

# Political Compass AI

## "Fake News," Lies and Propaganda: How to Sort Fact from Fiction

 Search[What is "Fake News"?](#)[Why is this important?](#)[Where do news sources fall on the political bias spectrum?](#)[How do you recognize bias in yourself and the media?](#)[How can you read, listen, and/or share news?](#)

### The Liberal Media

Does the media have a liberal bias?

Here's one research-based answer to the question of liberal bias:

The documentary *The Myth of the Liberal Media: The Propaganda Model of News* uses empirical evidence to look at ownership of the mainstream news media, filters that affect what news gets published, and examples of actual news coverage in order to show that conservative political and corporate interests significantly shape news coverage in the United States.

**The myth of the liberal media : the propaganda model of news**  
by Media Education Foundation  
Publication Date: 1997  
Edward Herman and Noam Chomsky

### New Sources on the Political Spectrum

#### What news sources are left-leaning, centrist, or right-leaning?

There is no completely clear answer to this question because there is no one exact methodology to measure and rate the partisan bias of news sources.

Here are a couple of resources that can help:

- [AllSides](#)

All Sides is a news website that presents multiple sources side by side in order to provide the full scope of news reporting.

The [Allsides Bias Ratings page](#) allows you to filter a list of news sources by bias (left, center, right).

AllSides uses a [patented bias rating system](#) to classify news sources as left, center, or right leaning. Components of the rating system include crowd-sourcing, surveys, internal research, and use of third party sources such as Wikipedia and research conducted by [Groseclose and Milyo](#) at UCLA. Note that while the Groseclose & Milyo results are popular, the methodology it is not without [critique](#).

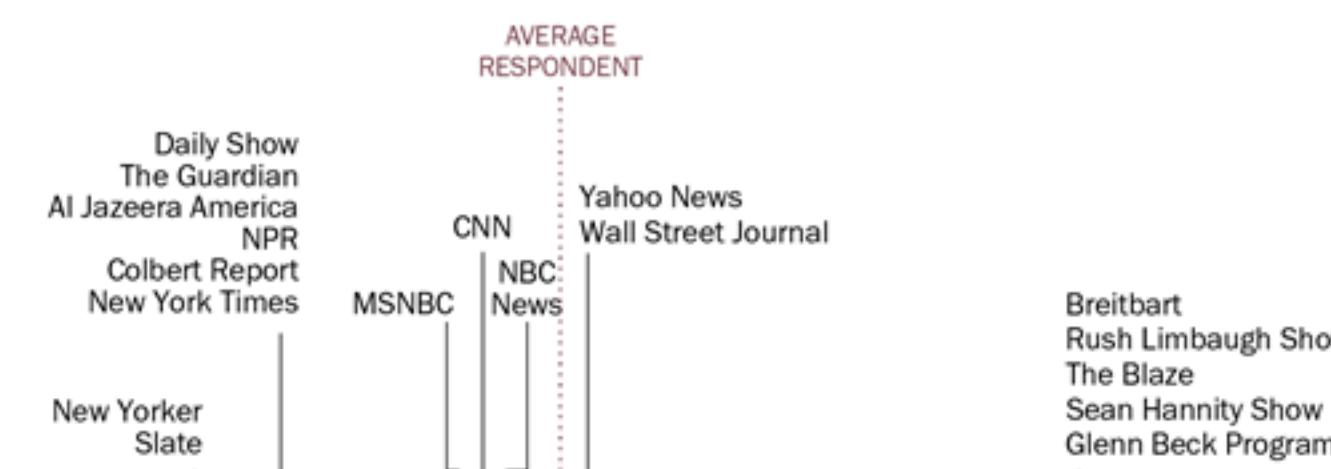
- [Pew Research Center - Political Polarization](#)

Survey data reveals the news source favored by people according to their political beliefs.

A [report based on a 2014 survey](#) shows which news sources are used and considered trustworthy based on individual's political values (liberal or conservative). Note that this report measures the political leanings of the audience rather than the source itself.

#### Ideological Placement of Each Source's Audience

*Average ideological placement on a 10-point scale of ideological consistency of those who got news from each source in the past week...*



## AllSides™ Media Bias Chart™

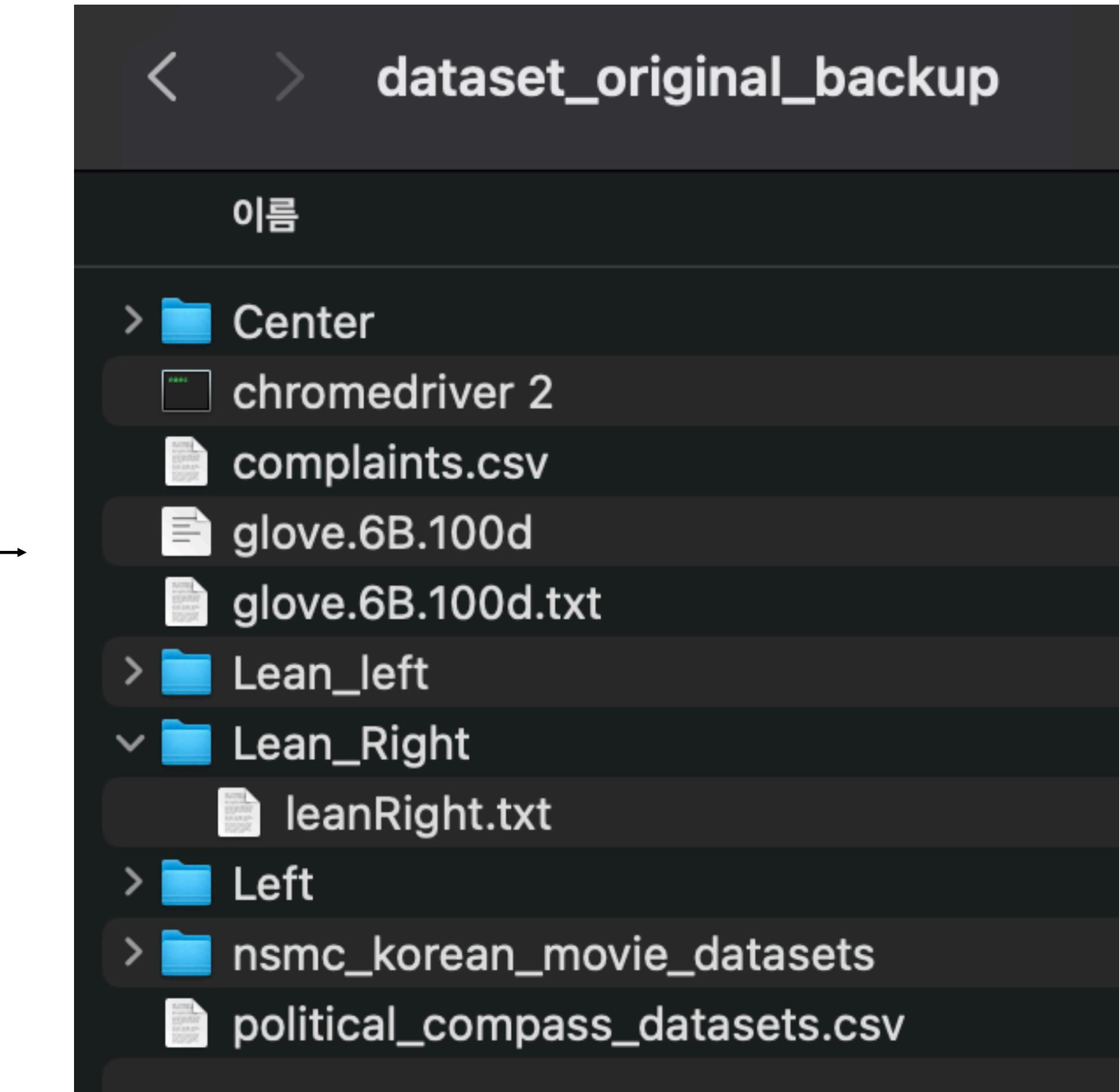
All ratings are based on online content only — not TV, print, or radio content.  
Ratings do not reflect accuracy or credibility; they reflect perspective only.



