

Scraping *derstandard.at*

TODO: talk about some legal stuff with scraping

First we define the function `get_standard_soup()`, which sends a request to derstandard, along with a cookie that derstandard checks if a user has accepted their DSGVO notice. If this cookie is not sent, a banner is displayed and the html is only partially loaded.

Secondly, we define `get_frontpage_articles()`, which expects a bs4 soup object of a frontpage. While derstandard no longer offers an archive, the frontpage articles of a given day can be conveniently accessed with the pattern `frontpage/y/m/d`. Each frontpage contains article sections which contain the (sub)heading, lead, number of comments, storylabels etc.

We will start by pulling the frontpage of december 22th 2023.

```
In [17]: from bs4 import BeautifulSoup
import requests

# fetch the html content of a derstandard.at page
def get_standard_soup(link):
    response = requests.get(link, cookies={'DSGVO_ZUSAGE_V1': 'true'})
    return BeautifulSoup(response.content, 'html.parser')

# generate a dictionary of articles with title as key and the bs4 element as value
def get_frontpage_articles(soup):
    articles_dict = {}
    articles = soup.select('div.chronological>section article')
    for article in articles:
        title_tag = article.find('a')
        if title_tag and title_tag.has_attr('title'):
            title = title_tag['title']
            articles_dict[title] = article
    return articles_dict

# Generate the articles dictionary for an arbitrary frontpage
soup = get_standard_soup('https://www.derstandard.at/frontpage/2023/12/22')
articles_dict = get_frontpage_articles(soup)

print(f'We have fetched {len(articles_dict)} articles\n')
```

We have fetched 137 articles

In the next step, let's look at the information we can get from those article sections on the frontpage. By inspecting the html, we have already identified various elements that we will use in the subsequent steps:

- title
- subtitle
- article type
- link
- datetime
- kicker (like an additional tag, not 100% about its meaning yet)
- postingcount
- storylabels

while playing with the data, we noticed that not every article contains storylabels. We will check this in the following step, as well as if every article tag has a type.

```
In [18]: # Function to analyze attributes of specified tags and their attributes
def analyze_tag_attributes(articles_dict):
    no_data_type = set()
    no_story_label = set()

    for title, article in articles_dict.items():
        # Check if every article tag has a data-type attribute - basically the type of the article
        if not article.has_attr('data-type'):
            no_data_type.add(title)
        # search for <div class="storyLabels"> in articles - the story labels
        if not article.find('div', class_='storylabels'):
            no_story_label.add(title)

    return no_data_type, no_story_label

no_data_type, no_story_label = analyze_tag_attributes(articles_dict)
print(f'Number of articles without data-type attribute: {len(no_data_type)}')
print(f'Number of articles without storylabels: {len(no_story_label)}')
```

```
# get articles that have a story label
has_label = set(articles_dict.keys()).difference(no_story_label)
print(f'Number of articles with story attribute: {len(has_label)}')

# a lot of articles do not have a story label, maybe an interesting goal for machine learning
```

Number of articles without data-type attribute: 0

Number of articles without storylabels: 104

Number of articles with story attribute: 33

All articles have a data-type, but only a few articles have story attributes. This could be an interesting labeling task for our machine learning project later. Next, we print out the html of two articles to show the data we are interested in.

```
In [19]: # example of an article without story label
print(f'No storylabel:\n{articles_dict[list(no_story_label)[0]]}\n')
print(80*'-')

# articles with story label
print(f'With storylabel:\n{articles_dict[list(has_label)[0]]}')
```

```
No storylabel:
<article class="fig" data-dg="p1-43" data-dt="7x2" data-mg="p1-43" data-mt="4x4" data-type="story">
<div class="teaser-inner">
<a href="/story/3000000200853/trump-uebte-laut-medienbericht-druck-auf-wahlpruefer-in-michigan-aus" title="Trump übte laut
Medienbericht Druck auf Wahlprüfer in Michigan aus">
<figure data-type="image">
<picture>
<source data-lazy-srcset="https://i.ds.at/V3WkIw/c:1200:800:fp:0.500:0.500/rs:fill:280:187/g:fp:0.54:0.29/plain/lido-image
s/2023/12/22/3f9fca60-5573-4475-8e22-0cd35d00484b.jpeg" media="(min-width: 960px)"/>
<source data-lazy-srcset="https://i.ds.at/KnrHlA/c:1200:800:fp:0.500:0.500/rs:fill:750:375/g:fp:0.54:0.29/plain/lido-image
s/2023/12/22/3f9fca60-5573-4475-8e22-0cd35d00484b.jpeg" media="(max-width: 959px)"/>

</picture>
</figure>
<header>
<time datetime="2023-12-22T14:38
">14:38
</time>
<p class="teaser-kicker">USA</p>
<div class="teaser-postingcount">25 <span>Postings</span>
</div>
<h1 class="teaser-title">Trump übte laut Medienbericht Druck auf Wahlprüfer in Michigan aus </h1>
<p class="teaser-subtitle">Der <span class="hyphenate">Ex-US-Präsident</span> soll versucht haben, zwei Wahlprüfer von ein
er Bestätigung der Wahlresultate abzuhalten </p>
</header>
</a>
</div>
</article>
```

With storylabel:

```
<article class="fig" data-dg="p1-132" data-dt="7x2" data-mg="p1-132" data-mt="4x4" data-type="story">
<div class="teaser-inner">
<a href="/story/3000000200661/baerige-weihnachten-in-schoenbrunn" title="Bärige Weihnachten in Schönbrunn">
<figure data-type="image">
<picture>
<source data-lazy-srcset="https://i.ds.at/oT-PeQ/c:1620:1080:fp:0.500:0.500/rs:fill:280:187/plain/lido-images/2023/12/21/0
dc7ec6e-fd95-4be8-94ca-801ed067dfb6.jpeg" media="(min-width: 960px)"/>
<source data-lazy-srcset="https://i.ds.at/xmQSZa/c:1620:1080:fp:0.500:0.500/rs:fill:750:375/plain/lido-images/2023/12/21/0
dc7ec6e-fd95-4be8-94ca-801ed067dfb6.jpeg" media="(max-width: 959px)"/>

</picture>
<time class="duration">01:03</time>
</figure>
<header>
<time datetime="2023-12-22T06:00
">06:00
</time>
<p class="teaser-kicker">Weihnachten</p>
<div class="teaser-postingcount">1 <span>Posting</span>
</div>
<h1 class="teaser-title">Bärige Weihnachten in Schönbrunn </h1>
<p class="teaser-subtitle">Für die Schönbrunner Brillenbären hat es bereits vor den Feiertagen eine weihnachtliche Überras
chung vom <span class="hyphenate">Tierpflegerteam</span> gegeben: einen Christbaum geschmückt unter anderem mit <span clas
s="hyphenate">Karottenlametta</span> und <span class="hyphenate">Paprikaringerln</span> </p>
<time class="duration">01:03</time>
<div class="storylabels">
<span data-label-category="indepth" data-label-name="Video">Video</span>
</div>
</header>
</a>
</div>
</article>
```

This is the information one can get from the frontpage. Later we will also follow the links and scrape additional data from the articles, but lets first focus on getting a dataset just based on the information that can be attained from the frontpage.

To this end, we define `extract_article_data()` , which uses the dictionary following the pattern `article_title:` `article_section_soup` . From the html (soup), we will extract and clean:

- title
- teaser-subtitle
- link
- time
- teaser-kicker
- n_posts
- storylabels

In [20]:

```
# Function to extract specific data from each article
def extract_article_data(articles_dict):
    HOST = 'https://www.derstandard.at'
    article_data = []

    for title, article in articles_dict.items():
        data = {
            'title': title,
            'teaser-subtitle': None,
            'link': None,
            'time': None,
            'teaser-kicker': None,
            'n_posts': None,
            'storylabels': None
        }

        # most links are relative, so we need to add the host
        link = article.find('a')['href']
        if not link.startswith(HOST):
            link = HOST + link
        data['link'] = link

        # for live articles, there is a second time tag with the duration of the live post
        # however, we only care about the time of publication here
        time = [tag for tag in article.find_all('time') if 'datetime' in tag.attrs][0]
        data['time'] = time['datetime'].rstrip('\r\n')

        # if there are no comments, the string is empty so set it to 0
        n_posts = article.find('div', 'teaser-postingcount')
        try: data['n_posts'] = int(n_posts.get_text(strip=True).rstrip('Posting').replace('.', ''))
        except: data['n_posts'] = 0

        # Extracting other specified tags
        for tag, class_name in [('p', 'teaser-kicker'),
                                ('p', 'teaser-subtitle'),
                                ('div', 'storylabels')]:
            found_tag = article.find(tag, class_=class_name)
            if found_tag:
                data[class_name] = found_tag.get_text(strip=True)

        article_data.append(data)

    return article_data

article_data = extract_article_data(articles_dict)
# last 5 articles, of which some have a story label
article_data[-5:]
```

```
Out[20]: [{ 'title': 'Jesus-Geburt unter Palmen',
  'teaser-subtitle': 'Auch der Koran erzählt über die Geburt von Jesus Christus. Nicht als Sohn Gottes, aber als Sohn Marias, die als Frau im Koran eine besondere Stellung einnimmt',
  'link': 'https://www.derstandard.at/story/3000000200744/jesus-geburt-unter-palmen',
  'time': '2023-12-22T06:00',
  'teaser-kicker': 'Wussten Sie schon?',
  'n_posts': 1,
  'storylabels': None},
{'title': '"Aquaman and the Lost Kingdom" scheitert an versuchter Schadensbegrenzung',
  'teaser-subtitle': 'Der Erfolg der Superheldenfilme leidet immer öfter unter den privaten Problemen ihrer Hauptdarsteller. Der Fortsetzung von "Aquaman" droht auch deshalb ein Flop',
  'link': 'https://www.derstandard.at/story/3000000200724/aquaman-and-the-lost-kingdom-scheitert-an-versuchter-schadensbegrenzung',
  'time': '2023-12-22T06:00',
  'teaser-kicker': 'Im Kino',
  'n_posts': 242,
  'storylabels': None},
{'title': 'One-Man-Show mit fatalen Folgen für die Demokratie in Serbien',
  'teaser-subtitle': 'Die Wahlen in Serbien waren eine Farce. Die Europäische Union hat bisher auf einen pragmatischen Kurschwechsel gesetzt. Damit könnte es jetzt aber vorbei sein. Für eine härtere Gangart wäre es auch höchste Zeit',
  'link': 'https://www.derstandard.at/story/3000000200698/one-man-show-mit-fatalen-folgen-fuer-die-demokratie-in-serbien',
  'time': '2023-12-22T06:00',
  'teaser-kicker': 'Vedran Džihic',
  'n_posts': 112,
  'storylabels': 'Kommentar der anderen'},
{'title': 'David Alaba zum zehnten Mal zu Österreichs Fußballer des Jahres gekürt',
  'teaser-subtitle': 'Zehn von zwölf Trainern wählten den derzeit verletzten Real-Verteidiger auf Rang eins. Für den Wiener ist es der vierte Erfolg in Serie',
  'link': 'https://www.derstandard.at/story/3000000200745/david-alaba-zum-10-mal-fuessballer-des-jahres-riesige-ehre',
  'time': '2023-12-22T05:46',
  'teaser-kicker': 'Fußball',
  'n_posts': 33,
  'storylabels': None},
{'title': 'Kreuzworträtsel F 10571',
  'teaser-subtitle': 'Täglich neu, exklusiv für Smart-Abonnent:innen: Das knifflige Phoenixen-Rätsel des STANDARD',
  'link': 'https://www.derstandard.at/story/3000000199163/kreuzwortraetsel-f-10571',
  'time': '2023-12-22T00:01',
  'teaser-kicker': 'Kreuzworträtsel',
  'n_posts': 4,
  'storylabels': 'Spiel'}}
```

The various attributes will be analyzed once we convert our data to a dataframe.

But before we start scraping like mad, let's check robots.txt such that we can comply with derstandard's scraping policies.

```
In [21]: print(get_standard_soup('https://www.derstandard.at/robots.txt'))
```

User-agent: *

Disallow: /profil/

Sitemap: https://www.derstandard.at/sitemaps/news.xml

Sitemap: https://www.derstandard.at/sitemaps/sitemap.xml

Crawl-delay: 1

C:\Users\Paul\AppData\Local\Temp\ipykernel_24276\4133154167.py:7: MarkupResemblesLocatorWarning: The input looks more like a filename than markup. You may want to open this file and pass the filehandle into BeautifulSoup.
return BeautifulSoup(response.content, 'html.parser')

Crawl-delay: 1, so let's be nice and wait 1 second between requests. Then we'll scrape the frontpage of every day in 2023 until the 20th of december and save the data as a csv.

Caution: you might not want to run this cell, as it takes about ~45 minutes to run. The data has already been extracted once and has been saved to `data/derstandard_frontpage_data.csv`.

```
In [28]: from datetime import datetime, timedelta
from time import sleep
from tqdm import tqdm

def scrape_frontpage(start_date: str, end_date: str, logging=False):
    # Validate that dates follow the pattern YYYY-MM-DD
    try:
        start = datetime.strptime(start_date, '%Y-%m-%d')
        end = datetime.strptime(end_date, '%Y-%m-%d')
    except ValueError:
        print("Invalid date format. Please use YYYY-MM-DD.")
        return
```

```

data = []
# all dates between start and end (inclusive)
delta = end - start
for i in tqdm(range(delta.days + 1)):
    # generate link for each day
    day = start + timedelta(days=i)
    date = day.strftime('%Y/%m/%d')
    link = f'https://www.derstandard.at/frontpage/{date}'
    # make a request to the link and extract the data
    article_dict = get_frontpage_articles(get_standard_soup(link))
    articles = extract_article_data(article_dict)
    if logging:
        print(f'Fetched {len(articles)} articles from {date}')
    data += articles
    # wait almost a second before next request, our data processing takes a bit of time as well
    sleep(0.8)

return data

# scrape the data for 4 years
data = scrape_frontpage('2019-12-22', '2023-12-22')

```

100%|██████████| 1462/1462 [46:44<00:00, 1.92s/it]

Okay, this cell took a while to run obviously, but we finally have our precious data. Lets convert it to a dataframe and see what we have.

```

In [29]: import pandas as pd

df = pd.DataFrame(data)
df.columns = df.columns.str.replace('teaser-', '')
df.rename(columns={'time': 'datetime'}, inplace=True)
df

```

Out[29]:

| | title | subtitle | link | datetime | kicker |
|--------|---|---|---|------------------|--------------------|
| 0 | Real Madrid stolpert mit Aluminiumpech im Tite... | Die Königlichen können Bilbao daheim nicht bes... | https://www.derstandard.at/story/2000112599363... | 2019-12-22T23:44 | Primera Division |
| 1 | Bolivien weist venezolanische Diplomaten aus | InterimspräsidentinJeanine Áñez wirft denBotsc... | https://www.derstandard.at/story/2000112598924... | 2019-12-22T22:50 | Übergangsregierung |
| 2 | Erdoğan warnt vor neuer Flüchtlingswelle aus S... | Türkischer Präsident: "80.000 Menschen Richtun... | https://www.derstandard.at/story/2000112598130... | 2019-12-22T21:43 | Bürgerkrieg |
| 3 | Massenkarambolage mit 63 Fahrzeugen in Virginia | Autos stießen auf vereister Brücke zusammen | https://www.derstandard.at/story/2000112597972... | 2019-12-22T21:29 | Weihnachtsverkehr |
| 4 | Salzburg schlägt Caps, Meister KAC mit vierter... | Die Bullen sind damit der Gewinner der Runde: ... | https://www.derstandard.at/story/2000112595206... | 2019-12-22T20:54 | Eishockey |
| ... | ... | ... | ... | ... | ... |
| 182102 | Wer braucht die Kirche? | Dass sich die Kirche nach soschwerwiegendenVer... | https://www.derstandard.at/story/3000000200743... | 2023-12-22T06:00 | Dominik Straub |
| 182103 | Sonderregelung verlängert: Mehr als 1.000 Ärzt... | Der"Pandemieparagraf"im Ärztegesetz hat mehr a... | https://www.derstandard.at/story/3000000200621... | 2023-12-22T06:00 | Pandemieparagraf |
| 182104 | Stadtforscher: "Architektur ist Teil unserer W... | Jetzt anhören: In Zukunft müssen Städte wieder... | https://www.derstandard.at/story/3000000200499... | 2023-12-22T06:00 | Edition Zukunft |
| 182105 | David Alaba zum zehnten Mal zu Österreichs Fuß... | Zehn von zwölf Trainern wählen den derzeit ve... | https://www.derstandard.at/story/3000000200745... | 2023-12-22T05:46 | Fußball |
| 182106 | Kreuzworträtsel F 10571 | Täglich neu, exklusiv fürSmart-Abonent:innen:... | https://www.derstandard.at/story/3000000199163... | 2023-12-22T00:01 | Kreuzworträtsel |

182107 rows × 7 columns



over 182 thousand rows, this should give us plenty data to analyze!

