution matrix  $\delta(\mathcal{D})$ 
4:   Initialize an  $M \times L$  matrix  $\delta(\mathcal{D})$  with zeros
5:   for each  $d \in \mathcal{D}$  do
6:     for each  $(t, \rho) \in g(d)$  do
7:       for each  $w^d \in \text{win}_{d,\rho}$  do
8:          $\delta(\mathcal{D})[t][w^d] += I(w^d, \text{win}_{d,\rho})$ 
9:   return  $\delta(\mathcal{D})$ 
```

Kuhr et al. use a sliding window instead of our previously described sentence-based approach. The authors assume an SCD generates the words in a certain radius around the SCD's location, while we assume an SCD generates the words of the sentence at the SCD's location. The sentence-wise approach is required in this thesis due to the comparability to BERT working on whole sentences. Furthermore, a sliding window results in more computations and as we use larger corpora as Kuhr et al. sentence-wise iteration allows us to keep the computations sufficiently low.

After Algorithm 1 has finished, the SCD matrix needs to be normalized row-wise to meet the requirements of a probability distribution. However, we skip the normalization because multiple calculations on small decimal values on a computer reduce the accuracy. Later, we use the cosine similarity with the rows of the matrix and the cosine similarity does a normalization by definition. Thus, by skipping the normalization we save computational resources and get slightly more accurate results.

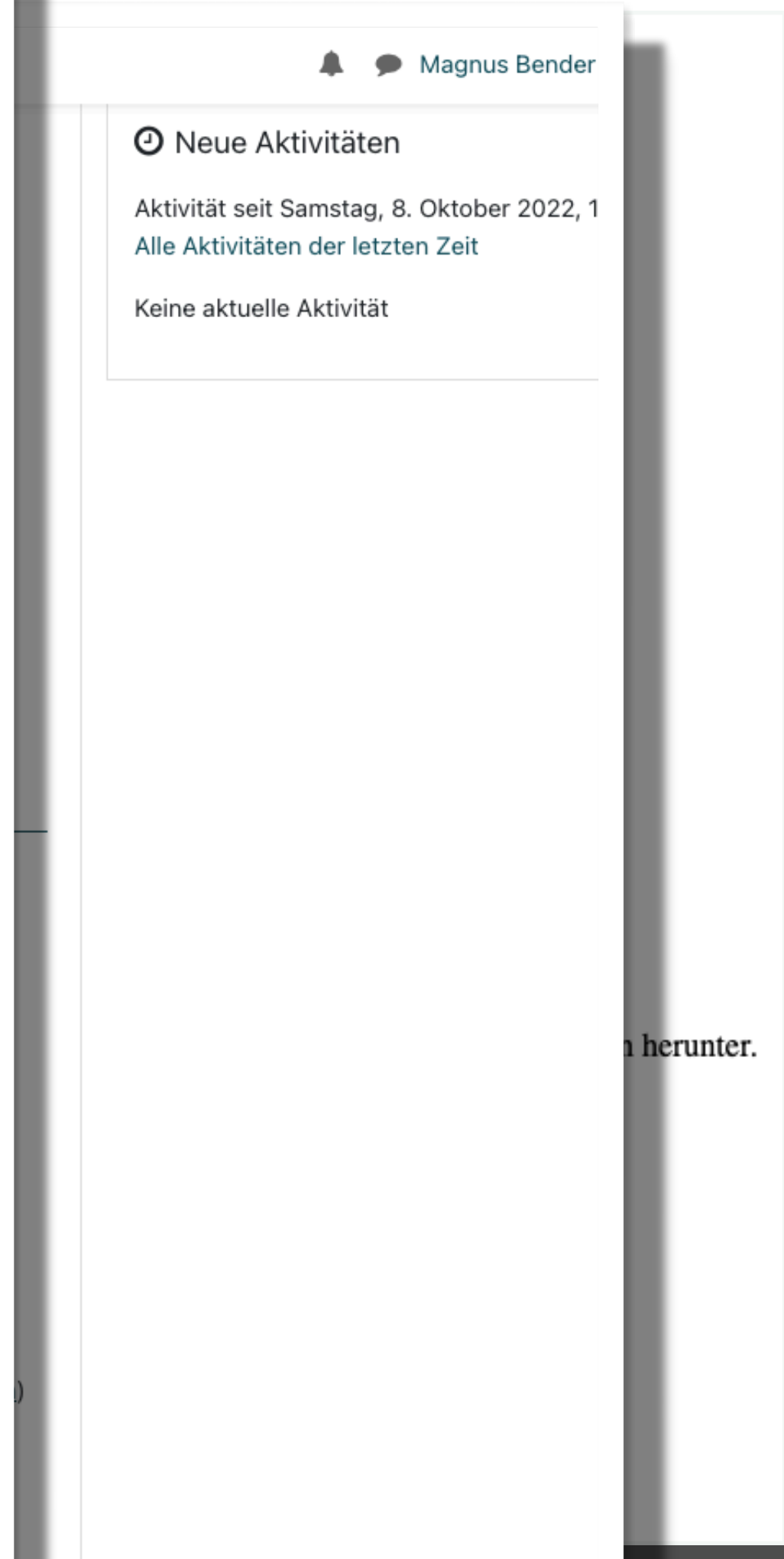
2.1. Most Probably Suited Subjective Content Descriptions

The previously described and trained SCD matrix can be used to estimate SCDs for a document without associated SCDs. First we formalize the MPS²CD problem and afterwards solve the problem by Algorithm 2 using the SCD matrix [KBBM19].

The MPS²CD problem asks for the M most probably suited SCDs t_1, \dots, t_M for a document d' given the SCD matrix $\delta(\mathcal{D})$:

$$g(d') = \arg \max_{t_1, \dots, t_M \in g(\mathcal{D})} P(t_1, \dots, t_M | d', \delta(\mathcal{D}))$$

The definition of the MPS²CD problem does not consider the sentence-wise iteration used while training the SCD matrix. We can reformulate the MPS²CD problem to



Semantische Auszeichnung

- Trennung zwischen Inhalt und Form
- Angabe des Inhalt mit Struktur
- Anschließend Formatierung der Struktur

Überschrift „Meine Übungsaufgabe“

Überschriften seien „**groß und fett**“

Struktur



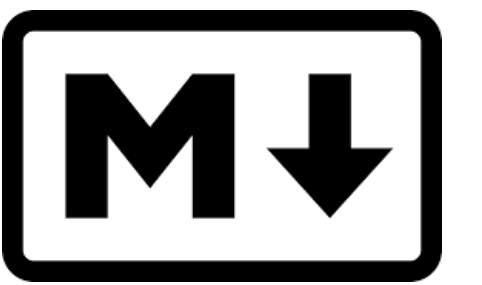
HTML, Markdown, LaTeX

Word, Libre Office, Pages

PDF, Vektorgrafiken

Pixelgrafik

Form



Markdown

Markdown

> From [Wikipedia](https://en.wikipedia/wiki/Markdown), the free encyclopedia

Article

Markdown is a lightweight markup language for creating formatted text using a plain text editor. John Gruber and Aaron Swartz created Markdown in 2004 as a markup language appealing to human readers in its source code form.

Paragraphs are separated by a blank line.

Two spaces at the end of a line produce a line break.

Text can be styled *italic*, **bold**, or `monospace`.

Markdown

From Wikipedia, the free encyclopedia

Article

Markdown is a lightweight markup language for creating formatted text using a plain-text editor. John Gruber and Aaron Swartz created Markdown in 2004 as a markup language that is appealing to human readers in its source code form.

Paragraphs are separated by a blank line.

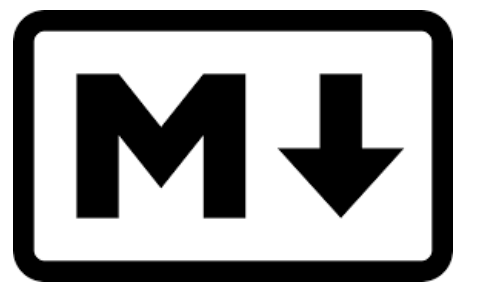
Two spaces at the end of a line

produce a line break. Text can be styled *italic*, **bold**, or `monospace`.

Inhalt und Struktur,
aber keine Form

Verschiedene
Formatierungen

Markdown




Es gibt auch Fußnoten^[1].

A	B	C
1	2	3