AllSides Media Bias Ratings™ are based on multi-partisan, scientific analysis.

Visit AllSides.com to view hundreds of media bias ratings.

Version 4.1 | AllSides 2021



Finder 창을 통해 파일 목록을 살펴보겠습니다.

이름	수정일	크기	종류
BERT Text Classification using Keras.ipynb	오늘 오후 1:51	48KB	문서
data	2021년 5월 24일 오후 10:36	--	폴더
complaints.csv	2021년 5월 23일 오전 4:19	1.35GB	CSV 문서
data_pc	어제 오후 5:17	--	폴더
data_political_compass	어제 오후 5:13	--	폴더
dataset_original_backup	어제 오후 5:12	--	폴더
Center	2021년 5월 28일 오후 4:58	--	폴더
chromedriver 2	2021년 3월 13일 오후 12:03	16.7MB	Unix 실행 파일
complaints.csv	2021년 5월 23일 오전 4:19	1.35GB	CSV 문서
glove.6B.100d	2021년 5월 17일 오전 7:21	713바이트	텍스트 클리핑
glove.6B.100d.txt	2020년 6월 14일 오후 6:44	347.1MB	일반 텍스트 문서
Lean_left	2021년 5월 16일 오후 4:06	--	폴더
Lean_Right	어제 오후 4:23	--	폴더
leanRight.txt	어제 오후 4:23	8KB	일반 텍스트 문서
Left	그저께 오후 9:27	--	폴더
nsmc_korean_movie_datasets	2021년 5월 30일 오후 1:28	--	폴더
political_compass_datasets.csv	어제 오후 5:12	209KB	CSV 문서
political_compass_datasets.csv	어제 오후 5:13	209KB	CSV 문서
datasets.ipynb	어제 오후 5:15	275KB	문서
imdb_movie_review_classification.ipynb	2021년 5월 29일 오전 9:43	60KB	문서
models	어제 오후 5:32	--	폴더
bert_model.h5	어제 오후 5:43	438.2MB	문서
multi_labeled_classification.ipynb	어제 오전 8:50	2.4MB	문서
preprocessing.ipynb	어제 오후 4:25	473KB	문서
saved_model	2021년 5월 24일 오후 10:35	--	폴더
tensorboard_data	어제 오후 5:18	--	폴더
Untitled.ipynb	오늘 오전 10:20	2KB	문서

→ 핵심 개발 코드

→ 데이터셋

→ 학습 결과 AI model

→ 데이터 전처리 코드

→ 학습 과정 모니터링

```

1 absl-py==0.12.0
2 appnope @ file:///opt/concourse/worker/volumes/live/4f734db2-9ca8-4d8b-5b29-
6ca15b4b4772/volume/appnope_1606859466979/work
3 argon2-cffi @ file:///opt/concourse/worker/volumes/live/4afd07c8-7fc3-4a09-6326-
d8c70269eb33/volume/argon2-cffi_1613037490059/work
4 astunparse==1.6.3
5 async-generator==1.10
6 attrs @ file:///tmp/build/80754af9/attrs_1620827162558/work
7 backcall @ file:///home/ktietz/src/ci/backcall_1611930011877/work
8 bleach @ file:///tmp/build/80754af9/bleach_1612211392645/work
9 cached-property==1.5.2
10 cachetools==4.2.2
11 certifi==2020.12.5
12 cffi @ file:///opt/concourse/worker/volumes/live/ca3c0fea-9d0f-4489-6889-
89b42078ed37/volume/cffi_1613246931142/work
13 chardet==4.0.0
14 click==8.0.1
15 cycler==0.10.0
16 decorator @ file:///tmp/build/80754af9/decorator_1621259047763/work
17 defusedxml @ file:///tmp/build/80754af9/defusedxml_1615228127516/work
18 entrypoints==0.3
19 filelock==3.0.12
20 flatbuffers==1.12
21 gast==0.4.0
22 google-auth==1.30.0
23 google-auth-oauthlib==0.4.4
24 google-pasta==0.2.0
25 grpcio==1.34.1
26 h5py==3.1.0
27 huggingface-hub==0.0.8
28 idna==2.10
29 importlib-metadata @ file:///opt/concourse/worker/volumes/live/4e25a0be-45cc-4b73-4314-
4604f00e30a4/volume/importlib-metadata_1617877365451/work
30 ipykernel @ file:///opt/concourse/worker/volumes/live/73e8766c-12c3-4f76-62a6-
3dea9a7da5b7/volume/ipykernel_1596206701501/work/dist/ipykernel-5.3.4-py3-none-any.whl
31 ipython @ file:///opt/concourse/worker/volumes/live/f33fa11b-d908-43e8-693a-
f4b58d8c695c/volume/ipython_1617120878195/work
32 ipython-genutils @ file:///tmp/build/80754af9/ipython_genutils_1606773439826/work
33 ipywidgets==7.6.3
34 jedi==0.17.0
35 Jinja2 @ file:///tmp/build/80754af9/jinja2_1621238361758/work
36 joblib==1.0.1
37 jsonschema @ file:///tmp/build/80754af9/jsonschema_1602607155483/work
38 jupyter-client @ file:///tmp/build/80754af9/jupyter_client_1616770841739/work
39 jupyter-core @ file:///opt/concourse/worker/volumes/live/a699b83f-e941-4170-5136-
bf87e3f37756/volume/jupyter_core_1612213304212/work
40 jupyterlab-pygments @ file:///tmp/build/80754af9/jupyterlab_pygments_1601490720602/work
41 jupyterlab-widgets==1.0.0
42 keras-nightly==2.5.0.dev2021032900
43 Keras-Preprocessing==1.1.2
44 kiwisolver==1.3.1
45 Markdown==3.3.4
46 MarkupSafe @ file:///opt/concourse/worker/volumes/live/2165c944-5d03-464f-57c1-
1df4e51b3ac3/volume/markupsafe_1621528146825/work
47 matplotlib==3.4.2
48 mistune==0.8.4
49 nbclient @ file:///tmp/build/80754af9/nbclient_1614364831625/work
50 nbconvert @ file:///opt/concourse/worker/volumes/live/d4b0787b-b6c8-4d28-5453-
3381885d5b33/volume/nbconvert_1601914848300/work
51 nbformat @ file:///tmp/build/80754af9/nbformat_1617383369282/work

```

# pip install -r requirements.txt



Files Running Clusters

Select items to perform actions on them.