Lets save it to a csv, totalling 57mb

```
In [30]: df.to_csv('data/4yrs_derstandard_frontpage_data.csv', index=False)
```

Data Processing Notebook

EDA

Lets load the frontpage data of 4 years going from december 22. 2019 to december 22. 2023.

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

df = pd.read_csv('./data/4yrs_derstandard_frontpage_data.csv')

# prepare time data
df['datetime'] = pd.to_datetime(df['datetime'])
df.head()
```

```
Out[1]:
```

| | title | subtitle | link | datetime | kicker | n_post |
|---|---|---|---|---------------------|--------------------|--------|
| 0 | Real Madrid stolpert mit Aluminiumpech im Tite... | Die Königlichen können Bilbao daheim nicht bes... | https://www.derstandard.at/story/2000112599363... | 2019-12-22 23:44:00 | Primera Division | 3 |
| 1 | Bolivien weist venezolanische Diplomaten aus | InterimspräsidentinJeanine Áñez wirft denBotsc... | https://www.derstandard.at/story/2000112598924... | 2019-12-22 22:50:00 | Übergangsregierung | 1 |
| 2 | Erdoğan warnt vor neuer Flüchtlingswelle aus S... | Türkischer Präsident: "80.000 Menschen Richtun... | https://www.derstandard.at/story/2000112598130... | 2019-12-22 21:43:00 | Bürgerkrieg | 10 |
| 3 | Massenkarambolage mit 63 Fahrzeugen in Virginia | Autos stießen auf vereister Brücke zusammen | https://www.derstandard.at/story/2000112597972... | 2019-12-22 21:29:00 | Weihnachtsverkehr | 3 |
| 4 | Salzburg schlägt Caps, Meister KAC mit vierer... | Die Bullen sind damit der Gewinner der Runde: ... | https://www.derstandard.at/story/2000112595206... | 2019-12-22 20:54:00 | Eishockey | |

182 thousand rows, lets inspect the data a bit

```
In [2]: print('All variables:\n', df.dtypes, '\n')
print('N_posts percentiles:\n', df['n_posts'].describe(), '\n')
print('n_posts 0-values:\n', len(df[df['n_posts']==0]))
```

```
All variables:
title           object
subtitle        object
link            object
datetime        datetime64[ns]
kicker          object
n_posts         int64
storylabels     object
dtype: object
```

```
N_posts percentiles:
count    1.821070e+05
mean     3.342677e+02
std      7.621049e+03
min      0.000000e+00
25%      1.400000e+01
50%      7.200000e+01
75%      2.490000e+02
max      3.000166e+06
Name: n_posts, dtype: float64
```

```
n_posts 0-values:
10417
```

Data correctness

While manually analyzing a sample of the over 10.000 articles with 0 posts, We noticed that there were a lot of articles that generated 0 posts. In the scraping script, we gathered the post count in the function `extract_article_data()` , which searched for the span containing the post count. Sometimes 0 posts is explicitly mentioned, sometimes the span is empty.

We analyzed a sample our data on articles with 0 posts and found that the information we could get from the frontpage seems to be correct and that the articles with 0 posts did indeed generate 0 posts. With one exception - the data on Live tickers was sometimes incorrect and indicated 0 posts when there were many posts in actuality.

The cause for this mismatch is unclear to us, but we were able to find a pattern which helps us find those wrongly labeled discussion forums.

```
In [3]: print('posts which contain ticker in the title:')
print(df[df['title'].str.contains('Ticker')]['storylabels'].value_counts(dropna=False), '\n')
print('posts which contain live in the title:')
print(df[df['title'].str.contains('Live')]['storylabels'].value_counts(dropna=False), '\n')
```

posts which contain ticker in the title:

```
storylabels
NachleseLiveticker      271
LivetickerNachlese      30
LiveLiveticker          6
NaN                     6
NachleseLivebericht     5
Forum+Nachlese          3
LivetickerLive          2
LiveberichtNachlese     2
NachleseForum+          1
LiveForum+              1
LiveLivebericht         1
Forum+Live              1
Name: count, dtype: int64
```

posts which contain live in the title:

```
storylabels
NaN                    502
NachleseLiveticker     64
NachleseLivebericht    12
LivetickerNachlese     9
Kolumne                7
Interview              5
LiveberichtNachlese    4
Video                  3
Analyse                3
Blog                   3
Podcast                2
Porträt                2
Bericht                2
Ansichtssache          2
Essay                  1
LiveLiveticker         1
Kopf des Tages         1
Kommentar              1
KolumneVideo           1
LivetickerLive         1
Rezension              1
Reportage              1
Name: count, dtype: int64
```

While there are still some NA values in the set of posts that contain the words Live or Ticker, closer inspection has shown that those articles were not actually discussion forums/live tickers but rather articles like **Livestreams, Ticker und Spielpläne: Wie man die EM 2021 online verfolgt**.

But as it turns out, livetickers seem to be among the rare data points that have all been assigned a storylabel. We will thus first identify the storylabels that live articles have been tagged with.

```
In [4]: # all the Live storylabels (I checked, they are all capitalized)
live_labels = df[(df['storylabels'].str.contains('Live')) & (~df['storylabels'].isna())]['storylabels'].value_counts()
live_labels = list(live_labels.index)
print('All Live storylabels:\n*', '\n* '.join(live_labels))
```

```

All Live storylabels:
* NachleseLivebericht
* NachleseLiveticker
* LiveberichtNachlese
* LivetickerNachlese
* LiveForum+
* LiveLiveticker
* Forum+Live
* LiveLivebericht
* LivetickerLive
* LiveberichtLive
* LiveberichtVorschau
* VorschauLiveticker
* Liveticker
* LiveberichtVideo
* Livebericht

```

We will now use this list of storylabels and simply drop all articles tagged with them that contain 0 posts. Of course this approach is just a heuristic and there is a very good chance there are still wrongly labeled datapoints, but for the purposes of this project we will leave it at that. Those >900 articles are very likely to be wrongly labeled, and making another request to fetch the actual comment count would be too time-consuming for the purposes of this project.

```

In [5]: forums = df[df['storylabels'].isin(live_labels)]
no_post_forums = forums[forums['n_posts'] == 0]
print(f'Forums wrongly labeled as 0 posts: {len(no_post_forums)} out of {len(forums)}')

# drop all rows from the df that are in no_post_forums
df.drop(no_post_forums.index, inplace=True)
print(f'Dropped {len(no_post_forums)} rows. {len(df)} rows remaining.')

```

Forums wrongly labeled as 0 posts: 907 out of 2422
Dropped 907 rows. 181200 rows remaining.

Post count distribution

Next we will look at the distribution of our intended target variable, `n_posts`. For a first look, we bin the articles in ranges of 10. The following scatterplot shows the number of articles that generate between n and $n+10$ posts on the y axis and the lower bound of the bin (n) of posts on the x-axis.

```

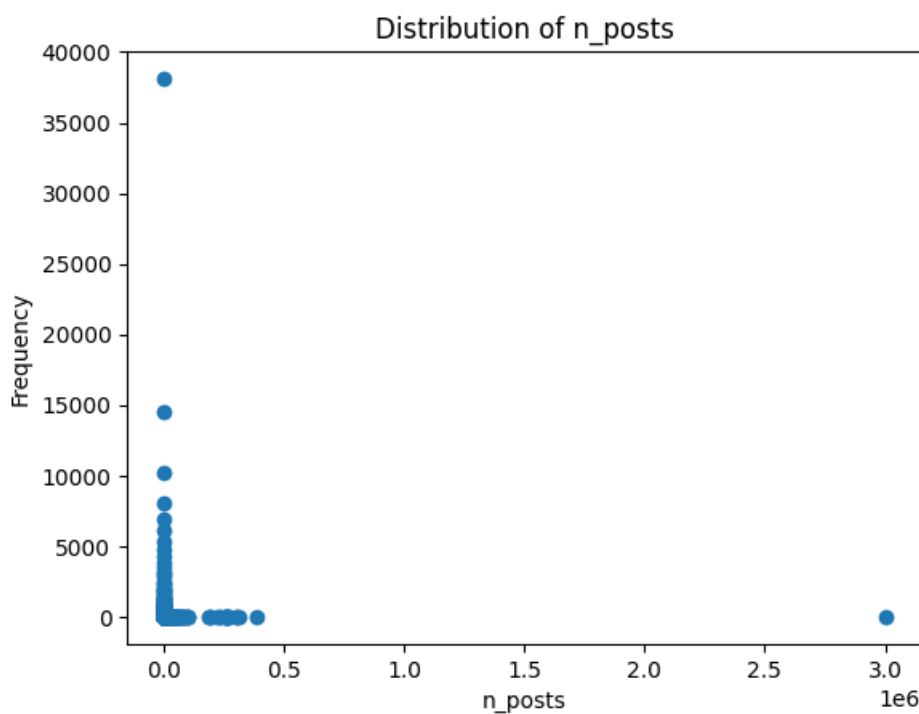
In [6]: # print n_posts distribution
max_posts = np.ceil(df['n_posts'].max() / 10) * 10 # rounded to nearest multiple of 10
bins = np.arange(0, max_posts + 10, 10) # steps of 10
binned_table = df.groupby(pd.cut(df['n_posts'], bins, include_lowest=True), observed=False).size().reset_index(name='counts')

binned_table = binned_table[binned_table['counts'] != 0]
binned_table['bin_middle'] = binned_table['n_posts'].apply(lambda x: x.mid)

# Scatterplot of table
plot = plt.scatter(binned_table['bin_middle'], binned_table['counts'])
plt.title('Distribution of n_posts')
plt.xlabel('n_posts')
plt.ylabel('Frequency')

```

Out[6]: Text(0, 0.5, 'Frequency')



This is a very zoomed out look. Before we focus this scatterplot, we'll check out the outlier that generated 3 million posts:

```
In [7]: print('Articles with the most posts')
print(df.sort_values(by='n_posts', ascending=False)[['title', 'n_posts']].head(10))
```

```
Articles with the most posts
```

| | title | n_posts |
|--------|--|---------|
| 108017 | Off-Topic-Ticker mit TickerOfLove | 3000166 |
| 108018 | Off-Topic-Ticker mit Ticker of Love and Laughter | 389259 |
| 102886 | Seuchenticker-Basislager | 314567 |
| 109509 | Seuchenticker-Basislager 3 | 308522 |
| 170225 | Ticker-Basislager 13 | 300140 |
| 149054 | (Seuchen)ticker-Basislager 10 | 275232 |
| 135239 | Seuchenticker-Basislager 7 | 270924 |
| 115763 | Seuchenticker-Basislager 4 | 268672 |
| 157329 | Ticker-Basislager 11 | 263017 |
| 149052 | Seuchenticker-Basislager 9 | 262740 |

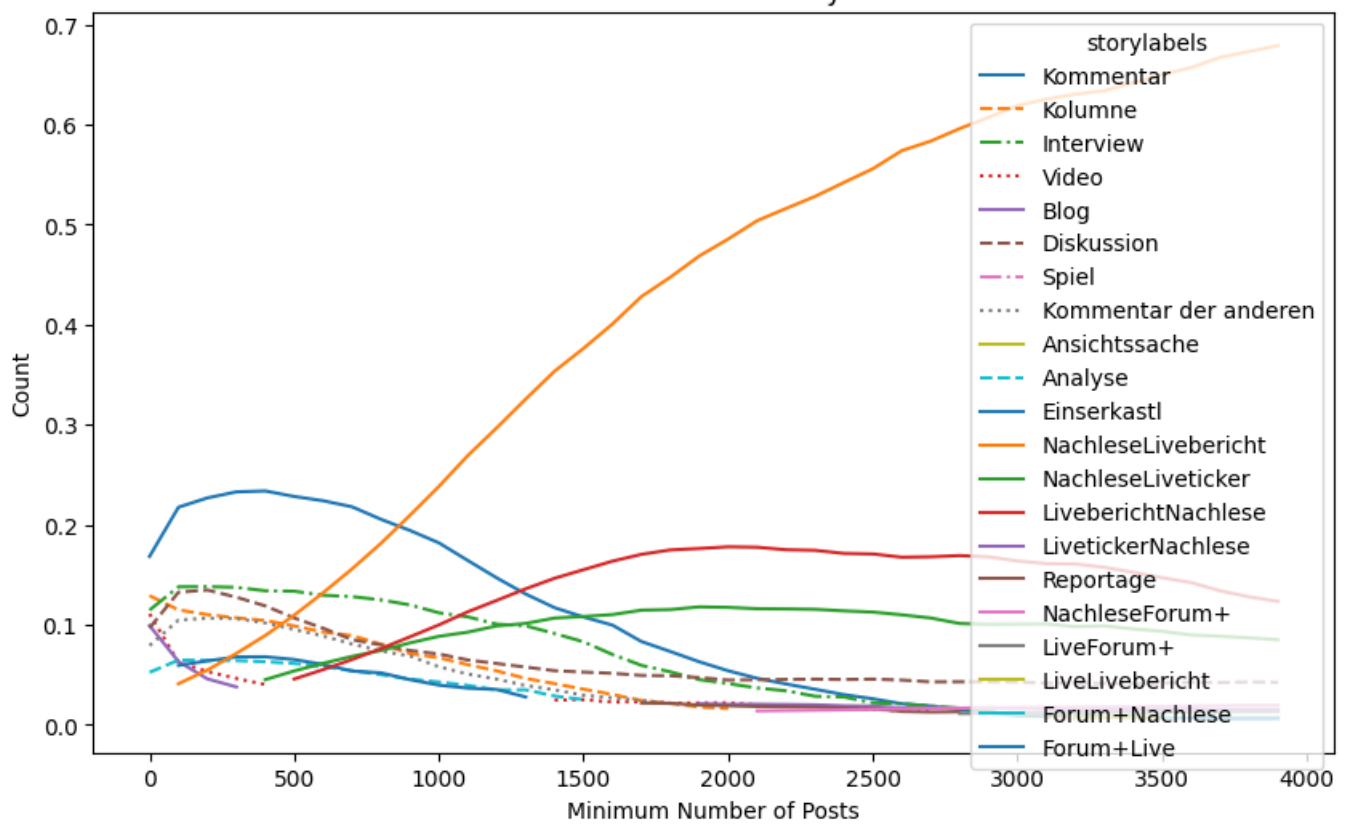
It is the love forum that generated 3 million posts, after that, the next post also deals with topics of love but has 400 thousand posts. The rest of the top 10 seem to also all be forum posts. Lets visualize the relative share of the different storylabels when the post count increases (keeping in mind that only a minority of posts has a storylabel):