```
print("Hallo ;)")
```

- $\frac{1}{2}^2$
- $\frac{x_1^2 + 5x_2}{x_1}$

1. Tatsächlich ist das hier erweitertes Markdown 

Überschrift

Unterüberschrift

Quote

- **Fett**
- *kursiv*
- `Code`

1. [Links](#)

2.



Einmalig zu aktivieren: Einstellungen →
Texteditor wählen → „Einfacher Text“
Unter den Eingabefeldern ist dann eine
Auswahl „Markdown“ möglich.

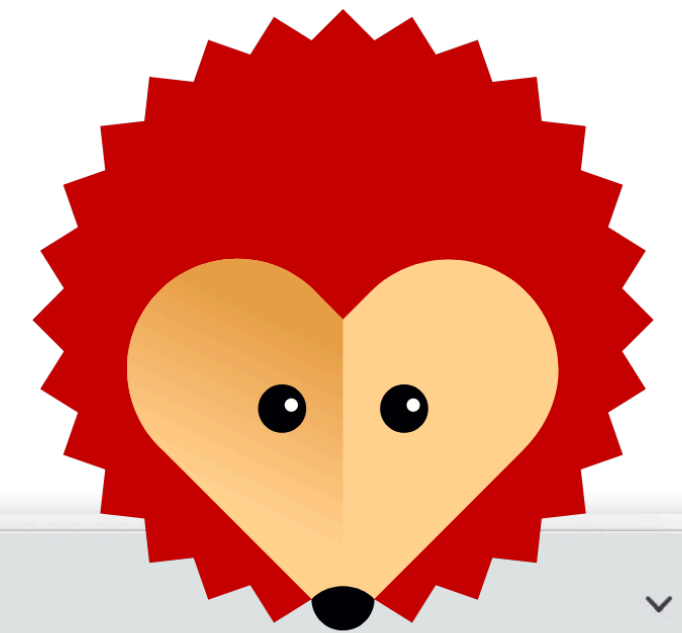
down



- Mittlerweile einer der Standards zum einfachen Formatieren von Texten
 - Kommunikation, z.B. *Moodle*, Discord, Slack
 - Code-Editoren, z.B. VS Code
 - Plattformen zur Versionsverwaltung, z.B. GitHub, GitLab
 - Plattformen zur Zusammenarbeit, z.B. HedgeDoc

„Oft lohnt es sich einfach mal
Markdown-Syntax zu testen.“

Beispiel: HedgeDoc



- Kollaborativ
- Markdown mit Erweiterungen
- Gleichungen (LaTeX)
- Diagramme
- Features & Demo

The screenshot displays the HedgeDoc web interface. The browser tab is titled "Markdown - HedgeDoc". The address bar shows "?both". The interface includes a top navigation bar with "HedgeDoc" logo, a toolbar with icons for eye, grid, edit, moon, and help, and buttons for "+ Neu", "Veröffentlichen", "Menü", and "1 ONLINE".

The editor area on the left shows the following Markdown code:

```
1 # Markdown
2
3 > From [Wikipedia].
   (https://en.wikipedia.org/wiki/Markdown),
   the free encyclopedia
4
5 ## Article
6
7 Markdown is a lightweight markup language
   for creating formatted text using a plain-
   text editor.
8 John Gruber and Aaron Swartz created
   Markdown in 2004 as a markup language that
   is appealing to human readers in its
   source code form.
9
10 Paragraphs are separated by a blank line.
11
12 Two spaces at the end of a line
13 produce a line break.
14 Text can be styled italic, bold, or
   `monospace`.
```

The rendered output on the right shows the following HTML:

VERÄNDERT VOR 2 MINUTEN FREELY

Markdown

From [Wikipedia](https://en.wikipedia.org/wiki/Markdown), the free encyclopedia

Article

Markdown is a lightweight markup language for creating formatted text using a plain-text editor.

John Gruber and Aaron Swartz created Markdown in 2004 as a markup language that is appealing to human readers in its source code form.

Paragraphs are separated by a blank line.

Two spaces at the end of a line produce a line break.

Text can be styled *italic*, **bold**, or `monospace`.

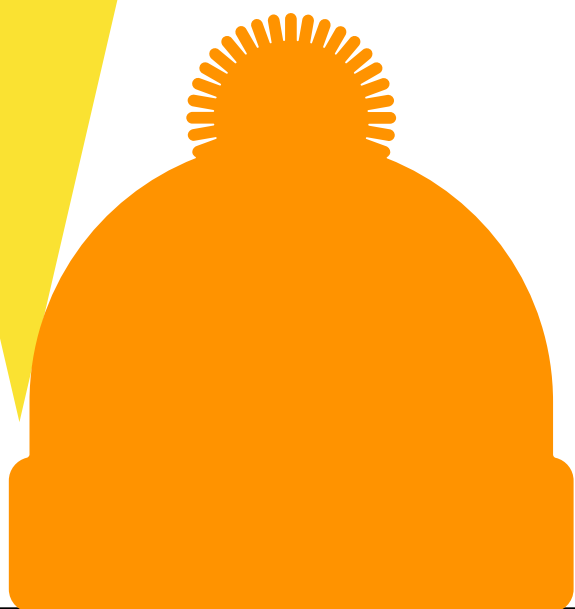
At the bottom, a status bar indicates "Line 1, Columns 1 - 16 Lines", "Spaces: 4", "SUBLIME", and "Length 495".

LaTeX

L^AT_EX

- Textsatzsystem
 - Geeignet für Übungszettel, Berichte, Zusammenfassungen, Abschlussarbeiten
 - Weniger für Notizen, Mitschriften (→ Markdown)
 - Insbesondere Unterstützung von Gleichungen, Literaturverzeichnisse, Inhaltsverzeichnisse

„LaTech nicht LaTeX!“



Wir sehen auch hier Trennung von Inhalt (im „document“) und Form in der Präambel.

Mein erstes Dokument

L^AT_EX

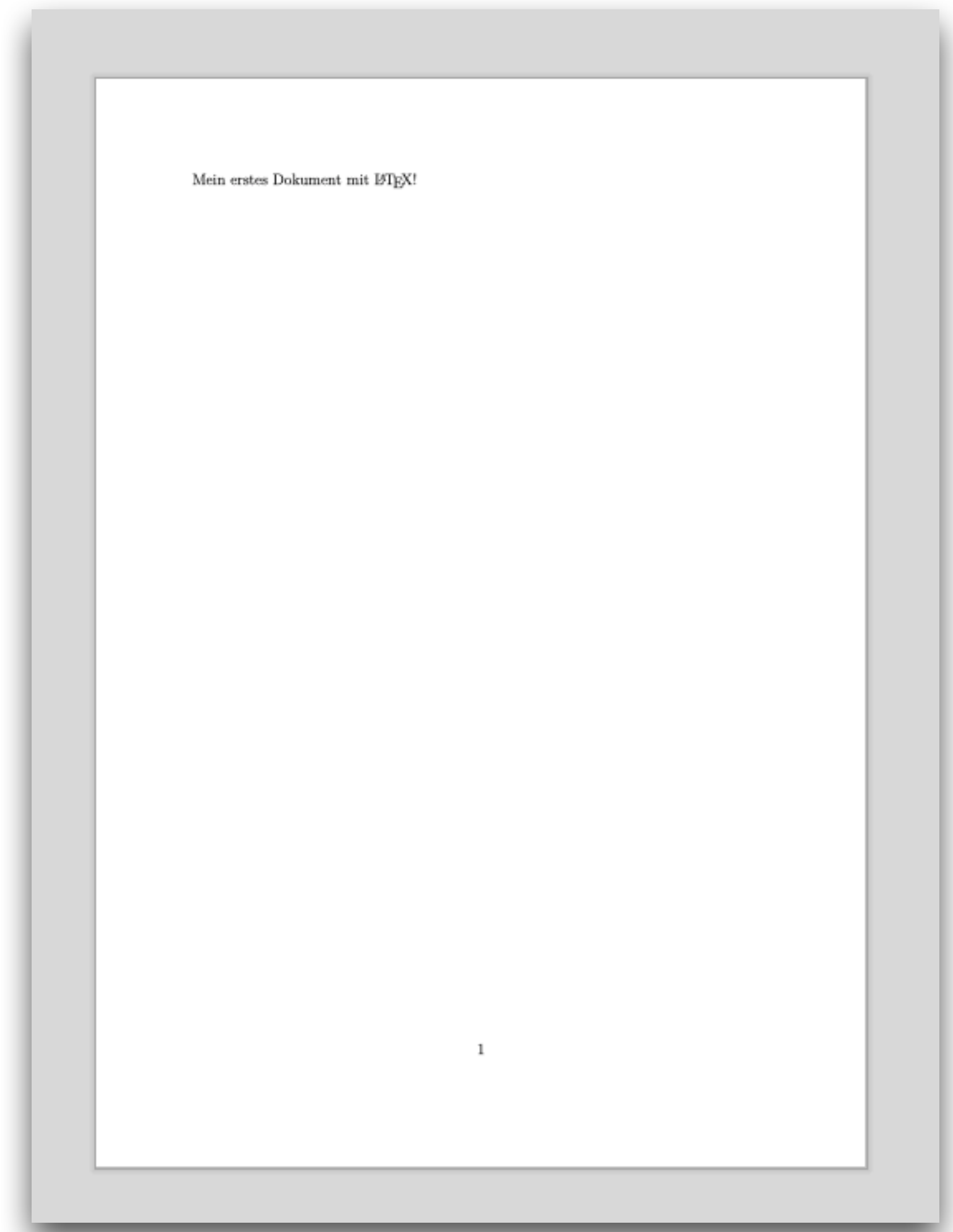
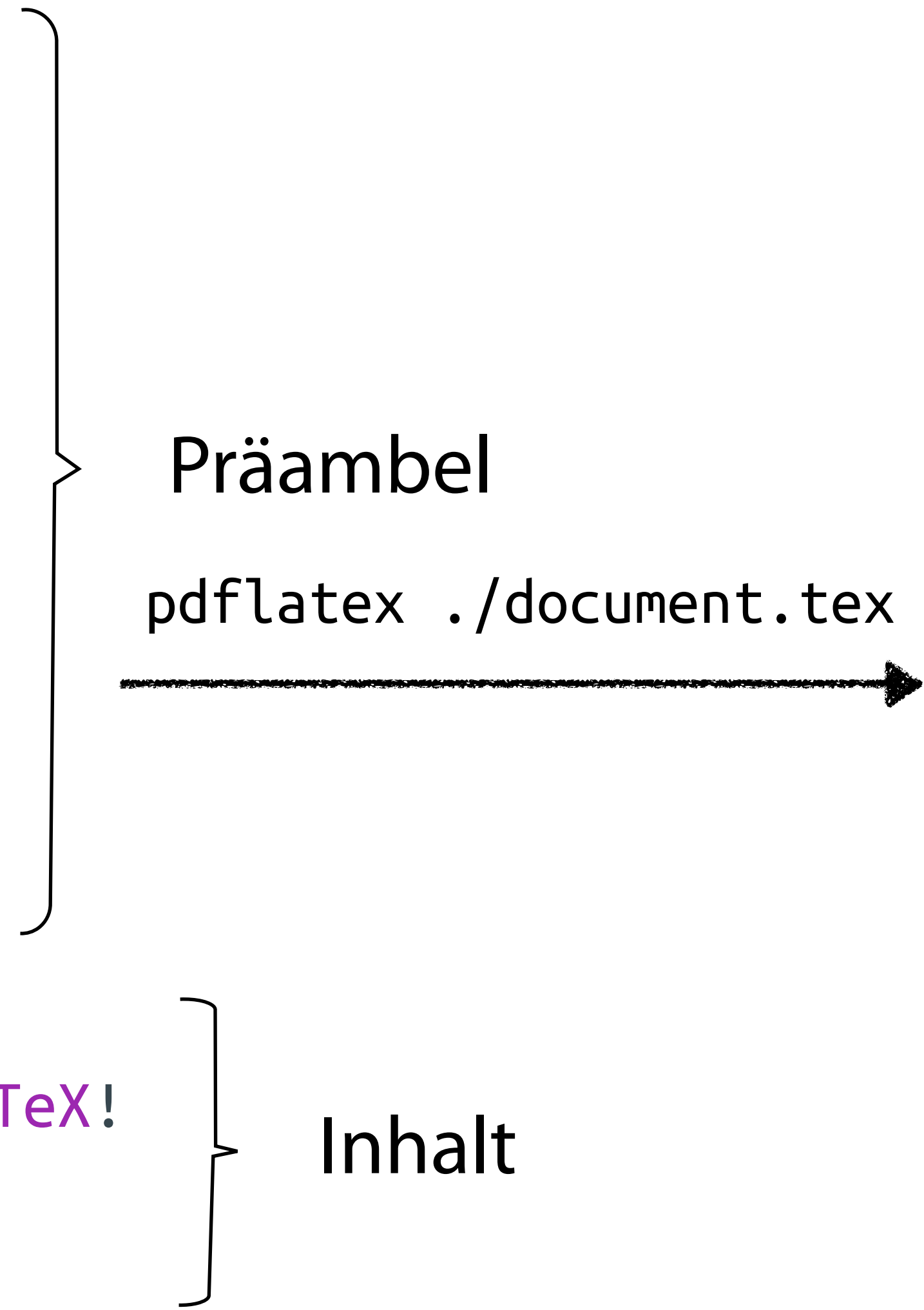
```
\documentclass[
  12pt, % Schriftgroesse
  a4paper, % Papier
  parskip=full % Absatzstil
]{scrartcl}

% Dateikodierung
\usepackage[utf8]{inputenc}
% Trennung "deutsch"
\usepackage[ngerman]{babel}
% Schriftart
\usepackage[T1]{fontenc}
\usepackage{lmodern}

\begin{document}

  Mein erstes Dokument mit \LaTeX!

\end{document}
```



3x ~~2x~~ pdflatex ./document.tex

L^AT_EX

```

\begin{document}
  \title{LaTeX Test}
  \author{Magnus Bender}
  \date{\today{} oder 15.10.2022}
  \maketitle