<input type="checkbox"/>	0	<input type="checkbox"/>	/ Documents / project / AI / BertMultiClassification / bert_finish	Na
			..	
<input type="checkbox"/>			data	
<input type="checkbox"/>			data_pc	
<input type="checkbox"/>			data_political_compass	
<input type="checkbox"/>			models	
<input type="checkbox"/>			saved_model	
<input type="checkbox"/>			tensorboard_data	
<input type="checkbox"/>			BERT Text Classification using Keras.ipynb	
<input type="checkbox"/>			datasets.ipynb	
<input type="checkbox"/>			preprocessing.ipynb	
<input type="checkbox"/>			requirements.txt	

jupyter datasets Last Checkpoint: 20시간 전 (autosaved)

Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [18]: `from pandas import DataFrame`

```
In [19]: center = ["ASIAN EXPORTERS FEAR DAMAGE FROM U.S.-JAPAN RIFT\n Mounting trade friction between the\n U.S. And Japan\n 'They told Reuter correspondents in Asian capitals a U.S.\n Move against Japan might boost protectionist sentiment\n 'But some exporters said that while the conflict would hurt\n them in the long-run, in the short-term Tokyo's loss\n 'The U.S. Has said it will impose 300 mln dlrss on\n imports of Japanese electronics goods on April 17,\n 'Unofficial Japanese estimates put the impact of the tariffs\n at 10 billion dlrss and spokesmen for major electron\n ''We wouldn't be able to do business," said a spokesman for\n leading Japanese electronics firm Matsushita Electr\n ''If the tariffs remain in place for any length of time\n beyond a few months it will mean the complete erosion of\n 'In Taiwan, businessmen and officials are also worried.',\n ''We are aware of the seriousness of the U.S.',\n 'Threat against\n Japan because it serves as a warning to us," said a senior\n Taiwanese trade official who asked\n 'Taiwan had a trade trade surplus of 15.6 billion dlrss last\n year, 95 pct of it with the U.S.',\n 'The surplus helped swell Taiwan's foreign exchange reserves\n to 53 billion dlrss, among the world's largest.",\n ''We must quickly open our markets, remove trade barriers and\n cut import tariffs to allow imports of U.S. Produc\n 'Retaliation," said\n Paul Sheen, chairman of textile exporters &lt;Taiwan Safe Group>.',\n 'A senior official of South Korea's trade promotion\n association said the trade dispute between the U.S. And Japa\n 'Last year South Korea had a trade surplus of 7.1 billion\n dlrss with the U.S., Up from 4.9 billion dlrss in 1985.'\n 'In Malaysia, trade officers and businessmen said tough\n curbs against Japan might allow hard-hit producers of\n 'In Hong Kong, where newspapers have alleged Japan has been\n selling below-cost semiconductors, some electronics\n 'But other businessmen said such\n a short-term commercial advantage would be outweighed by\n further U.S. Pressu\n ''That is a very short-term view," said Lawrence Mills,\n director-general of the Federation of Hong Kong Industry\n ''If the whole purpose is to prevent imports, one day it will\n be extended to other sources.',\n 'Much more serious for Hong Kong\n is the disadvantage of action restraining trade," he said.',\n "The U.S. Last year was Hong Kong's biggest export market,\n accounting for over 30 pct of domestically produced e\n 'The Australian government is awaiting the outcome of trade\n talks between the U.S. And Japan with interest and c\n ''This kind of deterioration in trade relations between two\n countries which are major trading partners of ours i\n 'He said Australia's concerns centred on coal and beef,\n Australia's two largest exports to Japan and also signif\n 'Meanwhile U.S.-Japanese diplomatic manoeuvres to solve the\n trade stand-off continue.',\n "Japan's ruling Liberal Democratic Party yesterday outlined\n a package of economic measures to boost the Japanese\n 'The measures proposed include a large supplementary budget\n and record public works spending in the first half o\n 'They also call for stepped-up spending as an emergency\n measure to stimulate the economy despite Prime Minister\n 'Deputy U.S. Trade Representative Michael Smith and Makoto\n Kuroda, Japan's deputy minister of International Trad\n 'VIEILLE MONTAGNE SAYS 1986 CONDITIONS UNFAVOURABLE\n A sharp fall in the dollar price of\n zinc and the deprecia\n 'It said in a statement that the two factors led to a\n squeeze on refining margins and an 18.24 pct fall in sales\n 'Vieille Montagne, which is actively pursuing a\n restructuring program, reported a 198 mln franc net loss, after\n 'BANKERS TRUST &lt;BT> PUTS BRAZIL ON NON-ACCRUAL\n Bankers Trust New York Corp said it has\n placed its approxim\n 'Brazil suspended interest payments on its 68 billion dlrss\n of medium- and long-term debt on February 22.',\n 'U.S. banking regulations do not require banks to stop\n accruing interest on loans until payments are 90 days ove\n 'Assuming no cash payments at current interest rates are\n received for the rest of 1987, Bankers Trust estimated
```

In [20]: `center_df = DataFrame(center, columns=['text'])  
center_df`

Out[20]:

	text
0	ASIAN EXPORTERS FEAR DAMAGE FROM U.S.-JAPAN RI...
1	They told Reuter correspondents in Asian capit...
2	But some exporters said that while the conflic...
3	The U.S. Has said it will impose 300 mln dlr ...
4	Unofficial Japanese estimates put the impact o...
...	...
59	cash flow from its low growth tobacco, but the...
60	It can't therefore use its funds to make acqu...
61	Analysts said the National Distillers' spir...
62	"The distilled spirits business has been in a ...
63	REUTER...