```
In [8]: r = range(0, 4_000, 100)
label_tracker = []
for min_posts in r:
    d = df[df['n_posts'] > min_posts].storylabels.value_counts()[:10]
    label_tracker.append(d / d.sum())

df_label_tracker = pd.concat(label_tracker, axis=1).T
df_label_tracker.reset_index(drop=True, inplace=True)
df_label_tracker['min_posts'] = r
df_label_tracker.set_index('min_posts', inplace=True)
styles = ['-', '--', '-.', ':', '-', '--', '-.', ':', '-', '--']
df_label_tracker.plot(kind='line', style=styles, figsize=(10, 6))

plt.title('Number of Posts vs Storylabels')
plt.xlabel('Minimum Number of Posts')
plt.ylabel('Count')
plt.show()
```

Number of Posts vs Storylabels

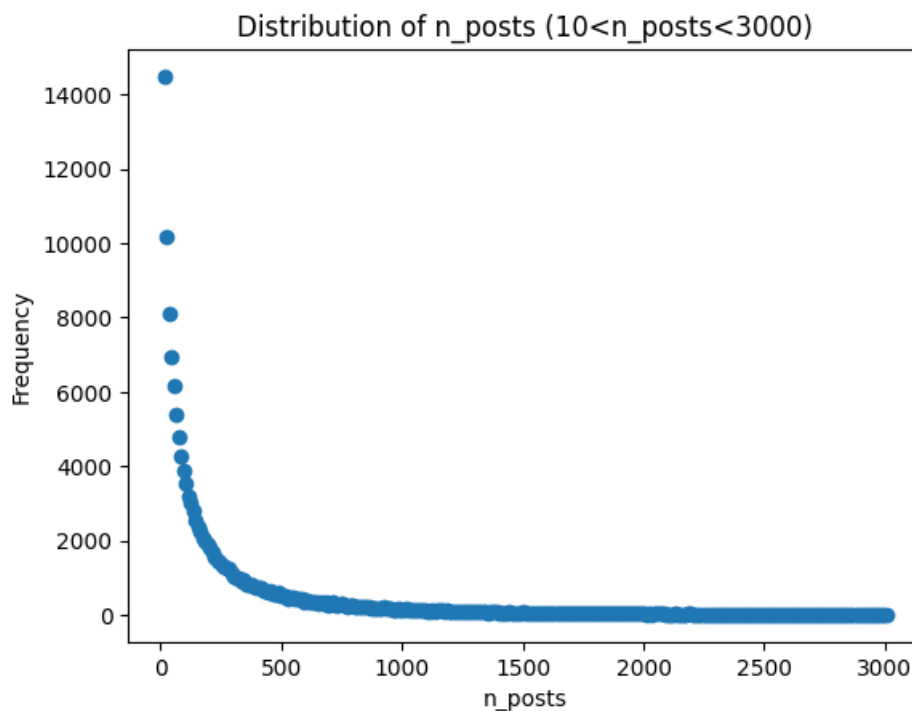


We can see that the relative share of articles with tagged with forum associated story labels increases above 1.500 posts. If we also plot the na, values, we can see them actually becoming the minority above ~2.700 posts, but we left out those Na in order to better illustrate the relative share of other storylabels decreasing relative to Liveposts.

But lets return to the distribution of posts and focus the scatterplot a bit. This time we will omit posts within the 0-10 range as that bucket dwarfs the scale of the following buckets (we're avoiding a log scale for now). We'll also only plot the buckets up to 3.000 posts:

```
In [9]: plt.scatter(binned_table['bin_middle'][1:301], binned_table['counts'][1:301])
plt.title('Distribution of n_posts (10<n_posts<3000)')
plt.xlabel('n_posts')
plt.ylabel('Frequency')
```

```
Out[9]: Text(0, 0.5, 'Frequency')
```



The distribution of posts seems to exhibit a power law distribution. The curve very beautifully follows some logarithmic function.

As seen in the first distribution plot, the distribution has an extremely long tail that gradually thins out. The 3million ticker is a very strong outlier here, but even up to 400 thousand posts there are still a few articles in our large data set. Before we decide on a strategy on how to deal with those, lets check out the boxplot:

```
In [10]: plt.boxplot(df['n_posts'], vert=False, )
print('Five-point summary:\n', df['n_posts'].describe(), '\n')
print('Increasingly smaller buckets:\n', df['n_posts'].quantile([.25, .5, .75, .9, .95, .99, .999, .9999]), '\n')
```

Five-point summary:

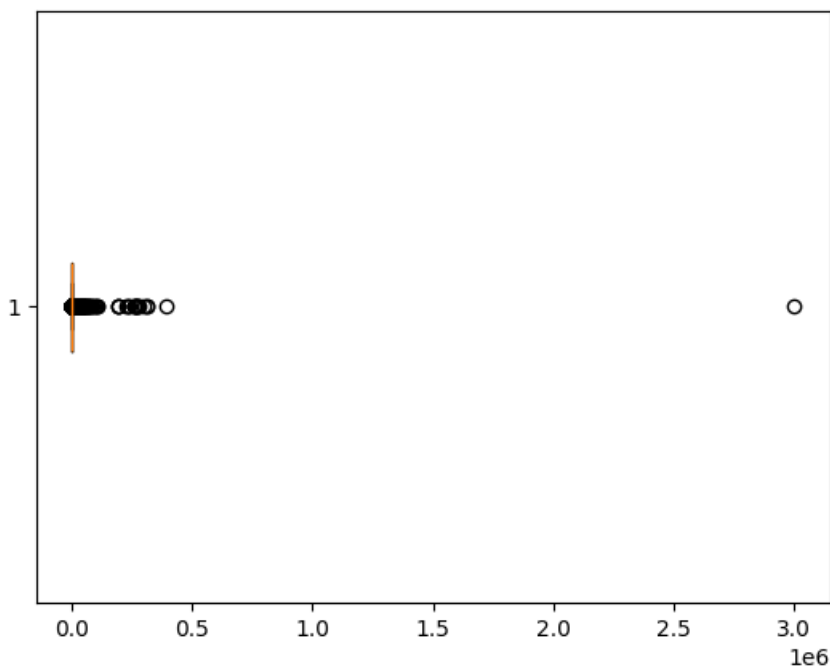
```
count    1.812000e+05
mean      3.359409e+02
std        7.640062e+03
min         0.000000e+00
25%        1.500000e+01
50%        7.300000e+01
75%        2.510000e+02
max        3.000166e+06
```

Name: n_posts, dtype: float64

Increasingly smaller buckets:

```
0.2500    15.000000
0.5000    73.000000
0.7500   251.000000
0.9000   637.000000
0.9500  1047.000000
0.9900  2751.020000
0.9990  20050.811000
0.9999  181071.088898
```

Name: n_posts, dtype: float64



This again confirms that there is a very long tail. For the purposes of our machine learning project, we have to consider how we'll deal with this long stretch. Initially the idea was to produce a regression neural network, with a single output neuron that would output continuous values. In this case, the outliers would have caused much trouble with scaling the data and we would have probably trimmed our dataset above a certain percentile-based threshold e.g. 99.99%.

We decided to simplify our regression problem by using (roughly) equally sized bins which are still an ordinal variable so we can still perform a regression without losing the valuable information for the gradient of how wrong a given prediction is. I.e. instead of just using categories, we can still perform regression over those bins and tell the network in our loss function how far a predicted bin is from the actual bin. Those bins should cover meaningful ranges in the number of posts, reflecting the logarithmic distribution of posts.

For this we use the `qcut` pandas function. This function however throws an error if the bins are not equally sized. Because there are so many articles with 0 posts (5% of our data), it will be impossible to make a bin that contains 1 % of the data withing those bounds. Similarly the articles with 1 post are also a bin that would be larger than the other bins. Thus we:

1. create a special bin for all articles with 0 posts
2. we set the `duplicates=drop` kwarg which allows for uneven sized bins (it will drop overlapping buckets)

Considering the power-law distribution of our data, it is inevitable that the first buckets consist of very small ranges. Even with manually assigning 0-posts to their own bucket, the first set of buckets are actually also just single-value (e.g. '2', '3') ranges with a disproportionally large number of observations in them. For training purposes it of course would be good if all our bins are equally sized - but maybe it is also important to make the network good at spotting unsuccessful posts.

First we tried 100 buckets but saw that the `qcut` function dropped a lot of buckets. Thus we decided for a smaller number of buckets that could still meaningfully dissect the data and just went with 64 buckets because it's a nice number. Furthermore you will notice that `q` is set

to 65, because there are 2 overlapping buckets which will be dropped.

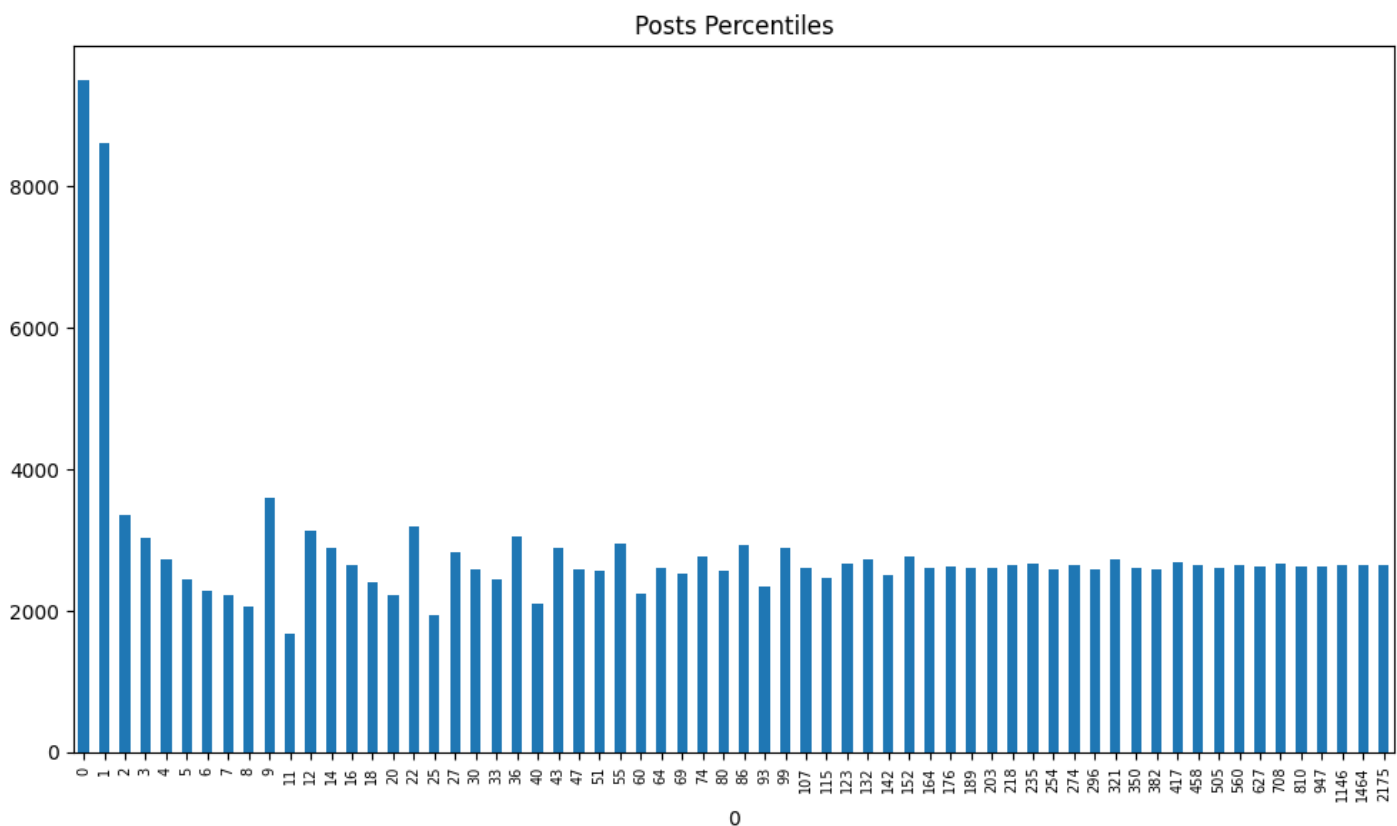
```
In [11]: nonzero = df['n_posts'] > 0

# Apply qcut to these rows and convert the result to string
df.loc[nonzero, 'n_posts_percentile'] = pd.qcut(df.loc[nonzero, 'n_posts'], q=65,
                                                duplicates='drop').astype(str)

# Set the 0 posts to '0'
df.loc[df['n_posts'] == 0, 'n_posts_percentile'] = '(0.0, 1.0]'

# extract the lower bound of the bins and convert to float
float_pattern = r'\d+\.\d+'
bounds = df['n_posts_percentile'].value_counts().index.str.extract(
    f'({float_pattern}), ({float_pattern})').astype(float)
bounds = bounds.round().astype(int)

s = pd.Series(df['n_posts_percentile'].value_counts().values, index=bounds[0])
s.sort_index().plot(kind='bar', figsize=(10, 6))
plt.xticks(fontsize=7)
plt.title('Posts Percentiles')
plt.tight_layout()
plt.show()
print(f'There are {len(df['n_posts_percentile'].value_counts())} buckets')
```



There are 64 buckets

The first set of buckets turn out to just be the range from 0 to 12, and we can see that the 0 and 1 buckets are both very large bins. We can also see that our percentiles cut off at around 2.100 posts, after which the long tail begins. Compared to the initial distribution plots though, we have mostly flattened the curve.

To produce the final target variable, we map those buckets to the range from 0 to 63. We won't transform this data any further as we hope that a sufficiently large network will be able to cope with that range without centering the values as well. As the buckets are roughly evenly distributed now (save for 0 and 1), we will not normalize our data.