  \begin{center}
    Mein \textbf{erstes} \textsc{
  \end{center}

  \tableofcontents

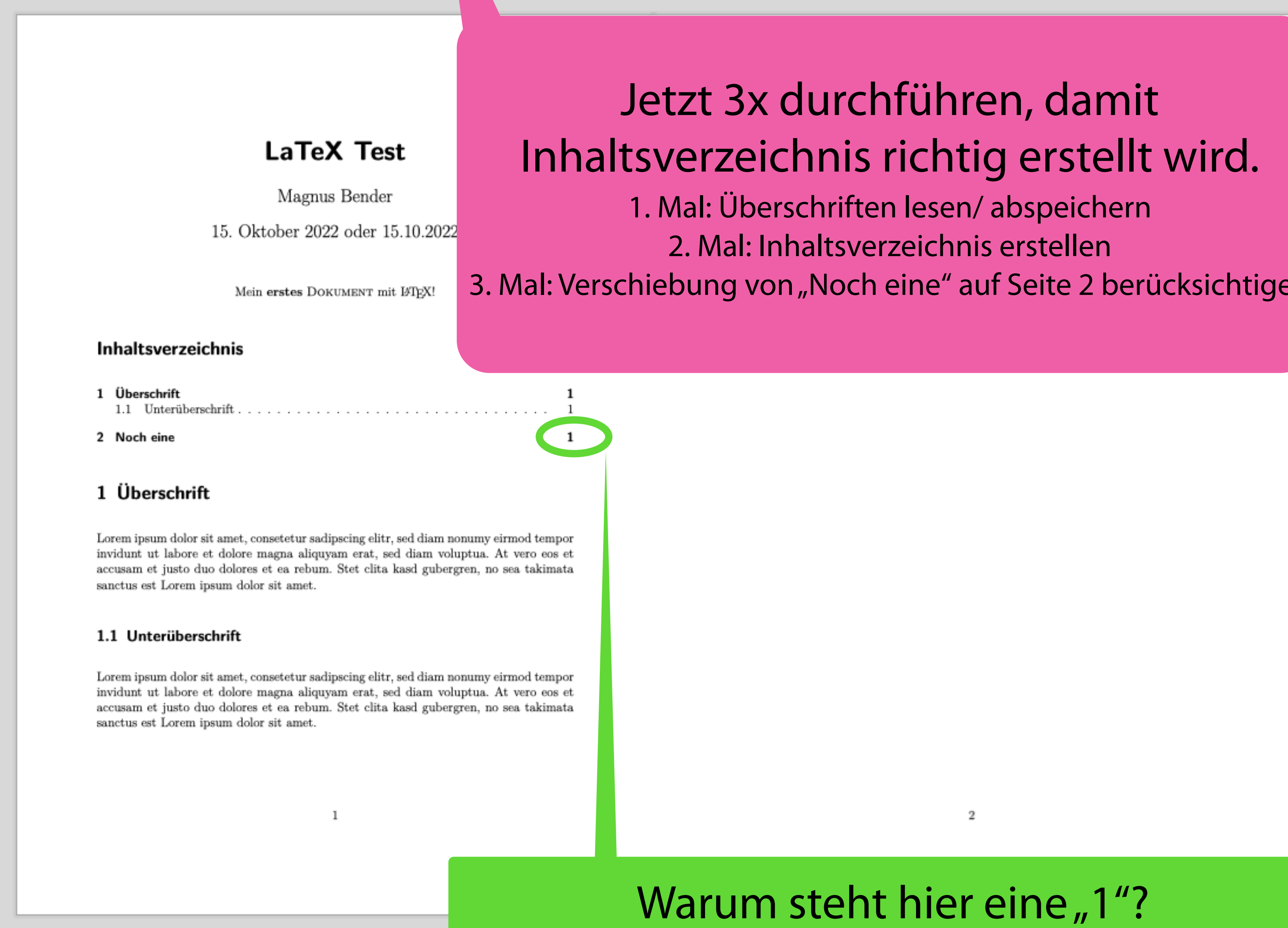
  \section{Überschrift}
  Lorem ipsum dolor sit amet,

  \subsection{Unterüberschrift}
  Lorem ipsum dolor sit amet
  At vero eos et accusam et

  \section{Noch eine}
  Lorem ipsum dolor sit amet,
  At vero eos et accusam et ..

  Lorem ipsum dolor sit amet,
  At vero eos et accusam et ..
\end{document}

```



Jetzt 3x durchführen, damit Inhaltsverzeichnis richtig erstellt wird.

1. Mal: Überschriften lesen/ abspeichern
2. Mal: Inhaltsverzeichnis erstellen
3. Mal: Verschiebung von „Noch eine“ auf Seite 2 berücksichtigen

Warum steht hier eine „1“?

2x pdflatex ./document.tex

Gleicher Inhalt, aber „Buch“, somit andere Form!

```
\documentclass[]{\book}
```

% ...