64 rows × 1 columns

In [21]: `# center_label = 0  
center_df['label'] = 'center'  
center_df`

Out[21]:

	text	label
0	ASIAN EXPORTERS FEAR DAMAGE FROM U.S.-JAPAN RI...	center
1	They told Reuter correspondents in Asian capit...	center
2	But some exporters said that while the conflic...	center
3	The U.S. Has said it will impose 300 mln dlr ...	center
4	Unofficial Japanese estimates put the impact o...	center
...	...	...
59	cash flow from its low growth tobacco, but the...	center
60	It can't therefore use its funds to make acqu...	center
61	Analysts said the National Distillers' spir...	center
62	"The distilled spirits business has been in a ...	center
63	REUTER... center	center

64 rows × 2 columns

```
In [22]: lean_left = ["The impeachment inquiry picks up tomorrow where it left off Friday, with a subpoena sent to a White Ho  
"Just this morning, the lawyer for the whistleblower, whose complaint is at the base of this inquiry, says he's rep  
"There's are a lot of moving parts.",  
"Fortunately, NPR's Mara Liasson is here to help.",  
'Good morning.',  
'Good morning, Lulu.',  
'All right.',  
"What's the latest?",  
'Well, the latest is that the lawyer for the first whistleblower is tweeting that he is now representing a second w  
'The first whistleblower only had second and thirdhand knowledge.',  
'We also know that there are subpoenas for White House documents and State Department documents and personnel.'  
"We don't know how cooperative the administration will be with those  
'We know that the former ambassador to Ukraine, Marie Yovanovitch, is  
"She's was removed because the president's personal lawyer Rudy Giulia  
'And we that the latest bombshell - the text messages between the thre  
'One of the people in those exchanges, Kurt Volker, has already testi  
'A second one, Gordon Sondland, is scheduled to testify on Tuesday.'  
'So lots and lots of action here.',  
'I mean, another whistleblower has got to be bad news for President T  
'It seems like we really just have not heard the whole story yet.',  
'No.',  
'This is kind of like a trail of bread crumbs.',  
'Each revelation leads to something else.',  
"Some of them are pretty dramatic, like those text messages, where we  
'It was U.S. diplomats actively pushing Ukraine to investigate Biden  
'And we have this cast of characters - Kurt Volker, special envoy to  
"We have William Taylor, who we haven't heard from yet, but he was th  
'He was uncomfortable about with the whole arrangement.',  
"And then Gordon Sondland, who's this millionaire donor to Trump and  
'What his happening to his support among Republican lawmakers, Mara?'  
'I mean, has anything changed?',  
'So far, most Republican lawmakers are sheltering in place, which is  
'Some of them, like Marco Rubio, suggested that the president was kid  
'Some of them are saying nothing.',  
'But there are a few cracks.',  
"Ben Sasse, Susan Collins have criticized the president's action.",  
'And Mitt Romney, who has been the voice of conscience in the Senate  
'Trump punched back at Romney.',  
"He said, Mitt Romney was a pompous ass who's been fighting me from t  
"And later, he said, Romney was a fool who's playing right into the h  
"So we don't know yet if Romney and Ben Sasse and Susan Collins are g  
'Or is this the start of something new?',  
"That's NPR political correspondent Mara Liasson.",
```

```
In [23]: lean_left_df = DataFrame(lean_left, columns=['text'])  
# lean_left label = -1  
lean_left_df['label'] = 'lean_left'  
lean_left_df
```

Out[23]:

	text	label
0	The impeachment inquiry picks up tomorrow wher...	lean_left
1	Just this morning, the lawyer for the whistlebl...	lean_left
2	There's are a lot of moving parts.	lean_left
3	Fortunately, NPR's Mara Liasson is here to help.	lean_left
4	Good morning.	lean_left
...	...	...
741	I appreciate it.	lean_left
742	The economy is humming, but we may be hearing ...	lean_left
743	The stock market fell 3% this past Wednesday b...	lean_left
744	While in the bond market, a key marker of a co...	lean_left
745	Now the payoff for long-term government debt f...	lean_left

746 rows × 2 columns

```
In [25]: lean_right = ['The Pentagon defended the military's diversity and inclusion training programs amid a torrent of crit  
'"We certainly respect the oversight that Congress provides.',  
'I'm not going to comment on any specific one initiative that members of Congress might be doing\u200b," Kirby said  
'"What I can speak to is what we're really focused on here at the department and that's defending the nation.',  
'And that means putting in place the right resources, the right strategies, the right operational concepts to do th  
'"And the \u200bsecretary has been very clear and fairly unapologetic about the fact that we want to get all the be  
'If you \u200bmeet the standards and you're qualified to be in the military – and you're willing to raise your hand  
(Sen. Tom Cotton (left) and Rep. Dan Crenshaw (Getty Images))\n\nRep. Dan Crenshaw (R-Texas) and Sen. Tom Cotton (L  
'began a campaign last week that seeks to have military members report th  
'"Enough is enough.',  
'We won't let our military fall to woke ideology.',  
'We have just launched a whistleblower webpage where you can submit your  
'Your complaint will be legally protected, and go to my office and @SenT  
'T\u200bhe two GOP lawmakers said they were prompted to act after Cotton  
'CLICK HERE TO GET THE FOX NEWS APP\nDefense Secretary Lloyd Austin ad  
'Cruz (R-Texas) accused the Army of portraying soldiers as "pansies" in a  
'Austin said America's foes "would like to capitalize on talking points" in  
'Rep. Alexandria Ocasio-Cortez posted a photo of her grandmother's Puerto  
'"Just over a week ago, my abuela fell ill.',  
'I went to Puerto Rico to see her- my 1st time in a year+ bc of COVID.',  
'This is her home.',  
'Hurricane María relief hasn't arrived.',  
'Trump blocked relief $ for PR," the New York Democrat wrote on Twitter.  
'"People are being forced to flee ancestral homes, & developers are taking  
'Conservative commentator Matt Walsh shot back: "Shameful that you live in  
'"You don't even have a concept for the role that 1st-gen, first-born da  
'"My abuela is okay.',  
'But instead of only caring for mine & letting others suffer, I'm calling  
'Tuesday the Supreme Court issued two more unanimous decisions in Garland  
'This follows two unanimous decisions last week.',  
'The weekly display of unanimity is notable given the calls by Democratic  
'This week I wrote on my blog about how the heavy-handed campaigns might  
'As we await important and likely divided decisions on issues like abort  
'In the Garland case, the court ruled (again) unanimously to reverse the  
'In Cooley, the court unanimously ruled in an opinion by Justice Stephen  
'SUPREME COURT PREPARES FINAL PUSH TO RELEASE HOT-BUTTON RULINGS, AMID RI  
'Justice Sonia Sotomayor wrote the opinion in United States v. Palomar-S  
'It also ruled unanimously in Territory of Guam v. United States, in an o  
'The court ruled in favor of Guam on the collection of funding from the I  
'Op-ed calls on Justice Breyer to retire so Biden can choose Supreme Cou  
'Recently, Breyer warned against any move to expand the Supreme Court.',
```

```
In [30]: lean_right_df = DataFrame(lean_right, columns=['text'])  
# lena_right_label = 1  
lean_right_df['label'] = 'lean_right'  
lean_right_df
```