```
In [12]: target_map = pd.DataFrame({
    'lower_b': bounds[0],
    'upper_b': bounds[1],
    'bounds': df['n_posts_percentile'].value_counts().index
})

target_map.sort_values(by='lower_b', inplace=True)
target_map.reset_index(drop=True, inplace=True)

# Map 'n_posts_percentile' to 'bounds' and assign the index of the matching row to 'target'
df['target'] = df['n_posts_percentile'].map(target_map.reset_index().set_index('bounds')['index'])
# save the target map
target_map.to_csv('./data/target_map.csv', index=False)
target_map
```

Out[12]:

| | lower_b | upper_b | bounds |
|-----|---------|---------|---------------------|
| 0 | 0 | 1 | (0.0, 1.0] |
| 1 | 1 | 2 | (0.999, 2.0] |
| 2 | 2 | 3 | (2.0, 3.0] |
| 3 | 3 | 4 | (3.0, 4.0] |
| 4 | 4 | 5 | (4.0, 5.0] |
| ... | ... | ... | ... |
| 59 | 810 | 947 | (810.0, 947.0] |
| 60 | 947 | 1146 | (947.0, 1146.0] |
| 61 | 1146 | 1464 | (1146.0, 1464.0] |
| 62 | 1464 | 2175 | (1464.0, 2175.0] |
| 63 | 2175 | 3000166 | (2175.0, 3000166.0] |

64 rows × 3 columns

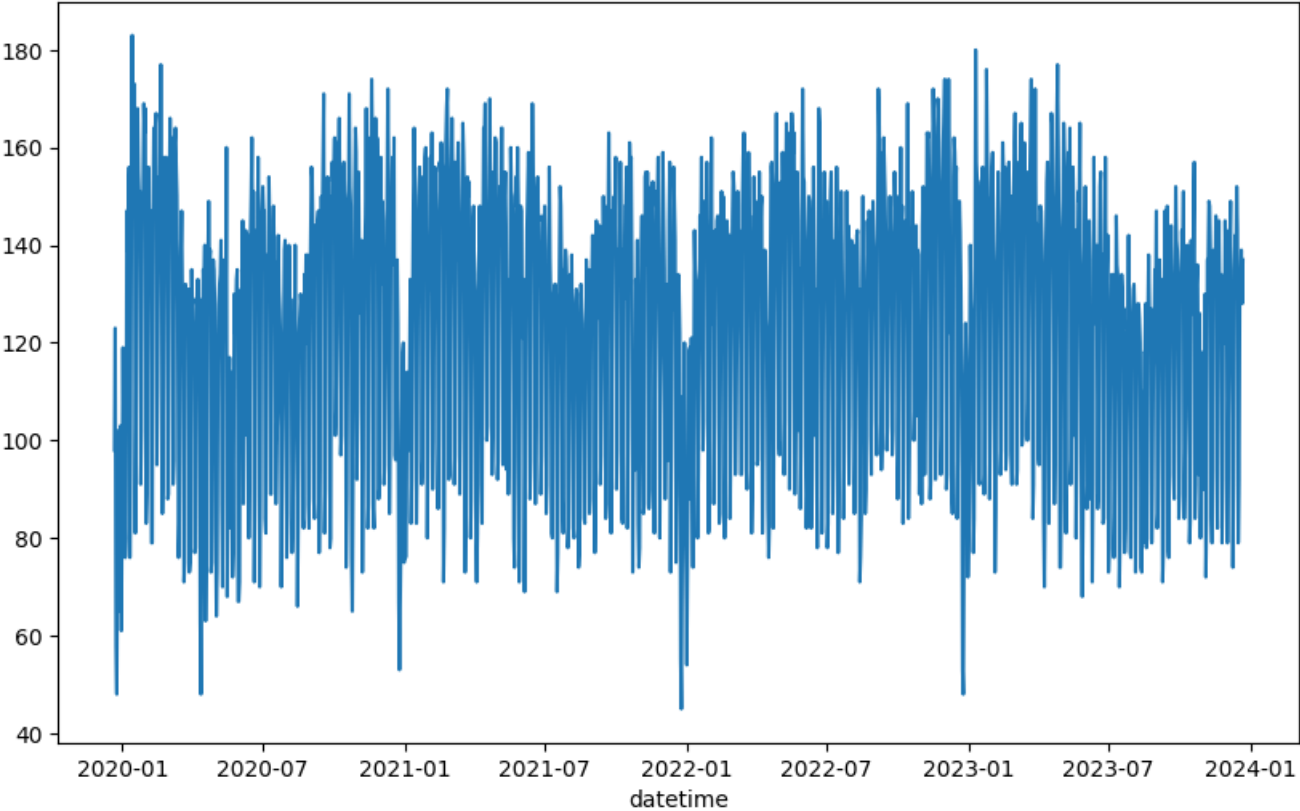
Time series analysis

First, lets' plot the number of articles on the frontpage per day

In [13]:

```
df.groupby(df['datetime'].dt.date).count()['title'].plot(figsize=(10, 6))
```

Out[13]: <Axes: xlabel='datetime'>

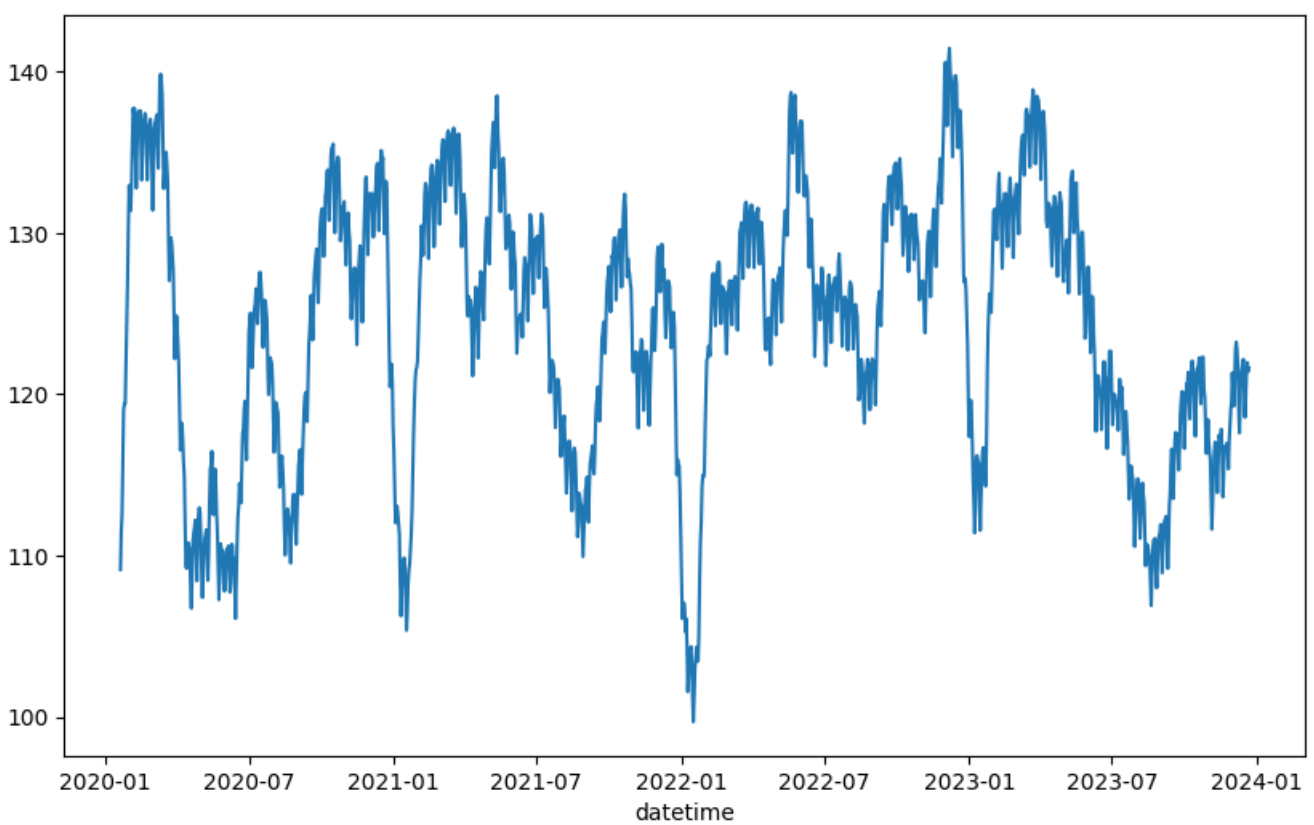


Too many fluctuations, let's look at the rolling 30 day average

In [14]:

```
df.groupby(df['datetime'].dt.date).count()['title'].rolling(30).mean().plot(figsize=(10, 6))
```

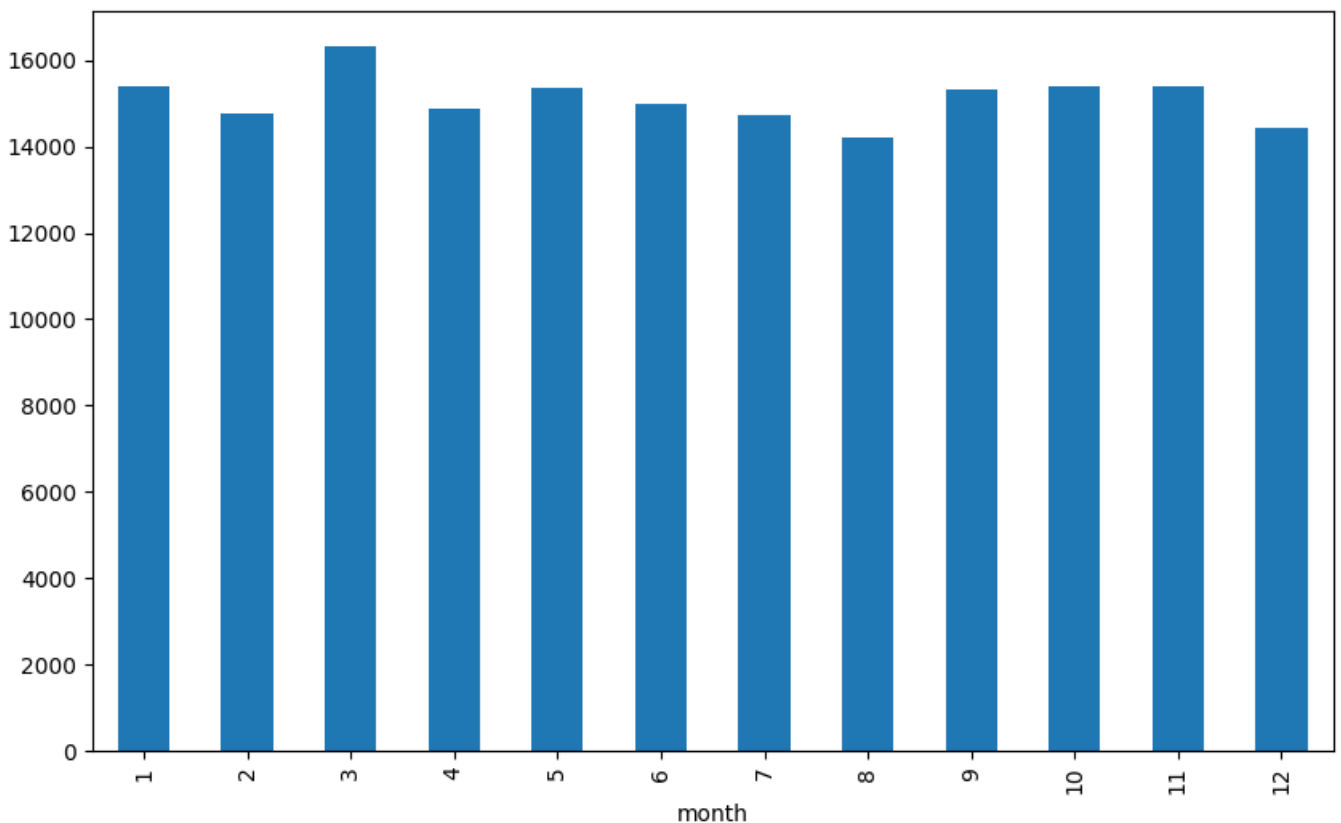
Out[14]: <Axes: xlabel='datetime'>



The number of articles per day does seem to follow some seasonal trends. Maybe at the end of the year the number of articles drops, as well as during the summer months due to the Sommerloch. Let's investigate this by aggregating the number of articles by month.

```
In [15]: df['month'] = df['datetime'].dt.month
df.groupby('month').count()['title'].plot(kind='bar', figsize=(10, 6))
```

Out[15]: <Axes: xlabel='month'>



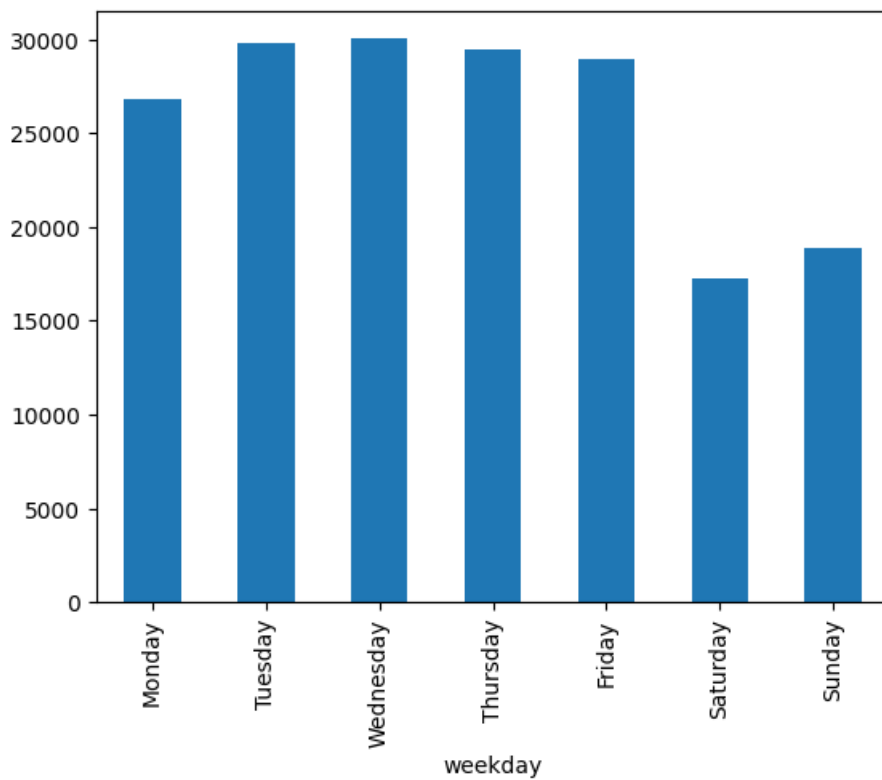
Well, that rejects that hypotheses, the articles seem to even out over 4 years by month.

Next, let's add a weekday column to our data.