```
\begin{document}  
  \title{LaTeX Test}  
  \author{Magnus Bender}  
  \date{\today}  
  \maketitle
```

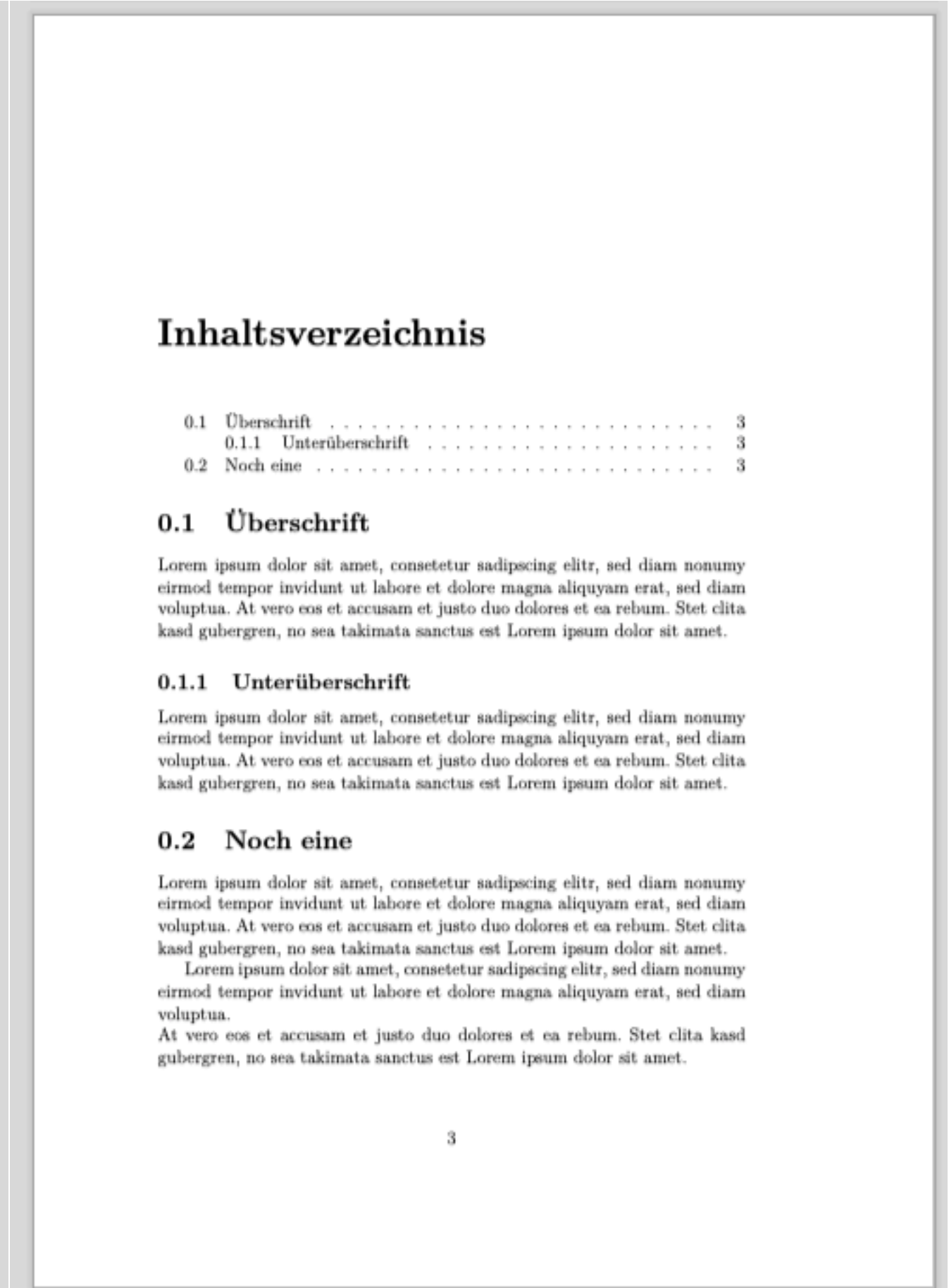
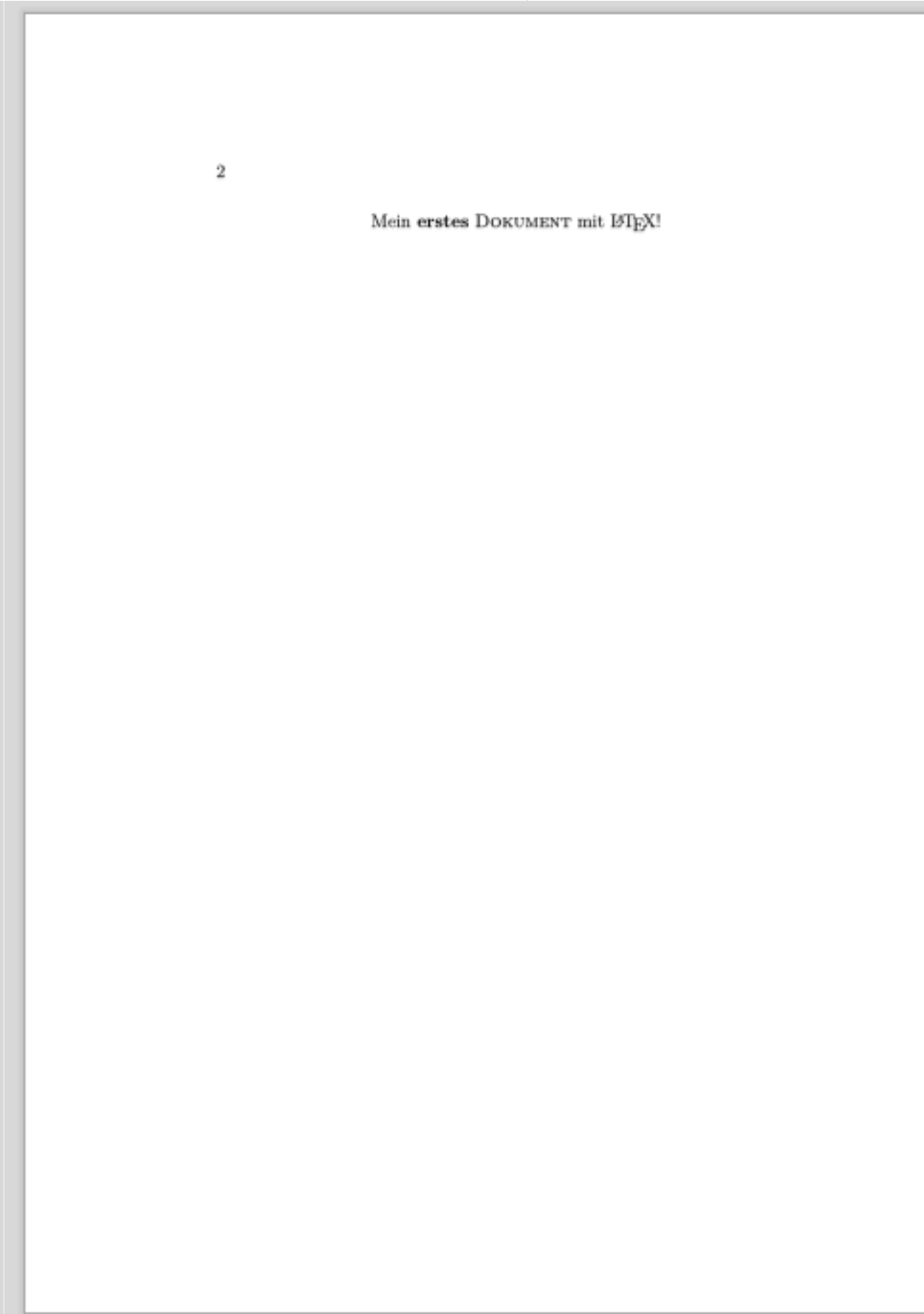
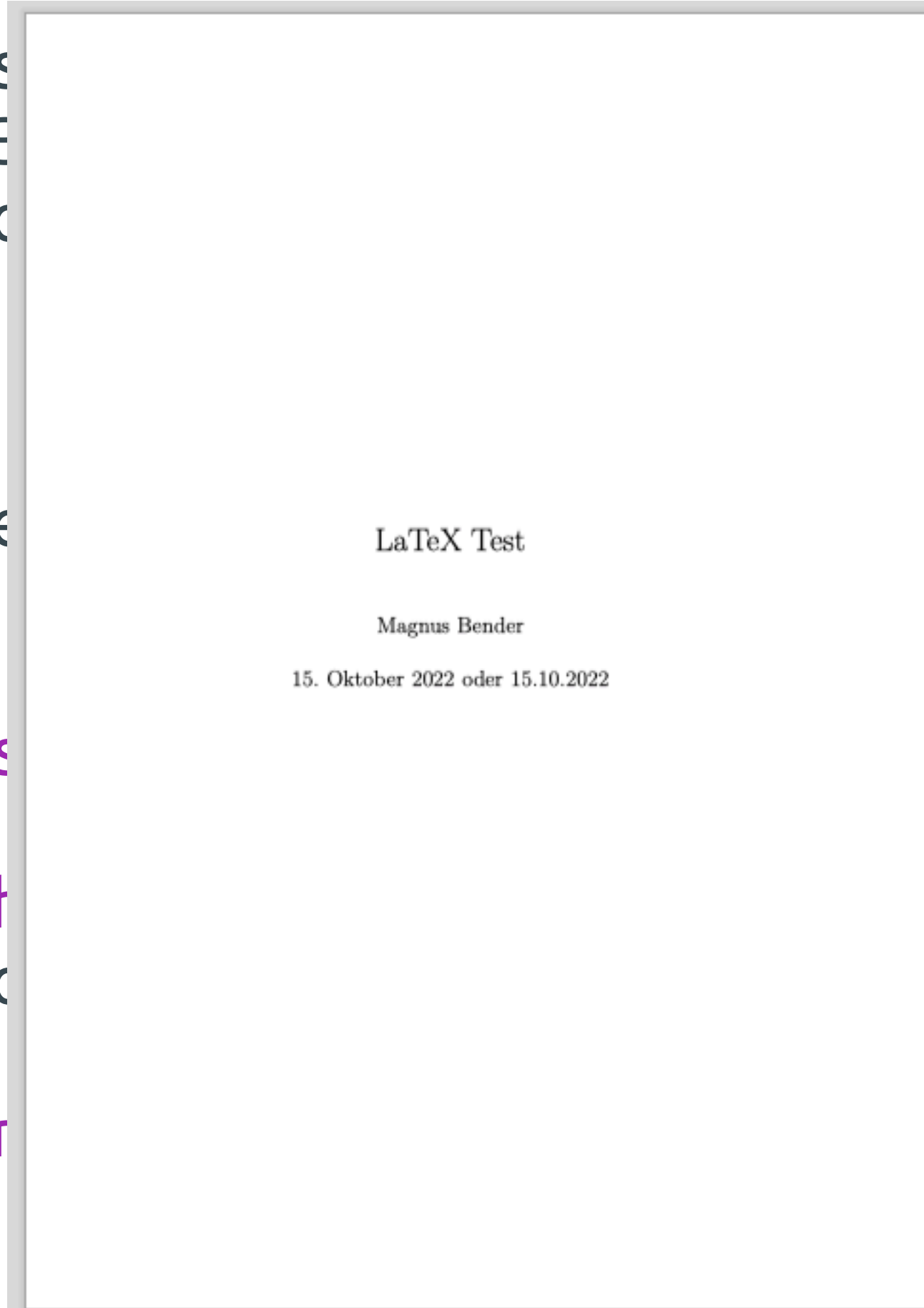
```
\begin{center}  
  Mein \textbf{e}  
\end{center}
```

```
\tableofcontents
```

```
\section{Überschrift}  
  Lorem ipsum do
```

```
\subsection{Unterüberschrift}  
  Lorem ipsum  
  At vero eos
```

% ...



Umwandlung von LaTeX

- Wie erstelle ich die PDF?
 - „`pdflatex ./document.tex`“
 - „`pdflatex ./document.tex`“ „`pdflatex ./document.tex`“
 - „`pdflatex ./document.tex`“ ...?
- LaTeX Mk
 - „`latexmk -pdf document.tex`“

L^AT_EX & Beamer

```
\documentclass{beamer}

% Dateikodierung
\usepackage[utf8]{inputenc}
% Trennung "deutsch"
\usepackage[ngerman]{babel}
% Schriftart
\usepackage[T1]{fontenc}
\usepackage{lmodern}

\title{LaTeX Test}
\author{Magnus Bender}
\date{\today{} oder 15.10.2022}

\usetheme{Luebeck}

\begin{document}

  \frame{\titlepage}

% ...
```

Präsentation, wieder andere Form, aber auch kleine Änderungen an der Struktur nötig!



L^AT_EX & Beamer

% ...