Out[30]:

	text	label
0	Over Memorial Day weekend, Michael Flynn, who ...	lean_right
1	During a Q&A session, an audience member asked... Flynn replied, "No reason.	lean_right
2	I mean, it should happen here."\\n\\nThe crowd c...	lean_right
3	I bring this up not to dwell on the fact that ...	lean_right
4	Flynn, a retired general, has been saying loon...	lean_right
5	No, I bring this up to ask a different questio...	lean_right
6	The timing of Flynn's remarks was darkly fortu...	lean_right
7	Flynn's comments also coincided with a spirite...	lean_right
8	Some even argue that this is what conservatism...	lean_right
9	It's certainly true that the Trump era has rev...	lean_right
10	But instead of going down various intellectual...	lean_right
11	The core problem afflicting the right—and to a...	lean_right
12	Definitions of populism vary, but for our purp...	lean_right
13	The defining emotion of populism and mobs alik...	lean_right
14	"Patriotism is when love of your own people co...	lean_right
15	"The people of Nebraska are for free silver, a...	lean_right
16	"I will look up the arguments later."\\n\\nBoth ...	lean_right
17	Andrew Jackson, William Jennings Bryan, Huey L...	lean_right
18	Conservatives deserve special criticism for fo...	lean_right
19	But that decision has less to do with conserva...	lean_right
20	Certainly, on paper, intellectual progressivis...	lean_right
21	The surrender to populism was greased by what ...	lean_right
22		

```
In [24]: left = ['No \'civil war\' here: Republicans are at peace in their embrace of Trump\nGreg Krieg\nAnalysis by Gregor  
"House Republicans booted Rep. Liz Cheney from their leadership trio this week, replacing the Wyoming conservati  
"Cheney's defenestration and Stefanik's subsequent ascent were an anticlimax, and not just because the switch-a-i  
"Unlike Cheney, Stefanik is a reliable messenger on the one issue above all else that unites Republican lawmakers  
'After the dust settled on Friday, House Republicans so  
"With Cheney sidelined, they argued, the party would be  
"Those imperatives, as Minority Leader Kevin McCarthy v  
'If the Cheney episode is instructive in any meaningful  
"Rather, to the extent there is conflict within the par  
'The Trumpist forces are on the march, in Congress and  
"Cheney's resistance, for all its sound and fury, attrac  
'Instead, it provided yet another platform for Trump to  
'For her part, Stefanik marched out to put a bow on the  
'The remark, in the context of what had unfolded over t  
'But beneath its apparent absurdity, a bright, glowing  
"Republicans are not only in near-perfect lockstep with  
'Trumpism without Trump\nIn a Friday evening interview  
'"We've got to be able to tell people you can trust us  
"That her replacement, Stefanik, is by any measure a le  
'Those willing to stand by Cheney mostly shared at leas  
'With the exception of Rep. Adam Kinzinger of Illinois,  
'A third bucket exists mostly in the states, where GOP  
"Asked in March why he supported the legislation, given  
'After conflating the two, Raffensperger claimed that t  
'In her interview with Tapper on Friday, Cheney describ  
'"Our system held," Cheney continued, "the institutions  
"That Raffensperger didn't bend to Trump's will is to h  
'That he is now parlaying that credibility to give cove  
"Rep. Claudia Tenney offered up a similar kind of misdi  
'"No one knows about what happened in the election," th  
'"We don't know if it was stolen or not, (Cheney) does  
'But what I know is we need to fix it."'  
"As anti-democratic dissembling goes, Tenney's existential journey might seem mild.  
'But it is, in many ways, the most dangerous kind.'  
'The loudest voices might get the most attention, but it is the measured ones, when peddling an outright lie abou  
"Marjorie Taylor Greene and the new GOP normal\nSome of Tenney's colleagues have been more aggressive -- and less  
"And they are not, as some apologists are quick to suggest, simply acting out of a fear of backlash from Trump's  
'That would be impossible, in a way, because they are, in the most literal sense, representative of them.',  
'None more, of course, than Rep. Marjorie Taylor Greene, a freshman Republican from Georgia.',  
"Earlier on Friday, two days after Greene personally confronted Democratic Rep. Alexandria Ocasio-Cortez of New Y  
"The New Yorker's door was locked, but the group camped out, rattling the letter slots, scribbling nasty messages  
'"We're going to go see, we're going to visit, Alexandria Ocasio-Cortez.',  
'Crazy eyes.',  
'Crazy eyes.',  
'Nutty.',  
"Cortez," Greene says at one point, mispronouncing "Ocasio," before continuing to taunt and yell at staffers ins:  
"By McCarthy's logic -- the one that guided his support for Cheney's ouster -- either of the two episodes involv  
...']
```

```
In [26]: left_df = DataFrame(left, columns=['text'])  
# left label = -2  
left_df['label'] = 'left'  
left_df
```

Out[26]:

	text	label
0	No 'civil war' here: Republicans are at peace ...	left
1	House Republicans booted Rep. Liz Cheney from ...	left
2	Cheney's defenestration and Stefanik's subsequ...	left
3	Unlike Cheney, Stefanik is a reliable messenge...	left
4	After the dust settled on Friday, House Republ...	left
...	...	...
470	A presidential task force on Covid-19 has alre...	left
471	Last month, Nigeria restricted travel from Bra...	left
472	AnNigeria's first batch of the AstraZeneca va...	left
473	(Abraham Archiga/Reuters)\n\nMuslims perform ...	left
474	(Sunday Alamba/AP	left

475 rows × 2 columns

IN THE END, THE WOMEN IN THE UNITED STATES HAVE ACCUSED THE GOVERNMENT OF THE

```
In [28]: right_df = DataFrame(right, columns=['text'])
          # right lable = 2
          right_df['label'] = 'right'
          right_df
```

Out[28]

		text	label
0		tomorrow that will boost federal spending by 2...	right
1	POLITICAL EDITOR FOR DAILYMAIL.COM\n\nPUBLISHE...		right
2	The budget would have the nation continue runn...		right
3	It shows the Biden administration's desire to ...		right
4	Taxes are also set to increase by \$3trillion o...		right
...	...	...	...
212	Alyssa McGrath (pictured) is one of two aides ...		right
213	She alleges Cuomo grabbed her face and kissed ...		right
214	Vill, who said she felt uncomfortable at the t...		right
215	The same photos appear on Cuomo's Flickr accou...		right
216	None of the women in the other photos have acc...		right

217 rows × 2 columns

데이터셋 최종 준비 완료!

```
In [38]: from pandas import DataFrame
import pandas as pd
dataset =pd.concat([center_df, lean_left_df, left_df, lean_right_df, right_df])
dataset
```