```
In [16]: df['weekday'] = df['datetime'].dt.weekday
df.groupby('weekday').count()['title'].plot.bar()
plt.xticks(np.arange(7), ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday'])
plt
```



```
Out[16]: <module 'matplotlib.pyplot' from 'c:\\Users\\Paul\\AppData\\Local\\Programs\\Python\\Python312\\Lib\\site-packages\\matplotlib\\pyplot.py'>
```

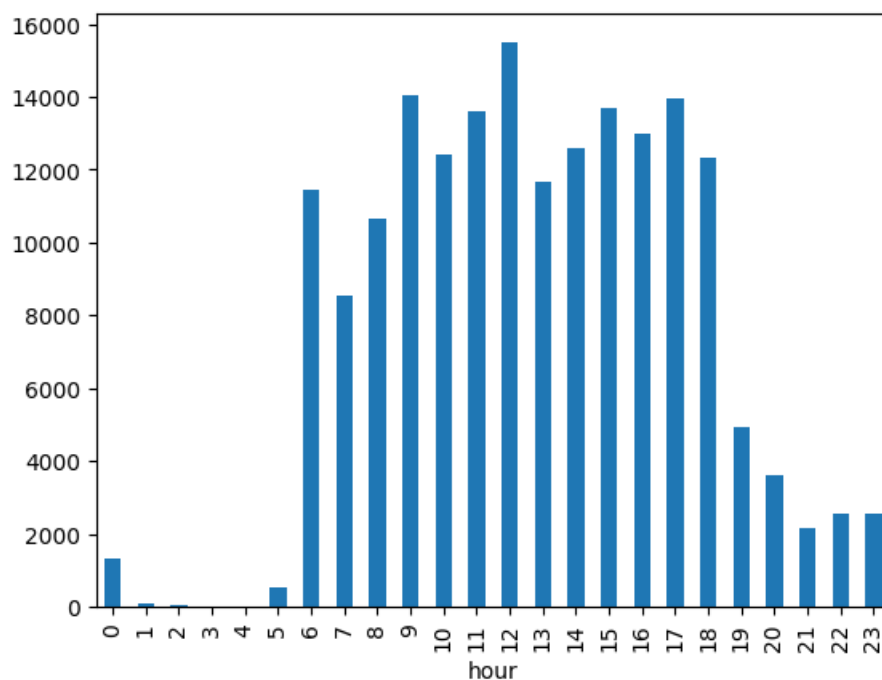


less datapoints for the weekends it seems. We will make use of this weekday variable in our predictions as well.

We'll also add a time of day column and plot this.

```
In [17]: # add a hour column
df['hour'] = df['datetime'].dt.hour
# which hour has the most articles?
df.groupby('hour').count()['title'].plot.bar()
```

```
Out[17]: <Axes: xlabel='hour'>
```



Most articles are posted during the day from 6am to 6pm. Continuously adding new articles during the day probably helps engagement as well.

Label analysis

We already mentioned in the scraping functions that most articles seem to lack storylabels, but let's investigate this a bit more.

```
In [18]: label_counts = df['storylabels'].value_counts(dropna=False)
print(label_counts)
print(f'\nPercentage of articles with a story label: {100-(label_counts.iloc[0]/len(df)) * 100:.2f}%')
```

```
storylabels
NaN      143784
Kommentar  4694
Kolumne   3608
Video     3238
Interview  3199
...
Liveticker      1
Userkommentar   1
FeatureReportage 1
Photoblog       1
Forum+Vorschau  1
Name: count, Length: 77, dtype: int64
```

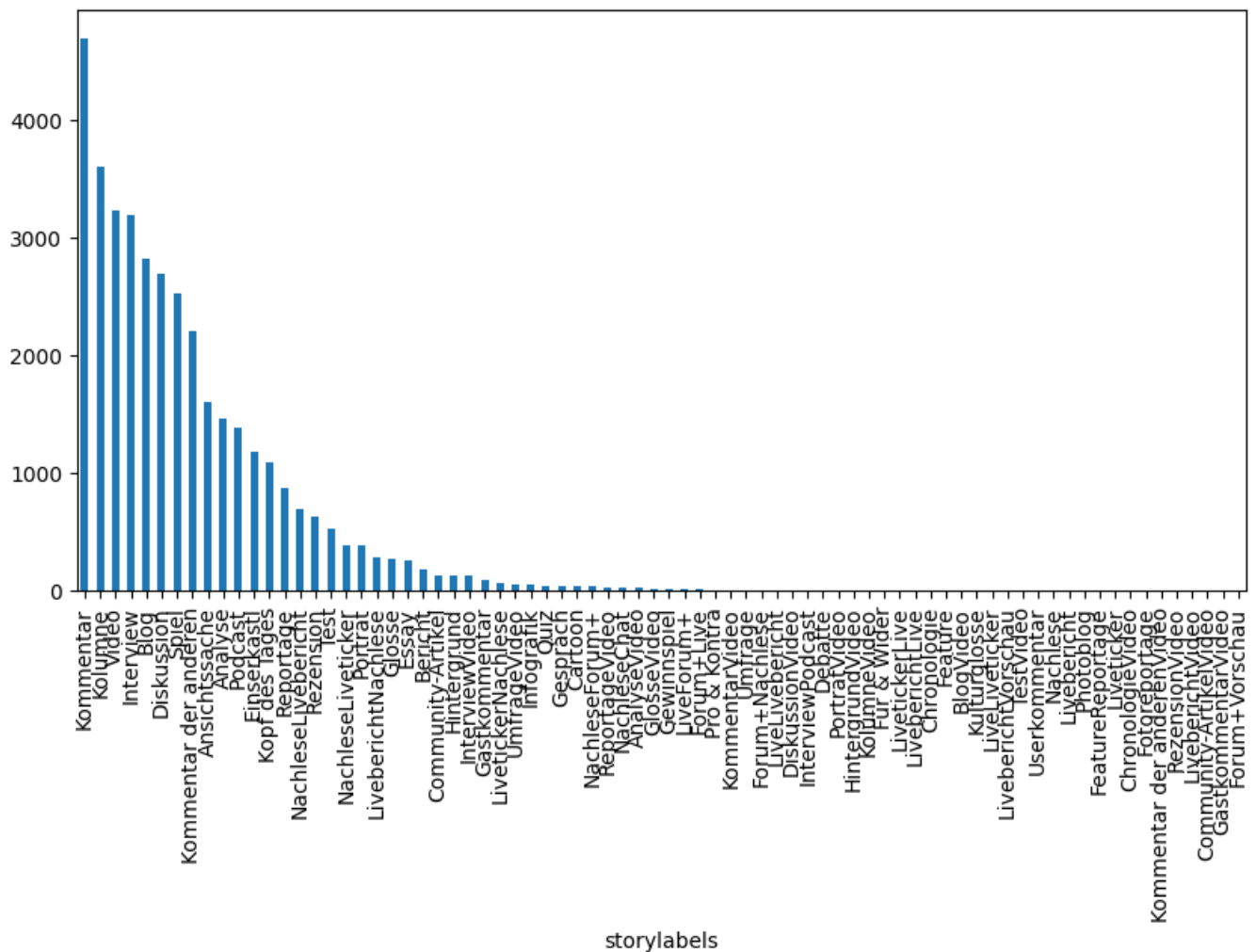
Percentage of articles with a story label: 20.65%

only around 20% of articles are labeled. We will not use those labels for our prediction but they were useful for filtering out Ticker posts before. As mentioned, it would also be an interesting labelling task to find labels for the remaining 80% but for now we will focus on the `n_posts` variable.

Next we will not look at the number of posts a given label is going to generate, just at which are the most common storylabels.

```
In [19]: # Plot the number of articles per story label
plt.figure(figsize=(10, 5))
df['storylabels'].value_counts().plot.bar()
```

```
Out[19]: <Axes: xlabel='storylabels'>
```



This time the Ticker labels are further behind the distribution, the list being headed by periodicals like `Kommentar`, `Kolumne` etc.

Lastly, lets look at the kicker labels, for which there are no na values.

```
In [20]: kickers = df['kicker'].value_counts(dropna=False)
print(f'Number of unique kicker labels: {len(kickers)}\n')
print(f'Most common kicker labels:\n{kickers[:20]}\n')

print(f'There are only {df['kicker'].isna().sum()} na kickers, lets drop them')
# drop kicker na values
df.dropna(subset=['kicker'], inplace=True)
```

Number of unique kicker labels: 44720

Most common kicker labels:

```
kicker
Fußball          3133
Nachrichtenüberblick  3099
Netzpolitik      2674
Sudoku           2414
Bundesliga       1711
Sport            1651
USA              1515
IT-Business      1464
Coronavirus      1375
Games            1356
Tennis           1244
Switchlist       1208
Krieg in der Ukraine 1203
Deutsche Bundesliga 1180
Etat-Überblick   1161
Hans Rauscher    1153
Wintersport      1123
TV-Tagebuch     1080
Eishockey        1058
Thema des Tages  1032
Name: count, dtype: int64
```

There are only 84 na kickers, lets drop them

Those also are dominated by periodicals like Fußball, Nachrichtenüberlick (daily summary of posted articles) etc.

The title along with the subtitle will be our most important predictors. In the following codecell, we concatenate those into a long string along with the kicker-label. We will later use embeddings to encode those strings as numeric vectors.

First, let's plot the distribution of string lengths.

```
In [21]: df['text'] = df['title'] + ' ' + df['kicker'].fillna('') + ' ' + df['subtitle'].fillna('')

print('Example of concatted string:\n', df['text'][0], '\n')
df['text'].str.len().hist(bins=100)
short_txt = df['text'].apply(len).nsmallest(3).index.map(df['text']).tolist()
print('\nShortest titles:\n', '\n'.join(list(short_txt)))

print(f'\nMean title length: {df["text"].str.len().mean():.2f} characters')
print(f'Standard deviation: {df["text"].str.len().std():.2f} characters')
print(f'Median title length: {df["text"].str.len().median():.2f} characters')
```

Example of concatted string:

Real Madrid stolpert mit Aluminiumpech im Titelrennen Primera Division Die Königlichen können Bilbao daheim nicht besiege
n

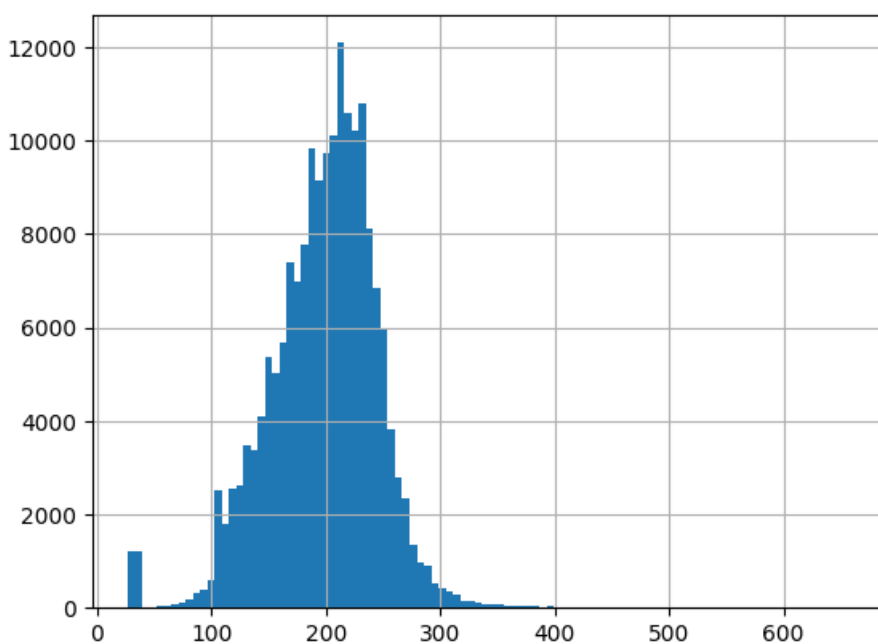
Shortest titles:

Sudoku mittel 4493a Sudoku
Sudoku mittel 4497a Sudoku
Sudoku mittel 4503a Sudoku

Mean title length: 197.88 characters

Standard deviation: 47.80 characters

Median title length: 203.00 characters



Again a distribution with a long tail that thins out very gradually. During our nlp processing, we will generate vectors of uniform length - thus there is no need to trim our texts.

The first large bin are the sudokus - short texts of roughly the same length. We will also keep them as they are also a part of the frontpage

Vectorizing our features

As inputs for our network, we need numeric values. The datetime information is already basically numeric, but the text will need some processing.

Text Processing

For the following steps, we will be using the spacy library which will help us with all the lemmatization and tokenization we need to turn our strings into tokens. We will also be using the german 'news' (large) pipeline, available [here](#) - Warning, it is 500mb of size. Its vocabulary consists of many newspaper headlines, so it should be perfect for our application. It also has embedding vectors for a large vocabulary (more on that in the next cell).

In case you want to run this locally, it can be installed using `python -m spacy download de_core_news_md`.

```
In [22]: import spacy

# nlp = spacy.load("de_core_news_md")
nlp = spacy.load("de_core_news_lg")

print("Pipeline Components:\n", ", ".join([n for n, p in nlp.pipeline]))
```

Pipeline Components:
tok2vec, tagger, morphologizer, parser, lemmatizer, attribute_ruler, ner

Generating document vectors

This step will add a new column `spacy_vector`, which contains the mean of all token vectors of a document (headline+subtitle+tags) generated by the embeddings of the spacy model. Spacy uses *Word2VecEmbeddings*, which encodes semantic meaning based on the word's contextual use in the `de_core_news` dataset.

These document vectors provide a numerical representation of our headlines, while still retaining information and at the same time producing vectors of uniform size.

```
In [23]: from tqdm import tqdm
tqdm.pandas() # for progress_apply

def process_text(text):
    doc = nlp(text)
    return doc.vector

df['doc_vector'] = df['text'].progress_apply(process_text)
```

0% | 0/181116 [00:00<?, ?it/s]

In case we want to work with the information of our model again, we save our nlp object as well.

```
In [24]: print('The vectors are of length ', df['doc_vector'][0].shape)
print('The vectors are of type', df['doc_vector'][0].dtype)
print('\n', df['doc_vector'].head())
```

The vectors are of length (300,)
The vectors are of type float32

```
0    [-0.37552637, 0.20481217, -0.216765, 0.0486718...
1    [-0.080757536, 0.54684085, -0.5815048, 0.41144...
2    [0.69385666, -0.0836207, -0.5870395, 0.8096802...
3    [-0.14959173, -0.72091854, -1.3704777, -1.0522...
4    [-0.021010712, -0.18174057, -0.4104177, -1.398...
```

Name: doc_vector, dtype: object

We will finally create our inputs column which also contains the month, weekday and hour. We think this information might be useful, e.g. if certain stories are posted in the summer months/ winter, if it is a particular day of the week, if it was posted in the morning or evening. We omit the year to limit overfitting, as well as it's limited usefulness when we consider that those headlines are newsstories and we now have 2024.

We will also normalize those date vectors which means their elements should be roughly the same size as the elements of the document vectors.

```
In [25]: import numpy as np

def process_row(row):
    date_vector = np.array([
        row['datetime'].month,
        row['datetime'].weekday(),
        row['datetime'].hour
    ])
    # Normalize the date_vector, ensure it is float32 like the doc vectors
    norm_date_vector = (date_vector / np.linalg.norm(date_vector)).astype(np.float32)
    return np.concatenate((norm_date_vector, row['doc_vector']))

df['inputs'] = df.apply(process_row, axis=1)

print(f'Sequence length: {len(df['inputs'])[0])}')
```

Sequence length: 303

Saving to parquet

The data processing is thus completed. We will save our data using the parquet file format which not only handles dtypes like lists natively, it is also a binary format that is much more efficient at storing all our data. You will need to install `pyarrow` if you want to run this yourself.