```
\begin{document}
```

```
\frame{\titlepage}
```

```
\begin{frame}
```

```
\begin{center}
```

```
Meine \textbf{erste} \alert{Präsentation} mit \LaTeX!
```

```
\end{center}
```

```
\end{frame}
```

```
\begin{frame}
```

```
\tableofcontents
```

```
\end{frame}
```

% ...



Überschrift
Überschrift 2

- 1 Überschrift
- 2 Überschrift 2

Magnus Bender LaTeX Test



Überschrift
Überschrift 2

Meine erste **Präsentation** mit \LaTeX !

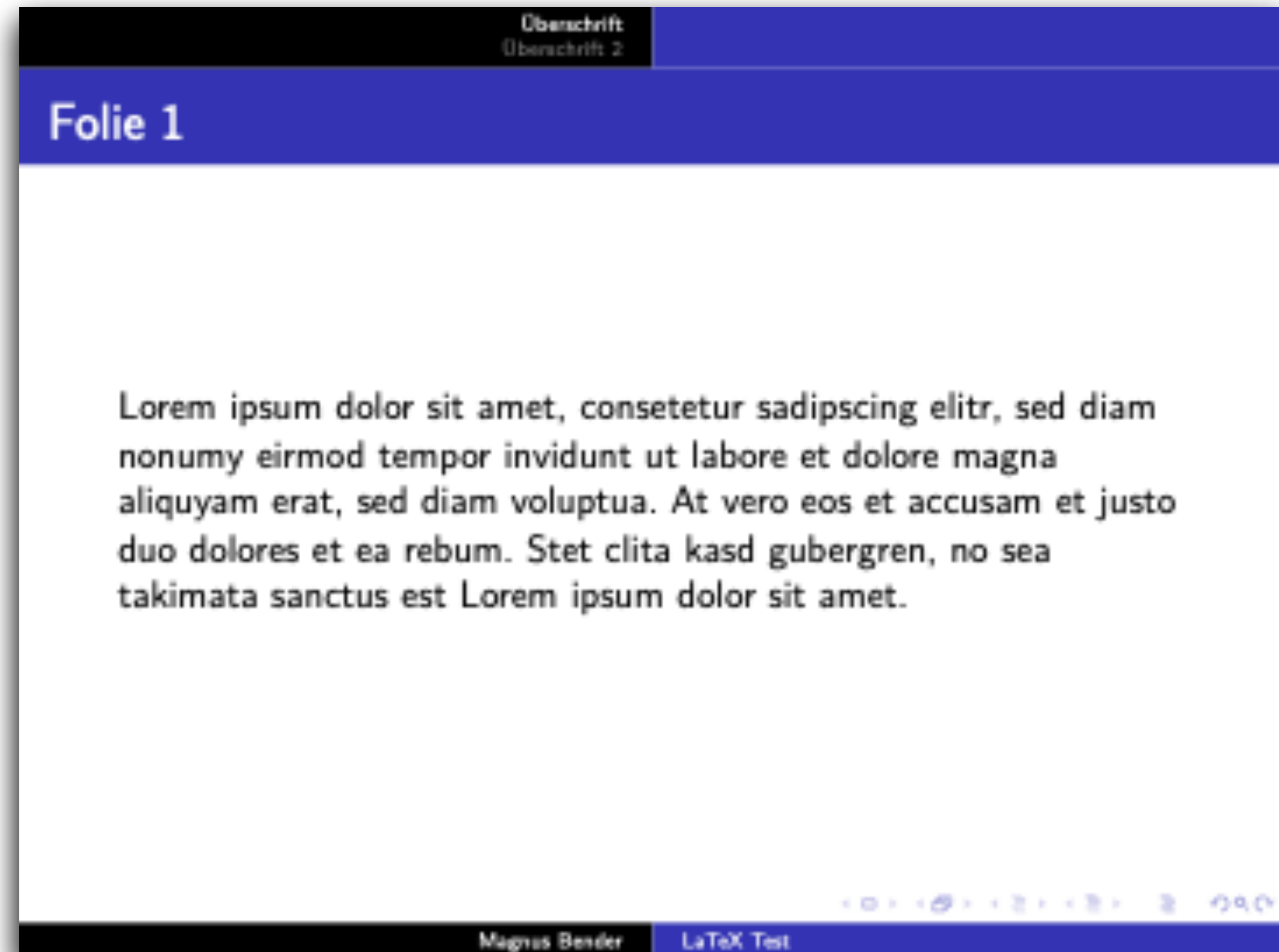
Magnus Bender LaTeX Test

L^AT_EX & Beamer

```
% ...  
  
\section{Überschrift}  
  
\begin{frame}{Folie 1}  
  Lorem ipsum dolor sit amet, ... \pause  
  At vero eos et accusam et ...  
\end{frame}
```

```
% ...
```

Folien mit `\pause` animieren,
erzeugt zwei Folien, eine nur mit
den Inhalten vor und auch mit
den Inhalten nach `\pause`.



L^AT_EX & Beamer

% ...

```
\section{Überschrift 2}
```

```
\begin{frame}{Folie 2}  
  \begin{block}{Hinweis}  
    Ein Text  
  \end{block}  
\end{frame}
```

```
\begin{alertblock}{Wichtig}<2->  
  Wieder ein Text  
\end{alertblock}  
\end{frame}
```

```
\section{Überschrift 3}
```

```
\end{document}
```

Es gibt verschiedene Blöcke in Beamer.

Neben `\pause` gibt es auch die `<x-y>` Syntax.



Fazit: Beamer

- *Beamer* erlaubt es Präsentationen mit LaTeX zu erstellen
 - Insbesondere ist die Umgebung `frame` neu, welche eine Folie darstellt
 - Verschiedene Themes für das Design der Folien
 - Weiterhin Trennung zwischen Form und Struktur

Form

Gegeben sei $x^2 + 5x - \alpha = 12$.

Wir nehmen an, dass die Gauß-Summe¹ $\sum_{i=1}^n i = \frac{n(n+1)}{2}$ bekannt ist.

Übrigens kann auch das Produkt verkürzt geschrieben werden:

```
\[
  y_1 \cdot \dots \cdot y_n
  = \prod_{i=1}^n y_i
\]
```

```
\begin{align*}
  3x^2 &+ 4x &= 0 & \\
  2x^2 &+ 10x &= 0 & \\
  4x^2 &+ &= 0 & \\
\end{align*}
```

Eine Übersicht über die verschiedenen Symbole, Klammern, etc. befindet sich z.B. hier: <http://tug.ctan.org/info/undergradmath/undergradmath.pdf>

Übrigens kann auch das Produkt verkürzt geschrieben werden:

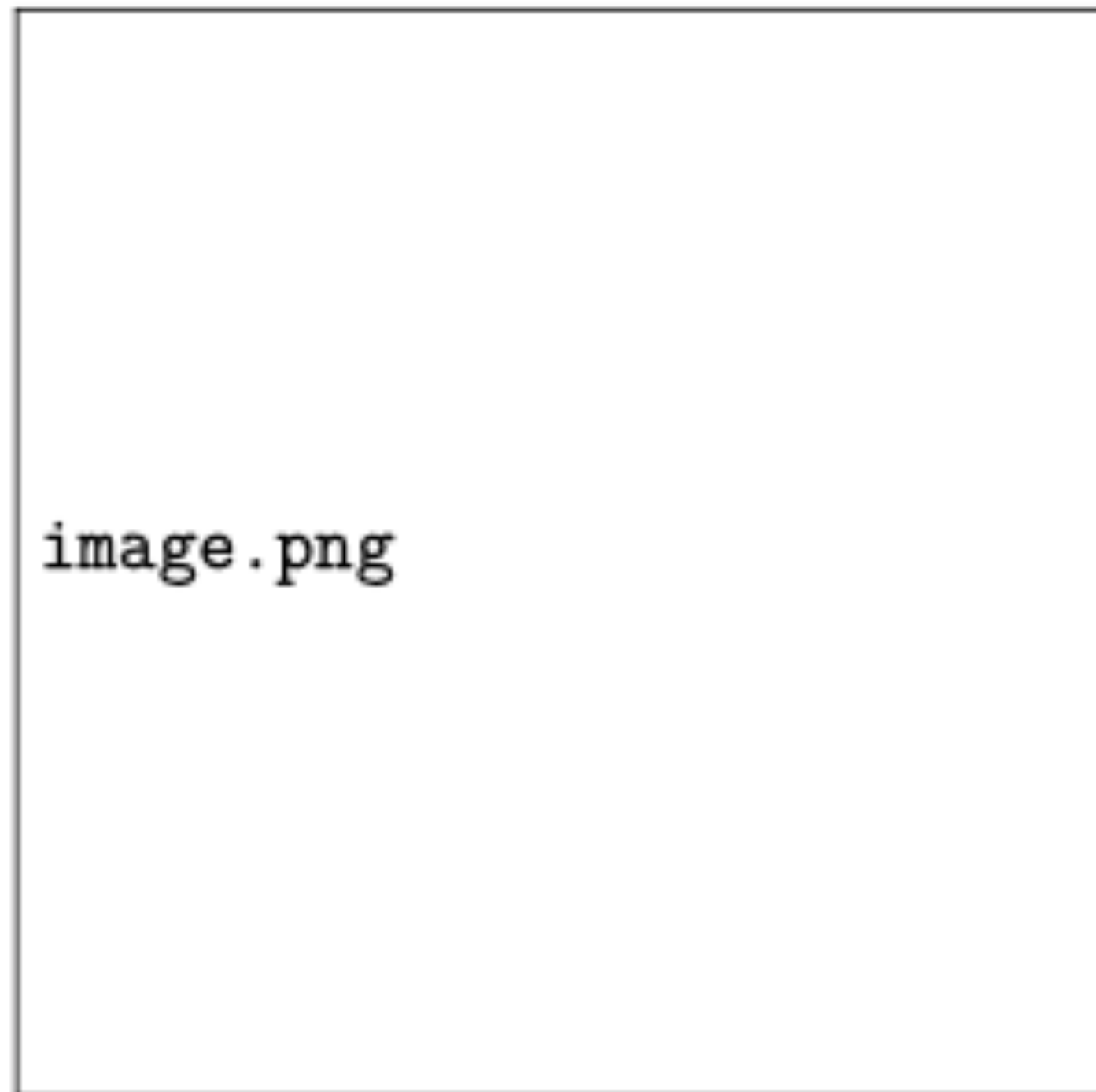
Die Formelsyntax von LaTeX wird auch an anderen Stellen genutzt, z.B. im Moodle oder in HedgeDoc.

¹siehe z.B. Wikipedia

Weitere Elemente in L^AT_EX

```
\begin{table}
\caption{Line}
\begin{table}
\caption{Line}
\end{table}
\end{table}
\begin{table}
\caption{Line}
\end{table}
\end{table}
```

Paket	Version	Anzahl
Numpy	1	12
Scipy	1.7	200
Gensim	2.4	30



i) Python

ii) Java

iii) LaTeX

- Auch Listen
- in Aufzählungen sind
- mit LaTeX
- möglich.

Ein Begriff der zu erklären ist.

Am Anfang wird der Begriff hervorgehoben und danach folgt der Text.

d danach folgt der Text.

Gleitumgeb

```
\begin{figure}  
  \centering  
  \includegraphics[width=\textwidth]{res/runtime_def}  
  \caption{  
    Time needed training the models for all scenarios  
  }  
  \label{fig:dur}  
\end{figure}
```

Besides the performance of all scenarios, also the runtime and the computational resources needed for training are relevant.

In `\ref{fig:dur}`, the duration for training each of the models is shown with a logarithmic scale.

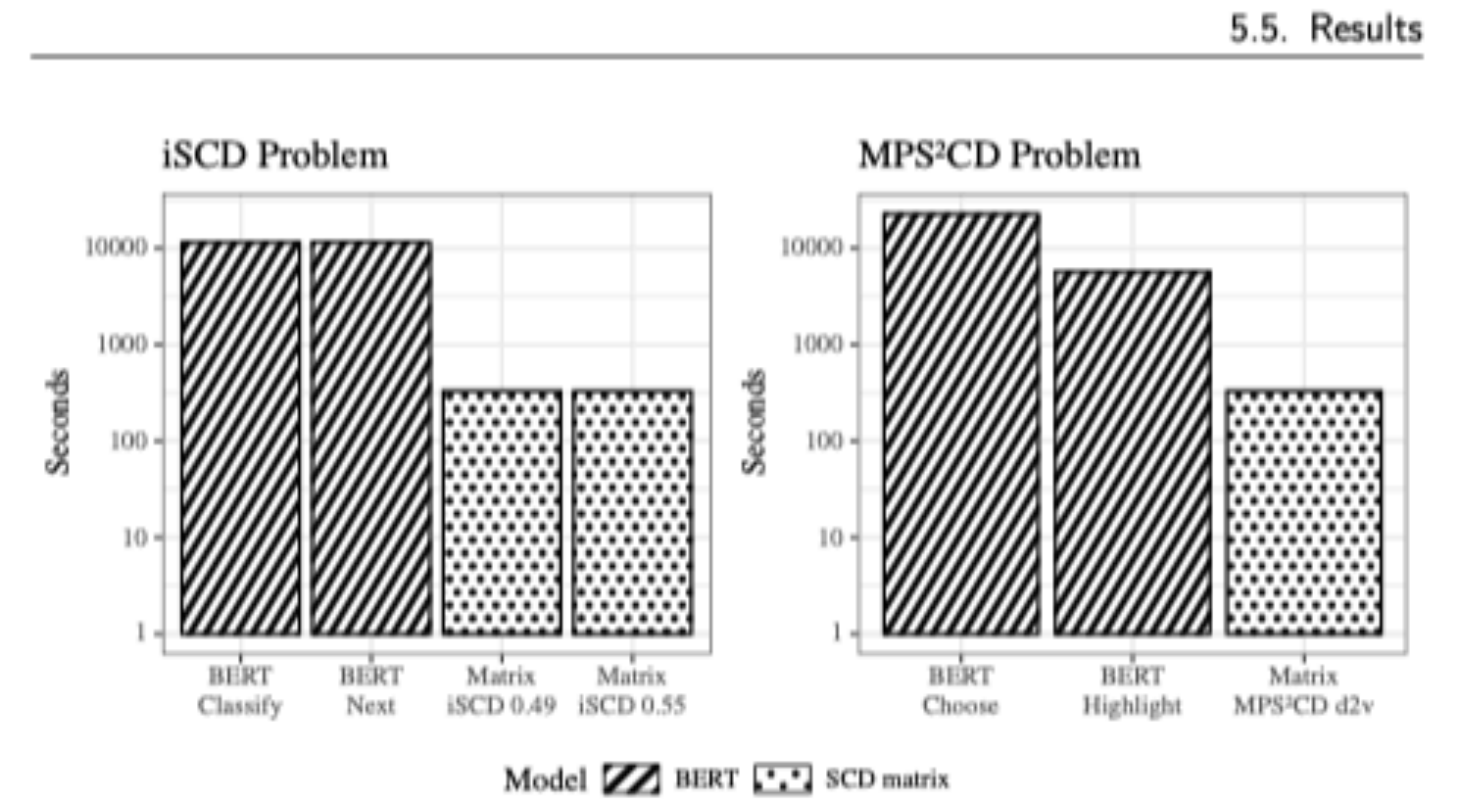


Figure 5.2.: Time needed training the models for all scenarios using the 20 newsgroups dataset and the Wiktionary annotation agent. There is no difference training the SCD matrix for Matrix MPS²CD ia or Matrix MPS²CD d2v.

Label und Referenz auf das Label.

In similar tasks. Only BERT Highlight with a disjoint set of SCDs achieves a very low accuracy. As BERT Highlight asks to highlight the matching SCD out of four SCDs, the accuracy of 0.25 is as worse as randomly highlighting an SCD. We simplify the problem for BERT Highlight and do not split the set of SCDs. Using BERT Highlight with the same set of SCD, then, shows a similar performance as BERT Choose and Matrix MPS²CD. However, for Matrix MPS²CD d2v we also have to use the same set of SCDs.

The best accuracy for the iSCD problem is yielded by BERT Next and for the MPS²CD problem by Matrix MPS²CD ia.

Besides the performance of all scenarios, also the runtime and the computational resources needed for training are relevant. In Figure 5.2, the duration for training each of the models is shown with a logarithmic scale. The training time of an SCD

Warum ist die Grafik im PDF an einer anderen Stelle?

BibTeX und TikZ

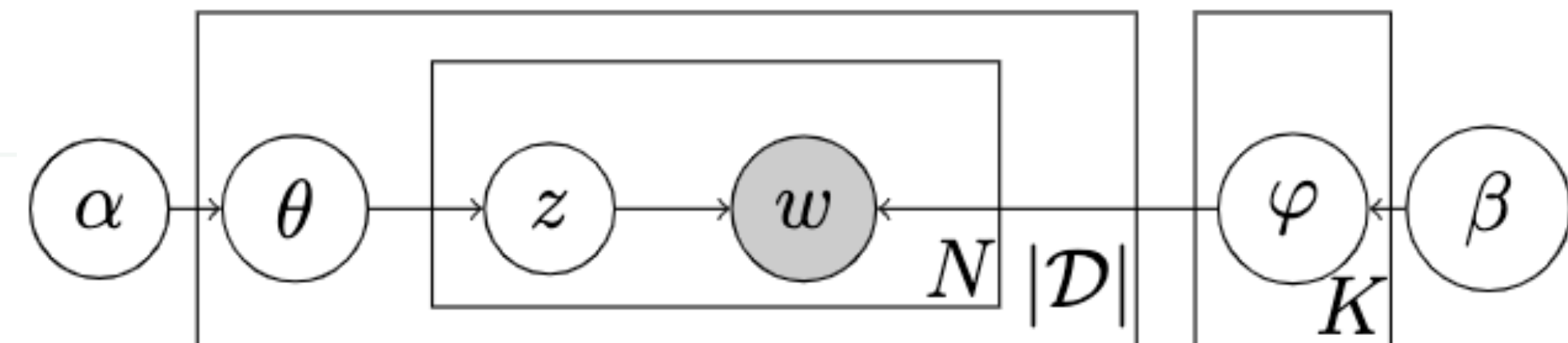
- BibTeX
- Literaturverzeichnisse und -zitate mit LaTeX
- TikZ ist kein Zeichenprogramm
- Grafiken und Abbildungen direkt in LaTeX erstellen

document.tex