Out [38]:

	text	label
0	ASIAN EXPORTERS FEAR DAMAGE FROM U.S.-JAPAN RI...	center
1	They told Reuter correspondents in Asian capit...	center
2	But some exporters said that while the conflic...	center
3	The U.S. Has said it will impose 300 mln dlrs ...	center
4	Unofficial Japanese estimates put the impact o...	center
...	...	...
212	Alyssa McGrath (pictured) is one of two aides ...	right
213	She alleges Cuomo grabbed her face and kissed ...	right
214	Vill, who said she felt uncomfortable at the t...	right
215	The same photos appear on Cuomo's Flickr accou...	right
216	None of the women in the other photos have acc...	right

1548 rows × 2 columns

```
In [40]: # 저장하기
dataset.to_csv("political_compass_datasets.csv")
```

# Keras를 사용한 BERT 텍스트 분류

```
In [4]: import tensorflow as tf
import tensorflow_hub as hub
import pandas as pd
from sklearn.model_selection import train_test_split
import numpy as np
import re
import unicodedata
import nltk
from nltk.corpus import stopwords
import keras
from tqdm import tqdm
import pickle
from keras.models import Model
import keras.backend as K
from sklearn.metrics import confusion_matrix,f1_score,classification_report
import matplotlib.pyplot as plt
from keras.callbacks import ModelCheckpoint
import itertools
from keras.models import load_model
from sklearn.utils import shuffle
from transformers import *
from transformers import BertTokenizer, TFBertModel, BertConfig
```

```
In [5]: def unicode_to_ascii(s):
    return ''.join(c for c in unicodedata.normalize('NFD', s) if unicodedata.category(c) != 'Mn')

def clean_stopwords_shortwords(w):
    stopwords_list=stopwords.words('english')
    words = w.split()
    clean_words = [word for word in words if (word not in stopwords_list) and len(word) > 2]
    return " ".join(clean_words)

def preprocess_sentence(w):
    w = unicode_to_ascii(w.lower().strip())
    w = re.sub(r"([?.!,¿])", r" ", w)
    w = re.sub(r'[^\w]+', " ", w)
    w = re.sub(r"^[^a-zA-Z?.!,¿]+", " ", w)
    w=clean_stopwords_shortwords(w)
    w=re.sub(r'@\w+', ' ',w)
    return w
```

```
In [6]: data = pd.read_csv('./data_political_compass/political_compass_datasets.csv')
```

```
In [7]: data.head()
```

```
Out[7]:
```

	Unnamed: 0	text	label
0	0	ASIAN EXPORTERS FEAR DAMAGE FROM U.S.-JAPAN RI...	center
1	1	They told Reuter correspondents in Asian capit...	center
2	2	But some exporters said that while the conflic...	center
3	3	The U.S. Has said it will impose 300 mln dlsr ...	center
4	4	Unofficial Japanese estimates put the impact o...	center

```
In [8]: print('File has {} rows and {} columns'.format(data.shape[0],data.shape[1]))
```

```
File has 1548 rows and 3 columns
```

```
In [9]: # Select required columns  
data = data[['text', 'label']]
```

```
In [10]: data.head()
```

```
Out[10]:
```

	text	label
0	ASIAN EXPORTERS FEAR DAMAGE FROM U.S.-JAPAN RI...	center
1	They told Reuter correspondents in Asian capit...	center
2	But some exporters said that while the conflic...	center
3	The U.S. Has said it will impose 300 mln dlsr ...	center
4	Unofficial Japanese estimates put the impact o...	center

```
In [11]: # Remove rows, where the label is present only ones (can't be split)  
data = data.groupby('label').filter(lambda x : len(x) > 1)  
#data = data.groupby('Product').filter(lambda x : len(x) > 1)
```

```
In [12]: data.head()
```

```
Out[12]:
```

	text	label
0	ASIAN EXPORTERS FEAR DAMAGE FROM U.S.-JAPAN RI...	center
1	They told Reuter correspondents in Asian capit...	center
2	But some exporters said that while the conflic...	center
3	The U.S. Has said it will impose 300 mln dlsr ...	center
4	Unofficial Japanese estimates put the impact o...	center

```
In [13]: # shuffle datasets  
data = data.sample(frac=1).reset_index(drop=True)  
data
```

```
Out[13]:
```

	text	label
0	In 2019, New York saw 265 million visitors wit...	right
1	Yes.	lean left

```
In [14]: data=data.dropna()  
data=data.reset_index(drop=True)  
data = shuffle(data)  
print('Available labels: ',data.label.unique())  
data['text']=data['text'].map(preprocess_sentence)  
  
# Drop NaN values, if any  
# Reset index after dropping the column  
# Shuffle the dataset  
# Print all the unique labels in the da  
# Clean the text column using preproces
```

Available labels: ['left' 'lean\_left' 'right' 'lean\_right' 'center']

```
In [15]:
```

```
print('File has {} rows and {} columns'.format(data.shape[0],data.shape[1]))  
data.head()
```

File has 1548 rows and 2 columns

```
Out[15]:
```

	text	label
1252	tianwen mars orbiter relay signal rover missio...	left
1407	know lulu every chance situation could increas...	lean_left
1394	last month south african health minister zweli...	left
936	recently also sat board log cabin republicans ...	lean_left
999	woman told colleague winter alleged encounter ...	right

```
In [16]:
```

```
1 num_classes=len(data.label.unique())  
2 num_classes
```

```
Out[16]:
```

5

```
In [17]:
```

```
bert_tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")  
bert_model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased',num_labels=num_classes)
```

All model checkpoint layers were used when initializing TFBertForSequenceClassification.

Some layers of TFBertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initialized: ['classifier']  
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
In [18]:
```

```
sent= 'how to train the model, lets look at how a trained model calculates its prediction.'  
tokens=bert_tokenizer.tokenize(sent)  
print(tokens)  
  
['how', 'to', 'train', 'the', 'model', ',', 'lets', 'look', 'at', 'how', 'a', 'trained', 'model', 'calculate', '##s  
' , 'its', 'prediction', '.']
```

```
In [19]:
```

```
tokenized_sequence= bert_tokenizer.encode_plus(sent,add_special_tokens = True,max_length =30,pad_to_max_length = Tru  
return_attention_mask = True)
```

Truncation was not explicitly activated but `max\_length` is provided a specific value, please use `truncation=True` to explicitly truncate examples to max length. Defaulting to 'longest\_first' truncation strategy. If you encode pairs of sequences (GLUE-style) with the tokenizer you can select this strategy more precisely by providing a specific strategy to `truncation`.