```
In [26]: df[['inputs', 'target', 'n_posts']].to_parquet('data/ml_data.parquet', index=False)
```

Machine Learning

Loading Data

We open our ML-ready dataset

- **inputs:** First we concatenated the title, subtitle and article labels. Then we tokenized & lemmatized individual words, and used Spacys large german model to look up embedding vectors for every token. Those vectors were originally generated using the word2vec algorithm, where the goal is to place words in a 300-dimensional space where words that are used in similar contexts have a smaller euclidean distance in the vector space. Then we average all word vectors for the tokens in our concatenated text, which for our short texts should roughly correspond to the overall position of the text in the 300-dimensional space expressed by the vector. To this we append the year, month and weekday as components of a normalized vector.
- **targets:** The targets correspond to 64 bins, each bin representing a range in the number of posts an article has generated. This was done in order to:

1. Have our target variable be roughly evenly distributed, as the bins grow in size to mitigate the power-law distribution present in the actual values.
2. Limit the range of values the neural network needs to predict.

In effect, we have turned a regression problem into a classification, and then back into a regression again. This way we can still preserve the meaning if a prediction closely misses the correct bin, which lets the neural network adjust more precisely than in a classification where this distance information would be lost.

```
In [1]: import pandas as pd

df = pd.read_parquet('data/ml_data.parquet') # requires pyarrow
df.head()
```

```
Out[1]:
```

| | inputs | target | n_posts |
|---|---|--------|---------|
| 0 | [0.45066947, 0.22533473, 0.8637831, -0.3755263... | 18 | 30 |
| 1 | [0.4656903, 0.23284516, 0.8537656, -0.08075753... | 12 | 16 |
| 2 | [0.48154342, 0.24077171, 0.84270096, 0.6938566... | 34 | 104 |
| 3 | [0.48154342, 0.24077171, 0.84270096, -0.149591... | 20 | 35 |
| 4 | [0.49827287, 0.24913643, 0.8304548, -0.0210107... | 3 | 4 |

Loading df into torch, train test split

Now we start using pytorch. I was able to finally play with cuda on my trusty old gtx970, but this code should be agnostic to the type of device.

In the following cell, we create tensors for our inputs and target values, and combine them into a tensor dataset, on which we perform an 80/20 test split. We chose batch size 64, meaning 64 datapoints are pushed to the gpu in one training run. This value utilizes my gpu slightly more (cuda usage according to task manager at 40%), while at the same time hopefully not being too large to cause overfitting. The train and test loader objects will be responsible for pushing the data to the gpu.

```
In [2]: import torch
from torch.utils.data import DataLoader, TensorDataset, random_split

# Using cuda
my_device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print('We are using ', my_device)

# Create a TensorDataset
inputs = torch.tensor(df['inputs'].tolist(), device=my_device)
targets = torch.tensor(df['target'].values, device=my_device)
```

```
dataset = TensorDataset(inputs, targets)

# Split into training and test set
train_size = int(0.8 * len(dataset))
test_size = len(dataset) - train_size
train_dataset, test_dataset = random_split(dataset, [train_size, test_size])

# Create DataLoaders
batch_size = 64
train_loader = DataLoader(train_dataset, batch_size=batch_size, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=batch_size, shuffle=False)
```

We are using cuda

```
C:\Users\Paul\AppData\Local\Temp\ipykernel_12152\2500532675.py:9: UserWarning: Creating a tensor from a list of numpy.ndarrays is extremely slow. Please consider converting the list to a single numpy.ndarray with numpy.array() before converting to a tensor. (Triggered internally at ..\torch\csrc\utils\tensor_new.cpp:278.)
    inputs = torch.tensor(df['inputs'].tolist(), device=my_device)
```

Model definitions

Deep regression network

Here we define our regression neural network. It takes as an argument the layer sizes and generates `len(layers_sizes)` layers. Furthermore we also add a dropout probability, meaning that during training, with a set probability, a neurons output will randomly be ignored, meaning its output will not be fed as input to the next layer. This pruning technique should hopefully make our model more robust to overfitting as the model should not become overly reliant on single neurons or pathways through the network, encouraging a more distributed and robust internal representation.

As activation functions, we use the nonlinear ReLU function, which worked decently in our testing.

The output layer consists of a single neuron, which should predict values in the range 0-63. We do not perform activation on this final output as to not squash its output.

```
In [3]: import torch.nn as nn
import torch.nn.functional as F

class RegressionNN(nn.Module):
    def __init__(self, drop, layer_sizes):
        super(RegressionNN, self).__init__()

        # Create layers dynamically
        self.layers = nn.ModuleList()
        for i in range(len(layer_sizes)-1):
            self.layers.append(nn.Linear(layer_sizes[i], layer_sizes[i+1]))

        self.dropout = nn.Dropout(drop)

        self.output_layer = nn.Linear(layer_sizes[-1], 1)

    def forward(self, x):
        for layer in self.layers:
            x = F.relu(layer(x))
            x = self.dropout(x)

        x = self.output_layer(x)
        return x
```

Training loop

Here we define the training loop, which takes as parameters the number of epochs and learning rate. An evaluation criterion can optionally be specified but we decided to go with the standard mean-squared-error approach as it is a staple in regression problems and maps well to our intention of minimizing the error of our predictions.

```
In [4]: import torch
import torch.optim as optim
```

```

from tqdm import tqdm

def train(model, epochs, learning_rate, criterion=nn.MSELoss()):
    optimizer = optim.Adam(model.parameters(), lr=learning_rate)
    losses = []
    for epoch in tqdm(range(epochs)):
        model.train()
        total_loss = 0

        for inputs, targets in train_loader:
            inputs, targets = inputs.to(my_device), targets.to(my_device)
            optimizer.zero_grad()
            outputs = model(inputs.float())
            loss = criterion(outputs.squeeze(), targets.float())
            loss.backward()
            optimizer.step()
            total_loss += loss.item()

        epoch_loss = total_loss/len(train_loader)
        losses.append(epoch_loss)
    return losses

```

Evaluation

On test set

This function will print the average mean squared error among all batches of the training set (20% of our data).

```

In [5]: def eval_on_test(model, criterion=nn.MSELoss()):
        model.eval()
        test_loss = 0
        with torch.no_grad():
            for inputs, targets in test_loader:
                inputs, targets = inputs.to(my_device), targets.to(my_device)
                outputs = model(inputs.float())
                loss = criterion(outputs.squeeze(), targets.float())
                test_loss += loss.item()
        print(f"Test Loss: {test_loss/len(test_loader)}")

```

Plotting the loss curve

This function simply plots the loss as a function of the number of epochs during training, allowing us to monitor our training efficiency and check once the model stops improving.

```

In [6]: import matplotlib.pyplot as plt

def plot_losses(losses):
    # Plot the losses
    plt.plot(losses)
    plt.xlabel('Epoch')
    plt.ylabel('Loss')
    plt.title('Loss Curve')
    plt.show()

```

Testing on the current frontpage

Finally, we thought it would be interesting to fetch a given day's actual frontpage and see how our model would fare.

For this, 4 functions are defined:

- `get_final_articles_df()`: Unfortunately I found no good way to import functions from other notebooks as you would normally do with python modules. Thus this method violates DRY and is a copy-paste of everything needed from the scraping and processing notebook to return a dataframe of a given day's frontpage, along with processing the text into a vector.
- `plot_preds()`: This plots our predicted bounds as a bar, along with the actual values as a red line, sorted in ascending order of actual number of posts.

- `print_preds()`: This prints all the articles that were within our predicted bounds, as well as the two articles furthest away from our prediction.
- `compare_prediction_for_date()`: This is just a wrapper for the above three functions, which takes a model as well as a date as parameters.

```
In [7]: from datetime import datetime, timedelta
import matplotlib.pyplot as plt
from bs4 import BeautifulSoup
from functools import cache
import numpy as np
import requests
import spacy

@cache # We cache the results of this function to avoid fetching the same data multiple times
def get_final_articles_df(date):
    link = f'https://www.derstandard.at/frontpage/{date.strftime('%Y/%m/%d')}'
    # fetch the html content of a derstandard.at page
    response = requests.get(link, cookies={'DSGVO_ZUSAGE_V1': 'true'})
    soup = BeautifulSoup(response.content, 'html.parser')
    # get the articles
    articles_dict = {}
    articles = soup.select('div.chronological>section article')
    for article in articles:
        title_tag = article.find('a')
        if title_tag and title_tag.has_attr('title'):
            title = title_tag['title']
            articles_dict[title] = article
    # make a list of the articles
    HOST = 'https://www.derstandard.at'
    article_data = []
    for title, article in articles_dict.items():
        data = {
            'title': title,
            'teaser-subtitle': None,
            'link': None,
            'time': None,
            'teaser-kicker': None,
            'n_posts': None,
            'storylabels': None
        }
        link = article.find('a')['href']
        if not link.startswith(HOST):
            link = HOST + link
        data['link'] = link
        time = [tag for tag in article.find_all('time') if 'datetime' in tag.attrs][0]
        data['time'] = time['datetime'].rstrip('\r\n')
        n_posts = article.find('div', 'teaser-postingcount')
        try: data['n_posts'] = int(n_posts.get_text(strip=True).rstrip('Posting').replace('.', ''))
        except: data['n_posts'] = 0
        for tag, class_name in [('p', 'teaser-kicker'),
                                ('p', 'teaser-subtitle'),
                                ('div', 'storylabels')]:
            found_tag = article.find(tag, class_=class_name)
            if found_tag:
                data[class_name] = found_tag.get_text(strip=True)
        article_data.append(data)
    # make a df
    df = pd.DataFrame(article_data)
    df.columns = df.columns.str.replace('teaser-', '')
    df.rename(columns={'time': 'datetime'}, inplace=True)
    df['datetime'] = pd.to_datetime(df['datetime'])
    df['text'] = df['title'] + ' ' + df['kicker'].fillna('') + ' ' + df['subtitle'].fillna('')
    # add embeddings
    nlp = spacy.load("de_core_news_lg")
    df['doc_vector'] = df['text'].apply(lambda t: nlp(t).vector)
    # add date
    def process_row(row):
        date_vector = np.array([
            row['datetime'].month,
            row['datetime'].weekday(),
            row['datetime'].hour
        ])
    # Normalize the date_vector, ensure it is float32 like the doc vectors
```

```

        norm_date_vector = (date_vector / np.linalg.norm(date_vector)).astype(np.float32)
        return np.concatenate((norm_date_vector, row['doc_vector']))
    df['inputs'] = df.apply(process_row, axis=1)
    return df

def plot_preds(preds):
    # Create a copy of preds and sort by actual value
    preds_copy = sorted(preds, key=lambda x: x[0])
    actual, lower_bound, upper_bound = zip(*preds_copy)
    # Create the plot
    plt.figure(figsize=(10, 6))
    plt.plot(actual, color='red', label='Actual')
    plt.fill_between(range(len(actual)), lower_bound, upper_bound, color='grey', alpha=0.5, label='Predi')
    plt.xlabel('Articles')
    plt.ylabel('number of posts')
    plt.title('Actual vs Predicted')
    plt.legend()
    plt.show()

def print_preds(preds, articles):
    within_bounds, outside_bounds = [], []
    for i, (act, lower, upper) in enumerate(preds):
        if lower <= act < upper: # upper bound is exclusive
            within_bounds.append(i)
        else:
            distance = min(abs(act - lower), abs(act - upper))
            outside_bounds.append((i, distance))
    # Sort the outside_bounds list by distance in descending order and get the first two indices
    furthest_indices = [i for i, _ in sorted(outside_bounds, key=lambda x: x[1], reverse=True)[:2]]
    # Fetch the rows from the articles DataFrame for within_bounds
    def print_in_color(indices, articles, preds, prediction_type):
        for i in indices:
            row = articles.iloc[i]
            print(f"{prediction_type} Prediction: {row['text']}")
            print(f"Number of Posts: {row['n_posts']}")
            print(f"Predicted Bounds: ({preds[i][1]}, {preds[i][2]})")
            print(f"Link: {row['link']}\n")
    # Print the correct predictions in green
    print("\033[92m")
    print_in_color(within_bounds, articles, preds, "Correct")
    # Print the incorrect predictions in red
    print("\033[91m")
    print_in_color(furthest_indices, articles, preds, "Incorrect")
    # Reset the color back to normal
    print("\033[0m")

@cache
def compare_prediction_for_date(model, date):
    # fetch articles for the date
    articles = get_final_articles_df(date)
    target_map = pd.read_csv('data/target_map.csv')
    predicted = [] # list of tuples (actual, pred_lower_bound, pred_upper_bound, row)
    for i, row in articles.iterrows():
        input_data = torch.tensor(row['inputs'], dtype=torch.float)
        input_data = input_data.to(next(model.parameters()).device)
        # Reshape the input data and pass it through the model
        input_data = input_data.unsqueeze(0) # Add a batch dimension
        with torch.no_grad():
            output = model(input_data)
        target_index = min(int(output.item()), 63) # ensure we don't go out of bounds
        tm = target_map.iloc[target_index]
        predicted.append((row['n_posts'], tm.lower_b, tm.upper_b))
    plot_preds(predicted)
    print_preds(predicted, articles)

```

Wrapping it all into a single function

And lastly, one function to wrap all the evaluations into one.

```
In [8]: def eval_trained_model(model, losses, date):
        plot_losses(losses)
        eval_on_test(model)
        compare_prediction_for_date(model, date)
```

Testing different models

We will be testing regression models of different sizes. For our evaluations of a particular frontpage, we picked a frontpage that is a week old (31.12.23), giving the articles enough time to gather posts (as in our training data).