```
Open Information Extraction
(OpenIE)~\cite{OpenIE}
extracts triples of
subject, predicate and
object.
```

literature.bib

```
@article{OpenIE,
  title = "Leveraging Linguistic Structure for
  author = "Angeli, Gabor and Johnson Pr
  journal = "Proceedings of the Associat
  year = "2015",
  publisher = "Association for Computati
  115/v1/P15-
```



Information Extraction. In: *Proceedings of the Association of Computational Linguistics (ACL)* (2015), 344–354. <https://doi.org/10.3115/v1/P15-1034>

[All83] ALLEN, James F.: Maintaining Knowledge about Temporal Intervals. In: *Commun. ACM* 26 (1983), November, Nr. 11, 832–843. <https://doi.org/10.1145/182.358434>

[BBG⁺21a] BENDER, Magnus ; BRAUN, Tanya ; GEHRKE, Marcel ; KUHR, Felix ; MÖLLER, Ralf ; SCHIFF, Simon: Identifying and Translating Subjective Content Descriptions Among Texts. In: *Int. J. Semantic Computing* 15 (2021). – Accepted for publication

[BBG⁺21b] BENDER, Magnus ; BRAUN, Tanya ; GEHRKE, Marcel ; KUHR, Felix ; MÖLLER, Ralf ; SCHIFF, Simon: Identifying Subjective Content Descriptions among Text. In: *Proceedings of the 15th IEEE International Conference on Semantic Computing (ICSC-21)* (2021). <https://doi.org/10.1109/ICSC50631.2021.00008>

[BMR⁺20] BROWN, Tom B. ; MANN, Benjamin ; RYDER, Nick ; SUBBIAH, Melanie ; KAPLAN, Jared ; DHARIWAL, Prafulla ; NEELAKANTAN, Arvind ; SHYAM, Pranav ; SASTRY, Girish ; ASKELL, Amanda ; AGARWAL, Sandhini ; HERBERT-VOSS, Ariel ; KRUEGER, Gretchen ; HENIGHAN, Tom ; CHILD, Rewon ; RAMESH, Aditya ; ZIEGLER, Daniel M. ; WU,

Links & Tools $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$

- <https://www.mlte.de/latex/>
 - Video und Unterlagen zu Vorlesungen über LaTeX, Beamer & TikZ
 - Klassen für Übungszettel
 - Vorlagen für Präsentationen und Abschlussarbeiten
- <http://en.wikibooks.org/wiki/LaTeX>
 - Wikibook
- <http://detexify.kirelabs.org/>
 - Sonderzeichen bestimmen
- <https://www.overleaf.com/>
 - Online und ohne Installation Dokumente in LaTeX verfassen
- <https://app.diagrams.net/>
 - Diagramme und Grafiken schnell und einfach erstellen
(kein Bezug zu LaTeX!)

Übungszettel mit L^AT_EX

- Erstellung eines Übungszettels mit LaTeX
- Wichtigste Pakete in der Präambel
- Name, Seitenzahl, Gruppe auf jeder Seite
- Alternativ Nutzung einer vorgefertigten Klasse, z.B.
- <https://www.mlte.de/latex/exercise-class/>



Live Demo

Zusammenfassung

- Auszeichnungssprachen
- Markdown
- LaTeX
 - Grundlagen
 - Beamer
 - Formeln
 - Beispiel: Übungszettel



Inhaltsübersicht

1. Programmiersprache Python
 - a) *Einführung, Erste Schritte*
 - b) *Grundlagen*
 - c) *Fortgeschritten*
2. Auszeichnungssprachen
 - a) *LaTeX, Markdown*
3. Benutzeroberflächen und Entwicklungsumgebungen
 - a) **Jupyter Notebooks lokal und in der Cloud (Google Colab)**
4. Versionsverwaltung
 - a) Git, GitHub
5. Wissenschaftliches Rechnen
 - a) NumPy, SciPy
6. Datenverarbeitung und -visualisierung
 - a) Pandas, matplotlib, NLTK
7. Machine Learning (scikit-learn)
 - a) Grundlegende Ansätze (Datensätze, Auswertung)
 - b) Einfache Verfahren (Clustering, ...)
8. DeepLearning
 - a) TensorFlow, PyTorch, HuggingFace Transformers

Anhang: Präambel Übungszettel

```
\documentclass[
  oneseide,      % Einzelne Seiten
  12pt,         % Schriftgroesse
  a4paper,      % Papier
  parskip=full  % Absatzstil
]{scrartcl}

% Mathematische Symbole, Umgebungen
\usepackage{amssymb, amsfonts, amsthm, amsmath}

% Auflistungen und Grafiken
\usepackage{paralist, graphicx}

% Dateikodierung
\usepackage[utf8]{inputenc}

% Fuss- und Kopfzeile
\usepackage[headsepline,footsepline]{scrlayer-scrpage}
\makeatletter
  \clearpaïrofpagestyles
  \ifoot{\@subtitle}
  \ofoot{\@title}
  \ihead{\@author}
  \ohead{\pagemark}
\makeatother

% Trennung "deutsch"
\usepackage[ngerman]{babel}

% Schriftart
\usepackage[T1]{fontenc}
\usepackage{lmodern}

% Anfuhrungszeichen
\usepackage[german=quotes]{csquotes}
```

Anhang: Inhalt Übungsblatt

```
% Angaben ueber Aufgabenblatt und Autor
\author{Magnus Bender (LaTeX Beispiel)}
\title{Einführung in Web and Data Science\\Übungsblatt 1}
\subtitle{Gruppe 5, Übung Mo. 10 Uhr}

\begin{document}
% Titel
% Kann man auch einfach auskommentieren, dann spart man
% den Platz auf der ersten Seite!
\maketitle

% Ueberschrift ohne Laufende Nummer
\section*{Aufgabe 1}
% Lösungen, ...

\section*{Aufgabe 2}
% Lösungen, ...

\end{document}
```