```
/Users/Kimkwangil/opt/anaconda3/envs/py38pytorch/lib/python3.8/site-packages/transformers/tokenization_utils_base.p  
y:2104: FutureWarning: The `pad_to_max_length` argument is deprecated and will be removed in a future version, use  
`padding=True` or `padding='longest'` to pad to the longest sequence in the batch, or use `padding='max_length'` to  
pad to a max length. In this case, you can give a specific length with `max_length` (e.g. `max_length=45`) or leave  
max_length to None to pad to the maximal input size of the model (e.g. 512 for Bert).  
warnings.warn(
```

In [20]:

tokenized\_sequence

```
In [21]: bert_tokenizer.decode(tokenized_sequence['input_ids'])
```

**Out[21]:** '[CLS] how to train the model, lets look at how a trained model calculates its prediction. [SEP] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD] [PAD]'

```
In [22]: data['gt'] = data['label'].map({'left':0,'lean_left':1, 'center':2,'lean_right':3, 'right':4})  
data.head()
```

Out [22]:

		text	label	gt
1252	tianwen mars orbiter relay signal rover missio...		left	0
1407	know lulu every chance situation could increas...	lean_left		1
1394	last month south african health minister zweli...		left	0
936	recently also sat board log cabin republicans ...	lean_left		1
999	woman told colleague winter alleged encounter ...		right	4

In [23]:

```
sentences=data['text']
labels=data['gt']
len(sentences), len(labels)
```

**Out[23]:** (1548, 1548)

In [24]:

```
attention_masks= []
```

**for** sent **in** sentences:

```
bert_inp=bert_tokenizer.encode_plus(sent,add_special_tokens = True,max_length =64,pad_to_max_length = True,return_tensors='pt')  
input_ids.append(bert_inp['input_ids'])  
attention_masks.append(bert_inp['attention_mask'])
```

```
input_ids=np.asarray(input_ids)
attention_masks=np.array(attention_masks)
labels=np.array(labels)
```

```
/Users/kimkwangil/opt/anaconda3/envs/py38pytorch/lib/python3.8/site-packages/transformers/tokenization_utils_base.py:2104: FutureWarning: The `pad_to_max_length` argument is deprecated and will be removed in a future version, use `padding=True` or `padding='longest'` to pad to the longest sequence in the batch, or use `padding='max_length'` to pad to a max length. In this case, you can give a specific length with `max_length` (e.g. `max_length=45`) or leave max_length to None to pad to the maximal input size of the model (e.g. 512 for Bert).  
    warnings.warn(
```

```
In [25]: len(input_ids), len(attention_masks), len(labels)
```

**Out [25]:** (1548, 1548, 1548)

```
In [26]: print('Preparing the pickle file.....')
pickle_inp_path='./data_pc/bert_inp.pkl'
pickle_mask_path='./data_pc/bert_mask.pkl'
pickle_label_path='./data_pc/bert_label.pkl'

pickle.dump((input_ids),open(pickle_inp_path,'wb'))
pickle.dump((attention_masks),open(pickle_mask_path,'wb'))
pickle.dump((labels),open(pickle_label_path,'wb'))

print('Pickle files saved as ',pickle_inp_path,pickle_mask_path,pickle_label_path)

Preparing the pickle file.....
Pickle files saved as ./data_pc/bert_inp.pkl ./data_pc/bert_mask.pkl ./data_pc/bert_label.pkl
```

```
In [27]: print('Loading the saved pickle files..')

input_ids=pickle.load(open(pickle_inp_path, 'rb'))
attention_masks=pickle.load(open(pickle_mask_path, 'rb'))
labels=pickle.load(open(pickle_label_path, 'rb'))

print('Input shape {} Attention mask shape {} Input label shape {}'.format(input_ids.shape,attention_masks.shape,lab

Loading the saved pickle files..
Input shape (1548, 64) Attention mask shape (1548, 64) Input label shape (1548,)
```

```
In [28]: train_inp,val_inp,train_label,val_label,train_mask,val_mask=train_test_split(input_ids,labels,attention_masks,test_s

print('Train inp shape {} Val input shape {}\\nTrain label shape {}\\nVal label shape {}\\nTrain attention mask shape {}\\nVal attention mask shape {}')

Train inp shape (1238, 64) Val input shape (310, 64)
Train label shape (1238,) Val label shape (310,)
Train attention mask shape (1238, 64) Val attention mask shape (310, 64)
```

```
In [29]: log_dir='tensorboard_data/tb_bert'
model_save_path='./models/bert_model.h5'

callbacks = [tf.keras.callbacks.ModelCheckpoint(filepath=model_save_path,save_weights_only=True,monitor='val_loss',m

print('\\nBert Model',bert_model.summary())

loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
metric = tf.keras.metrics.SparseCategoricalAccuracy('accuracy')
optimizer = tf.keras.optimizers.Adam(learning_rate=2e-5,epsilon=1e-08)