```
In [9]: # pick test date one week ago
        test_date = datetime.now() - timedelta(days=7)
```

Small regression Model

First is a small-ish Regression network.

- `drop = 0.3`, as we found that higher dropout values took extremely long to converge and did not really help with the overfitting problem. Our training loss was always around 20% higher than our test-set loss.
- `layer_sizes = [303, 128, 64]`. The idea here is to create a sort of information funnel, where each layer should hopefully pick up on higher-level features in the data, with fewer neurons being responsible for aggregating information from the previous layer. If I were to fantasize what those higher-level-features could be - the 300-dimensional position of the text could represent the overall sentiment, picked up by higher layers? But really we have no idea how the network interprets those vectors. The idea behind having 64 neurons is to have one neuron per bucket- each one maybe contributing `1` to the final output neuron. But again we are speculating, let's see how it performs.

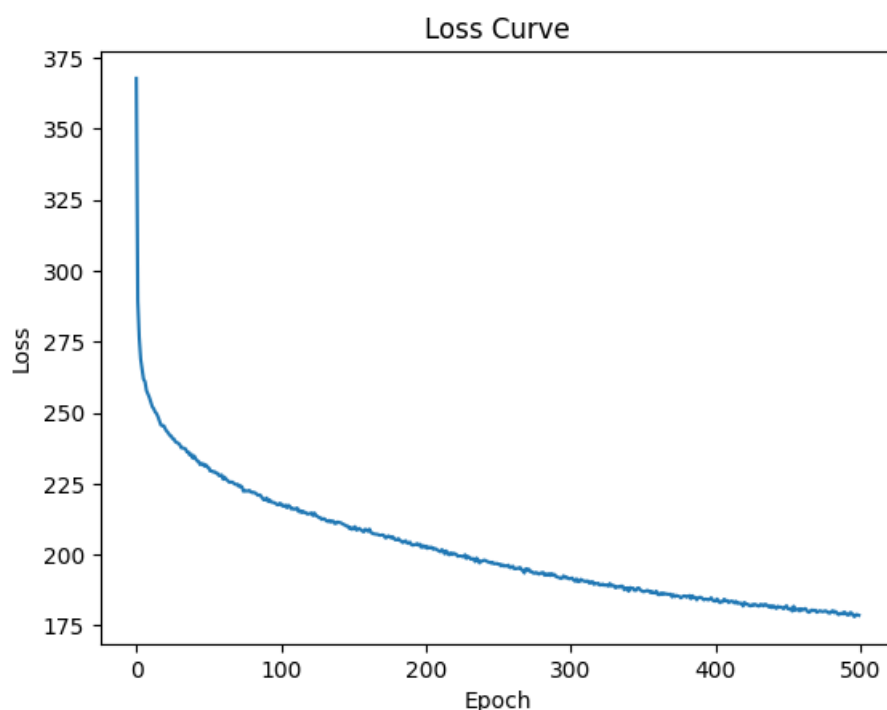
```
In [10]: small_model = RegressionNN(
        drop = 0.3,
        layer_sizes = [303, 128, 64]
        ).to(my_device)

        print('training small model')
        small_losses = train(small_model, epochs=500, learning_rate=0.0001)
        print('evaluating small model')
        eval_trained_model(small_model, small_losses, test_date)
```

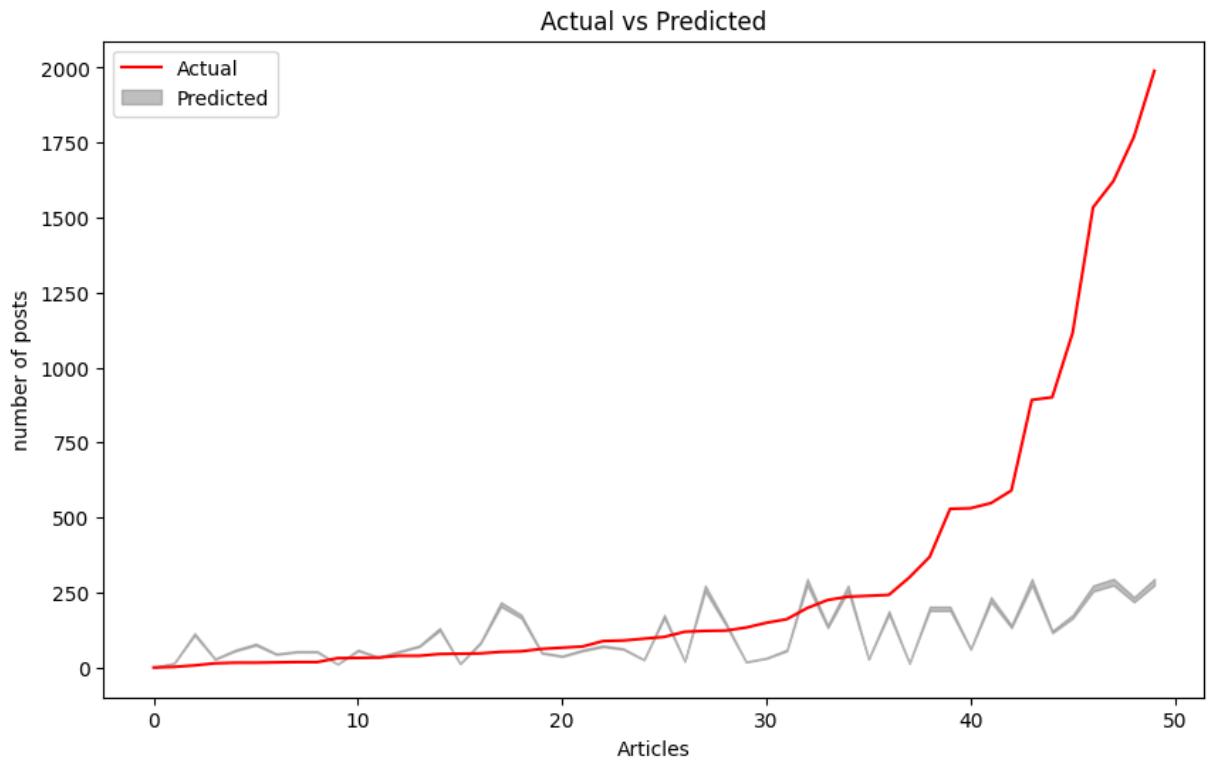
training small model

100%|██████████| 500/500 [1:02:39<00:00, 7.52s/it]

evaluating small model



Test Loss: 217.58059333828228



Correct Prediction: Web- und Games-News: Gebrauchtwagen-Besitzer könnten bald Abogebühren bezahlen Nachrichtenüberblick Das sind die aktuellen Schlagzeilen aus Web und Games

Number of Posts: 0

Predicted Bounds: (0, 1]

Link: <https://www.derstandard.at/story/3000000201431/web-und-games-news-gebrauchtwagen-und-abo>

Correct Prediction: Das Universum als Perspektivgeber Feelgood Manchmal hilft es, Dinge in Relation zu setzen - vom kleinsten Teilchen bis zum galaktischen Supercluster

Number of Posts: 33

Predicted Bounds: (33, 36]

Link: <https://www.derstandard.at/story/3000000200558/das-universum-als-perspektivgeber>

Incorrect Prediction: Klimawandel, Rechtsruck, Inflation: Die Sorgen der Jungen steigen weiter Zukunftsängste Serienweise Krisen, permanenter Leistungsdruck, das Streben nach Besonderem und der ständige Vergleich über soziale Medien belasten junge Menschen zunehmend. Unbeschwert fühlen sich nur wenige. Aber es gibt auch Grund für Optimismus

Number of Posts: 1989

Predicted Bounds: (274, 296]

Link: <https://www.derstandard.at/story/3000000201340/klimawandel-rechtsruck-inflation-die-sorgen-der-jungen-steigen-weiter>

Incorrect Prediction: Gebrauchtwagen-Besitzer könnten bald Abogebühren bezahlen Teure Zukunft Neue Autos verrechnen für bestimmte Features mittlerweile monatliche Kosten. Nun droht dem Gebrauchtwagenmarkt ähnliches

Number of Posts: 1769

Predicted Bounds: (218, 235]

Link: <https://www.derstandard.at/story/3000000201410/gebrauchtwagen-besitzer-koennten-bald-abogebuehren-bezahlen>

I am actually kind of happy with this result. Of course the mean squared error of 217 on the test set is not particularly amazing (seeing as our values only have a range of 64), but the loss curve indicates that the model was actually able to learn some things about the data and was not just randomly guessing. It seems like there is some gradient in the vector space that can be learned.

Furthermore, the plot of articles on (31.12.23) shows us that our bar of predictions is actually pretty slim. With fewer than 64 buckets (maybe only 10), the predictions would probably hit the correct neighborhood more often. The problem here remains the distribution of posts. I have experimented with a smaller number of buckets, but as the vast majority of articles generates between 0 and 10 posts (as indicated by those buckets not actually being buckets but discrete values) - these low values would still overrepresented in the data and the cutoff to articles with higher

post counts would be extremely early. The only alternative here would be to actually drop a lot of values from our test set - which is something that could be tested with more time.

Maybe dropping the majority of low-postcount observations would also help with the overall performance of the model as the current distribution of observations might lead the regression to trend lower than it should - as seen by the lineplot where the actual number of posts (red) trends much higher than the predictions past around ~250 posts.

Medium regression Model 'with a slimmer funnel'

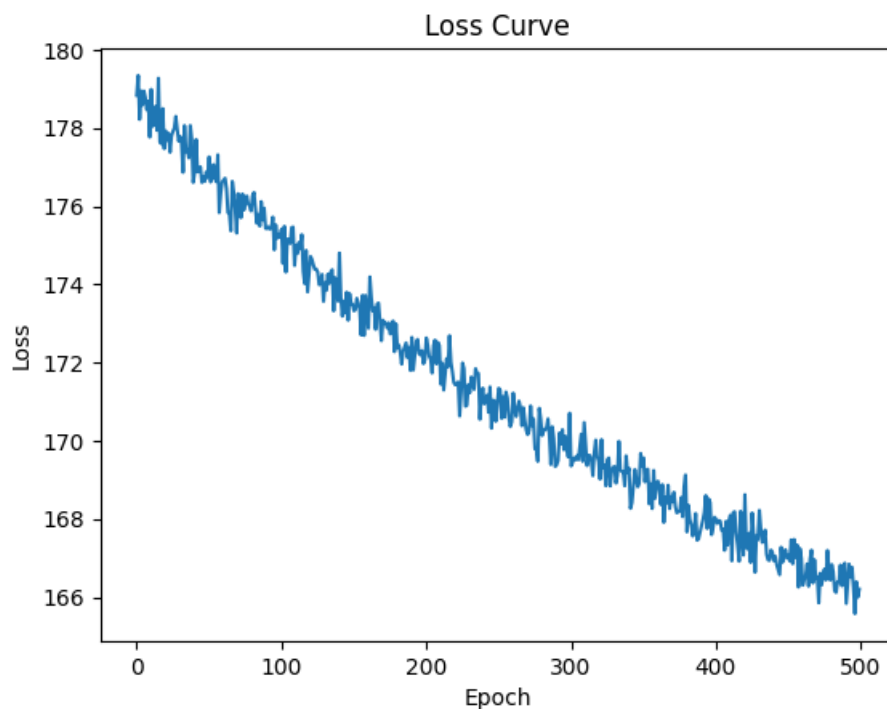
For the next experiment, we will append two final layers with 32 and 16 neurons respectively. Maybe the regression will work better with a longer funnel. To switch things up, we will be testing it on articles from (30.12.23)

```
In [11]: medium_model = RegressionNN(  
    drop = 0.3,  
    layer_sizes = [303, 128, 64, 32, 16]  
).to(my_device)  
  
print('training medium model')  
medium_losses = train(small_model, epochs=500, learning_rate=0.0001)  
print('evaluating medium model')  
eval_trained_model(small_model, medium_losses, test_date-timedelta(days=1))
```

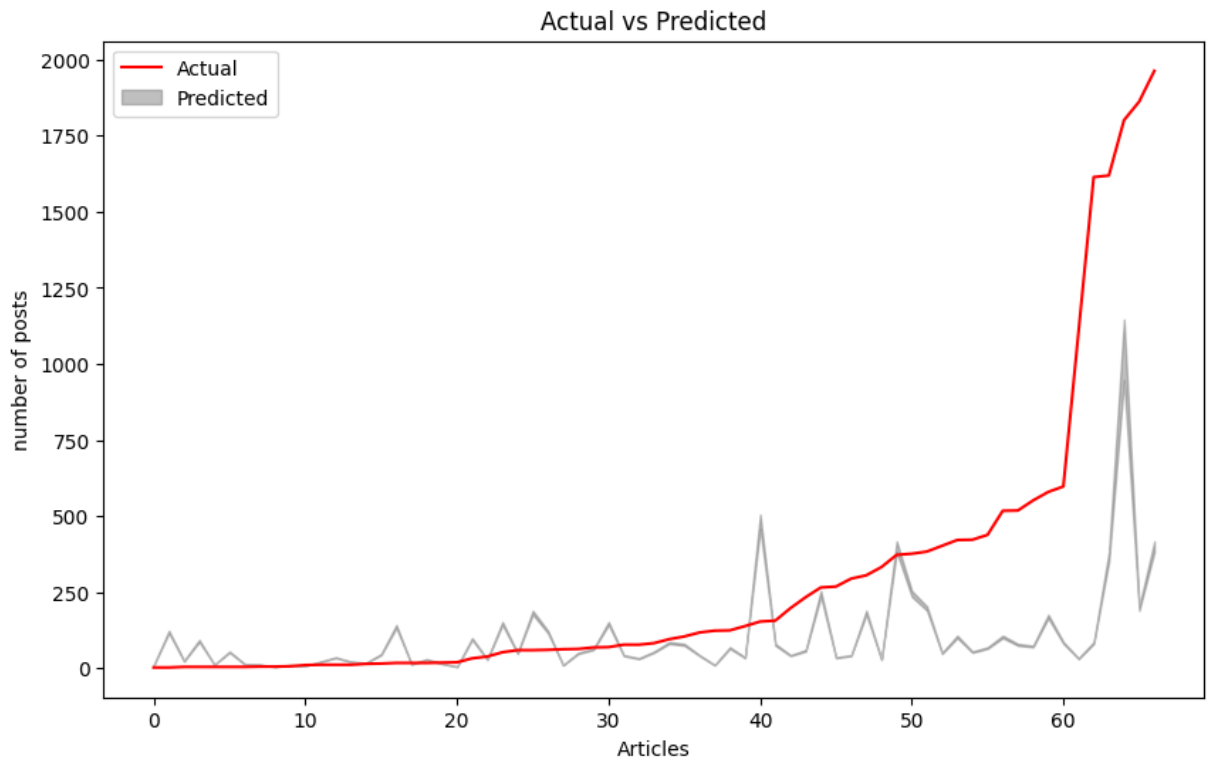
training medium model

100%|██████████| 500/500 [1:03:27<00:00, 7.61s/it]

evaluating medium model



Test Loss: 220.54502462919524



Incorrect Prediction: Signa-Kollaps: Gusenbauer tritt als Millionen-Gläubiger auf Wirtschaft Laut Recherchen von "Profil" und der "Süddeutschen Zeitung" hat der Ex-Kanzler Alfred Gusenbauer (SPÖ) Forderungen in Höhe von mehr als 6,3 Millionen Euro eingebracht

Number of Posts: 1862

Predicted Bounds: (189, 203]

Link: <https://www.derstandard.at/story/3000000201381/signa-kollaps-chef-der-finanzprokuratur-sorgt-sich-um-staatsgeld>

Incorrect Prediction: Frauen werden gleichgestellt, aber nicht gleich behandelt Petra Stuiber 2024 beginnt die Anpassung des Pensionsalters an das der Männer. Sonst aber bleibt alles beim Alten

Number of Posts: 1963

Predicted Bounds: (382, 417]

Link: <https://www.derstandard.at/story/3000000201355/frauen-werden-gleichgestellt-aber-nicht-gleich-behandelt>

Once more we need to accept that bigger is not necessarily better. In this case, extending our network by two smaller layers has not improved our test set performance at all, it is pretty much on par with the small model. Furthermore, for the 30.12., the model was not able to correctly categorize a single article :(

What is remarkable though is how different the loss curve looks this time. Instead of approaching some limit in a logarithmic fashion, the loss (on the training set) seems to decrease linearly. Also the loss starts at a much lower value, probably we just got lucky with the initialization of our weights. Normally this linear decrease should be a good sign that our model has still not found an optimum, but the test-set performance indicates that we were probably just overfitting to the training data (hard to call it overfitting with such a high loss).

We won't be discouraged by this, let's go one step bigger and make an even deeper network before we call it quits. This time we will still keep 64 neurons for the second-to-last layer.

Big regression Model

8 hidden layers, decreasing in size. With more parameters, I have decided to increase the dropout probability to 40%, decrease the learning rate by a decimal place and train for 1.000 epochs.

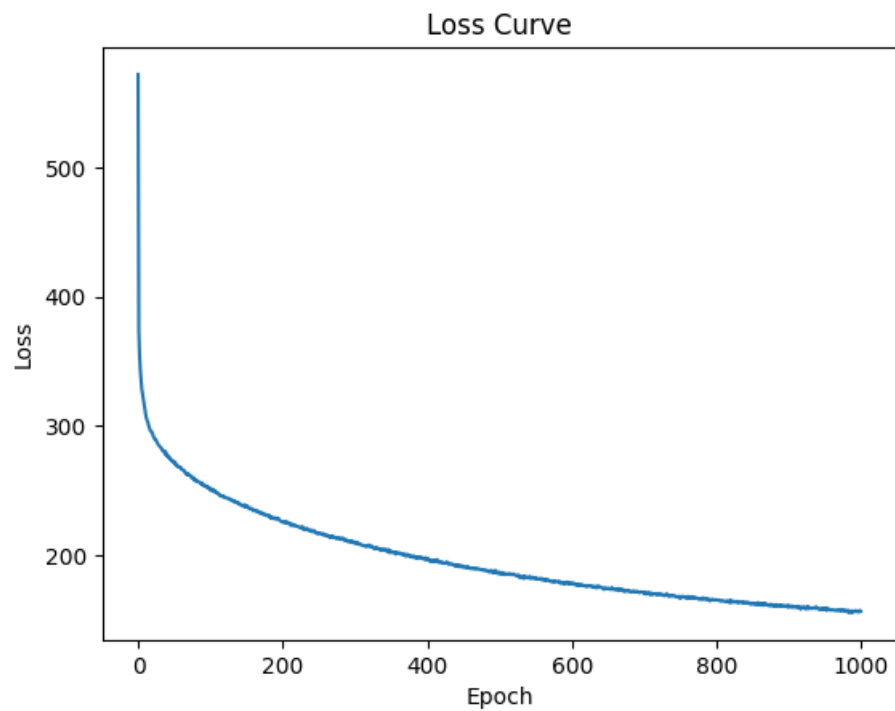
```
In [12]: big_model = RegressionNN(
    drop = 0.4,
    layer_sizes = [303, 303, 256, 256, 256, 128, 128, 64]
).to(my_device)
```

```
print('training big model')
big_losses = train(big_model, epochs=1000, learning_rate=0.00001)
print('evaluating big model')
eval_trained_model(big_model, big_losses, test_date-timedelta(days=2))
```

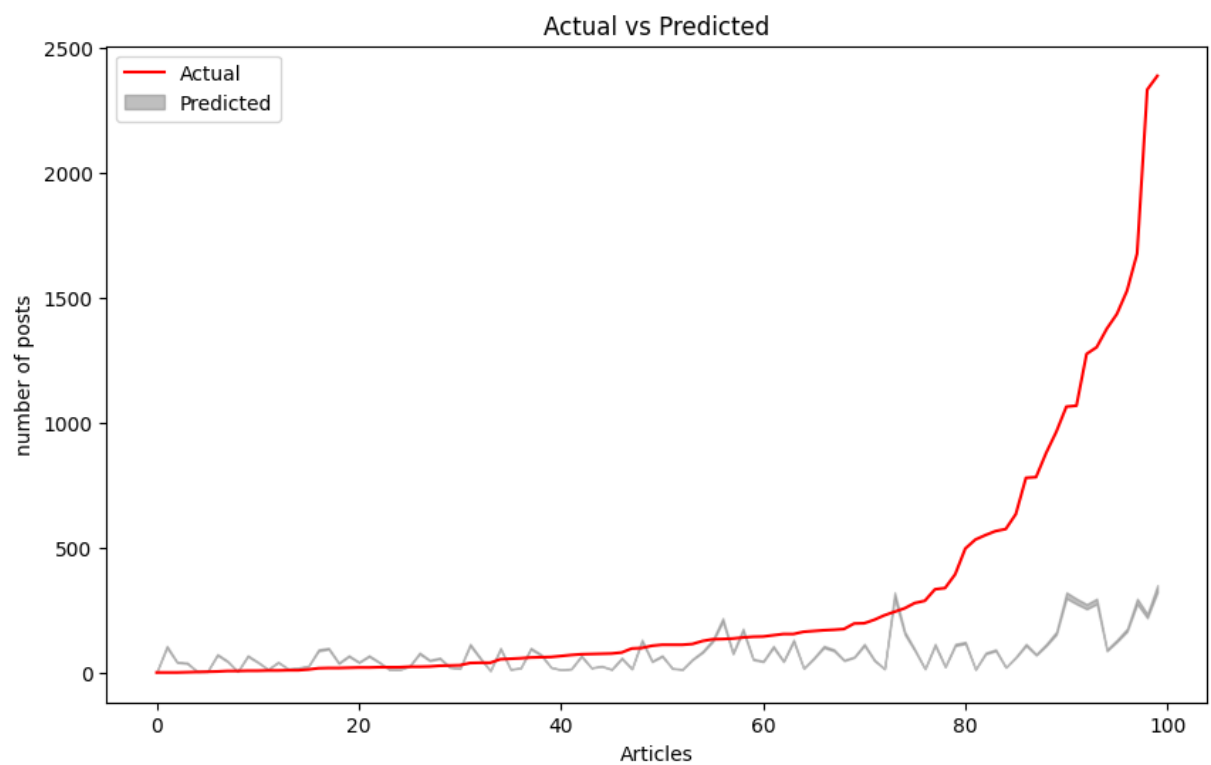
training big model

100%|██████████| 1000/1000 [3:18:36<00:00, 11.92s/it]

evaluating big model



Test Loss: 267.1881389280933



Incorrect Prediction: Bereits 31 Tote nach massiver Angriffswelle auf die Ukraine Krieg in der Ukraine Unter anderem wurden die Städte Kiew, Charkiw, Dnipro, Lwiw und Odessa angegriffen. Die EU beteuerte trotz ungarischen Widerstands ihre Unterstützung

Number of Posts: 2333

Predicted Bounds: (218, 235]

Link: <https://www.derstandard.at/jetzt/livebericht/3000000201236/russische-angriffe-auf-staedte-charkiw-und-lwiw>

Incorrect Prediction: Gegen Woke und Wärmepumpen: Monika Gruber wettet in Büchern – und auf Demos Ausweitung der Kampfzone Eine Bloggerin fühlt sich von Passagen des neuen Buches des Kabarett-Stars rassistisch beleidigt. Die Bayerin wettet auch auf Bühnen gegen grünen, woken Fortschritt

Number of Posts: 2388

Predicted Bounds: (321, 350]

Link: <https://www.derstandard.at/story/3000000201329/gegen-woke-und-waermepumpen-monika-gruber-wettet-in-buechern--und-auf-demos>

Even worse performance than the medium model on the test set. I have a feeling that the vector representation does not capture the message content well enough.

We will perform one last experiment. Scaling our neural networks hidden layers to be wider than the input.

Wide regression Model

Here we go both wide and deep. This model should have over 3 million parameters. Because the last run was overfitting so terribly, I decided to only train this network for 600 epochs, with the same learning rate as before and increased dropout to 50%.

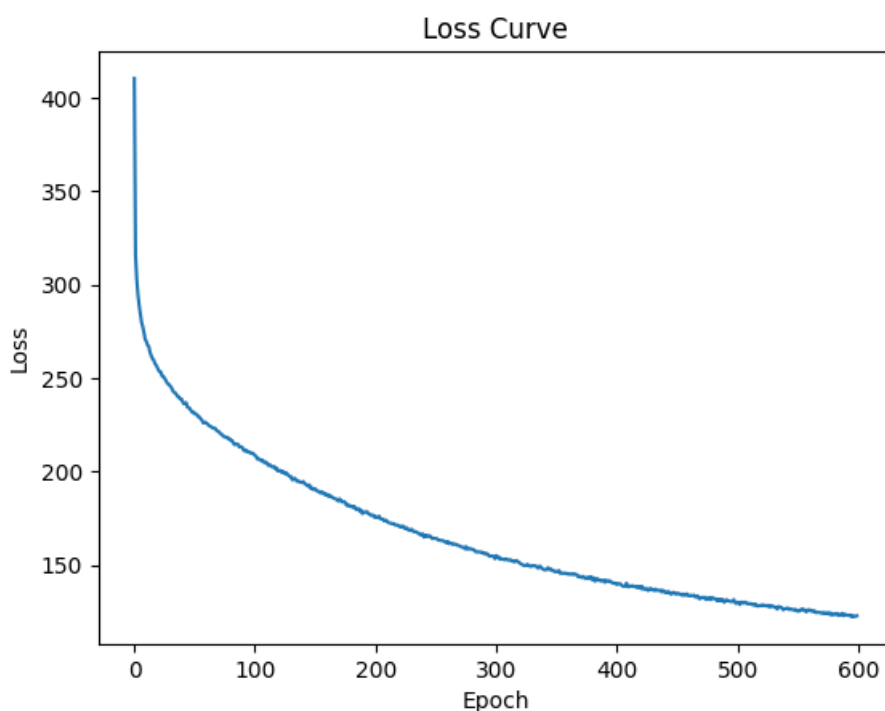
```
In [13]: wide_model = RegressionNN(
          drop = 0.5,
          layer_sizes = [303, 606, 909, 909, 909, 606, 303]
        ).to(my_device)

print('training wide model')
wide_losses = train(wide_model, epochs=600, learning_rate=0.00001)
print('evaluating wide model')
eval_trained_model(wide_model, wide_losses, test_date-timedelta(days=3))
```

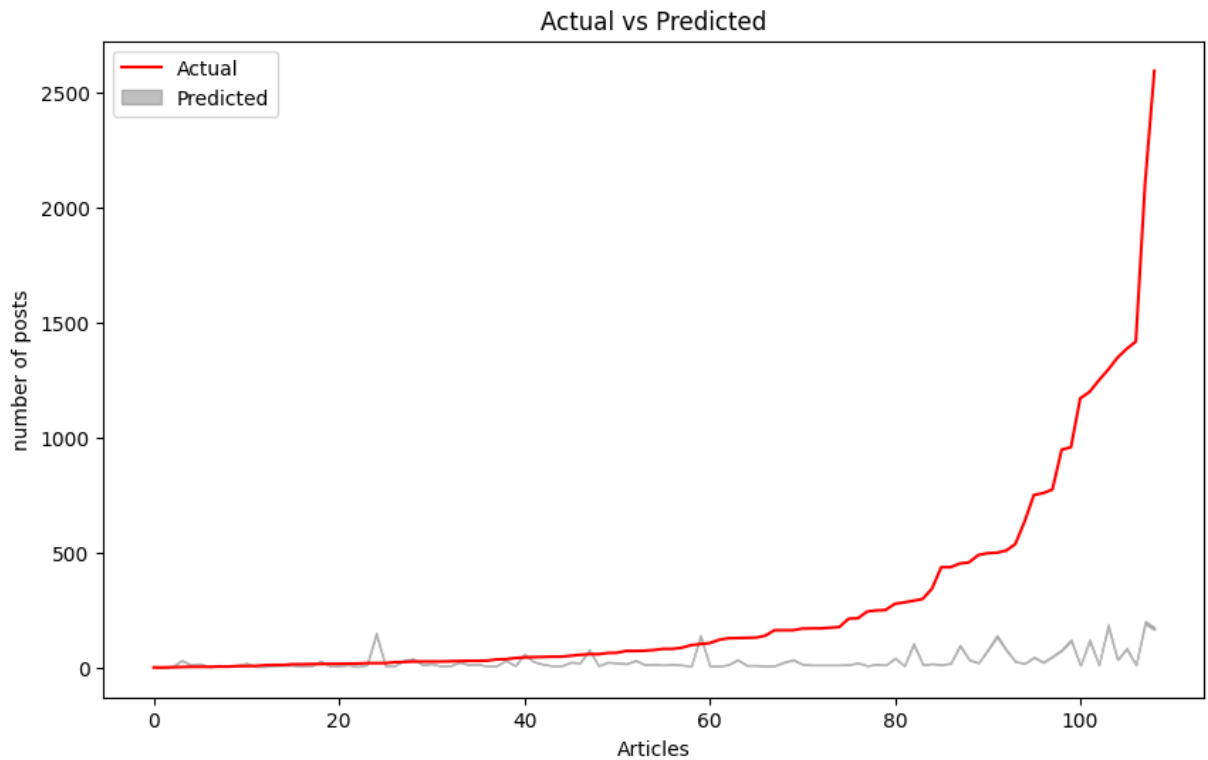
training wide model

100%|██████████| 600/600 [2:49:28<00:00, 16.95s/it]

evaluating wide model



Test Loss: 496.87058120565786



Correct Prediction: Web- und Games-News: Ein Onlyfans-Model gegen die NASA Nachrichtenüberblick Das sind die aktuellen Schlagzeilen aus Web und Games
 Number of Posts: 0
 Predicted Bounds: (0, 1]
 Link: <https://www.derstandard.at/story/3000000201208/web-und-games-news-das-comeback-eines-baby-ichhoernchens>

Correct Prediction: Kreuzworträtsel I 10574 Kreuzworträtsel Täglich neu, exklusiv fürSmart-Abonent:innen:Das kniffligephoenixen-Rätseldes STANDARD
 Number of Posts: 3
 Predicted Bounds: (3, 4]
 Link: <https://www.derstandard.at/story/3000000200358/kreuzwortraetsel-i-10574>

Incorrect Prediction: Deutschland sichert Ukraine weitere EU-Finanzhilfe zu Krieg in der Ukraine Russland hat ukrainischen Angaben zufolge das Land in der Nacht erneut mit Drohnen angegriffen. Zwei russische Dichter erhielten wegen kritischen Gedichts zum Ukrainekrieg lange Haftstrafen
 Number of Posts: 2594
 Predicted Bounds: (164, 176]
 Link: <https://www.derstandard.at/jetzt/livebericht/3000000201130/erneut-nacht-drohnen-angriff-russland>

Incorrect Prediction: Später in Pension – Fortschritt oder Falle für Frauen? Gleichberechtigung Ab dem neuen Jahr steigt das Pensionsalter der Frauen auf 65 Jahre wie bei den Männern. Das könnte eineWin-win-win-Situationwerden – aber auch viele zu Verliererinnen machen
 Number of Posts: 2104
 Predicted Bounds: (189, 203]
 Link: <https://www.derstandard.at/story/3000000200635/spaeter-in-pension-fortschritt-oder-falle-fuer-frauen>

A test set loss that is almost 2.5x that of the small model. Not very promising signs to just scale the network.

Conclusion

There may be countless reasons why this project was not successful in reliably predicting post counts.

In all fairness, I think this is a really difficult problem. Ultimately it should be impossible to predict how many people will open a given article and feel compelled to add their 5 cents to the discussion. However as Isaac Asimov's psychohistory postulates: *'the laws of statistics as applied to large groups of people could predict the general flow of future events'*, so we decided to give it a shot anyways.

Possible next steps

With even more time, I would try encoding the data differently. Maybe using an RNN architecture, feeding in embedding vectors one at a time. Or possibly instead of using Word2Vec, we could apply a transformer to get our embeddings for the whole set.

Perhaps traditional machine learning techniques could also work, maybe a simple clustering of the current vectors.

I am sceptical that simply changing the activation functions would do much, I experimented a bit with leaky ReLu but found no improvements.

As already mentioned, maybe filtering the data would prevent the network from trending so low for all posts.

Saving the 'best' model

To finish things up, we save the network that was most successful on the test data - the small model. It will also be available for download in the [onedrive link](#)

```
In [18]: torch.save(small_model.state_dict(), 'data/small_model.pth')
```