bert_model.compile(loss=loss,optimizer=optimizer,metrics=[metric])
```

Model: "tf\_bert\_for\_sequence\_classification"

Layer (type)	Output Shape	Param #
bert (TFBertMainLayer)	multiple	109482240
dropout_37 (Dropout)	multiple	0
classifier (Dense)	multiple	3845

Total params: 109,486,085  
Trainable params: 109,486,085  
Non-trainable params: 0

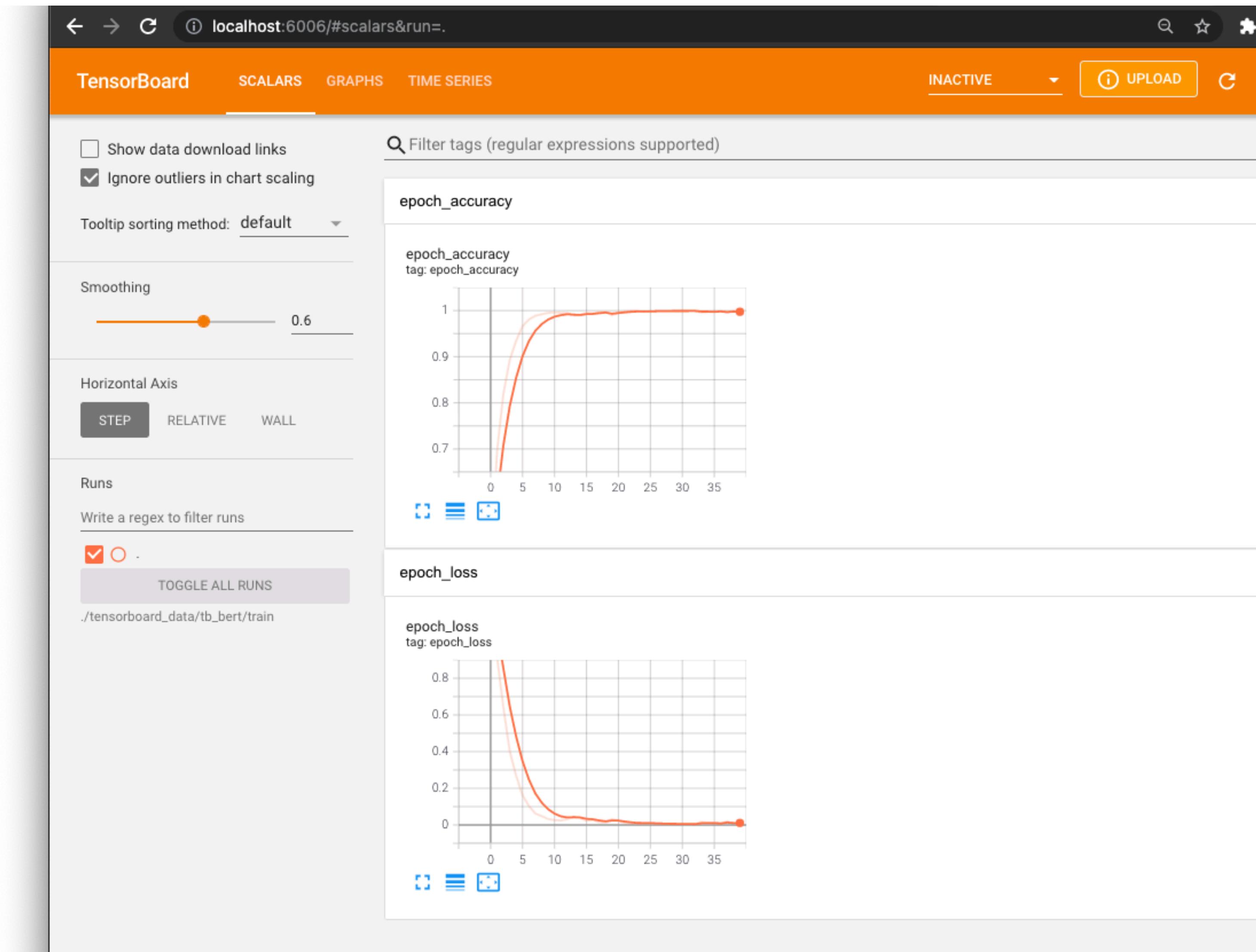
---

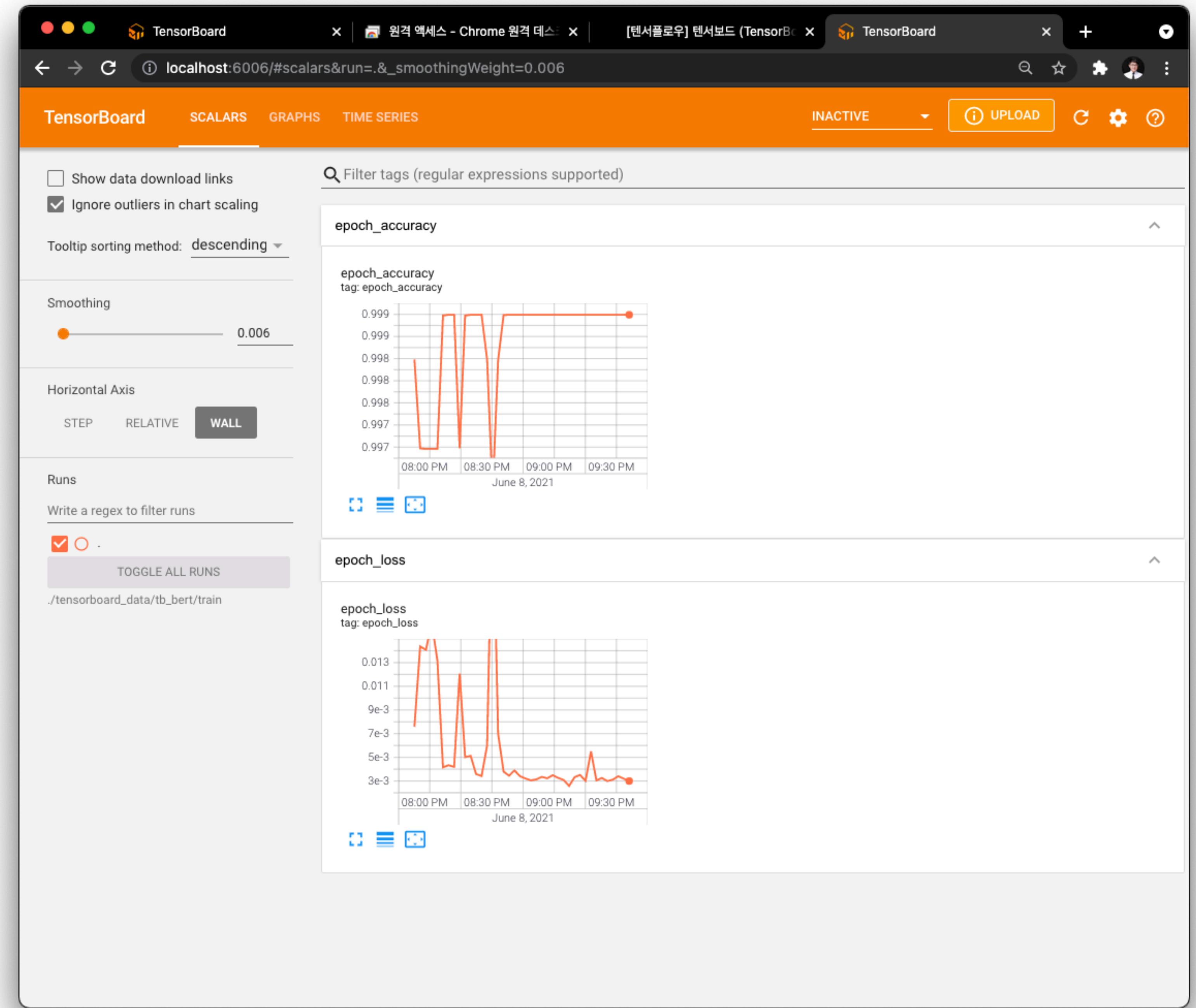
Bert Model None

In [30]:

```
history=bert_model.fit([train_inp,train_mask],train_label,batch_size=32,epochs=20,validation_data=[val_inp,val_mask]

Epoch 1/20
WARNING:tensorflow:The parameters `output_attentions`, `output_hidden_states` and `use_cache` cannot be updated when calling a model.They have to be set to True/False in the config object (i.e.: `config=XConfig.from_pretrained('name', output_attentions=True)`).
WARNING:tensorflow:AutoGraph could not transform <bound method Socket.send of <zmq.sugar.socket.Socket object at 0x7f8e52456460>> and will run it as-is.
Please report this to the TensorFlow team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH_VERBOSITY=10`) and attach the full output.
Cause: module, class, method, function, traceback, frame, or code object was expected, got cython_function_or_method
To silence this warning, decorate the function with @tf.autograph.experimental.do_not_convert
WARNING: AutoGraph could not transform <bound method Socket.send of <zmq.sugar.socket.Socket object at 0x7f8e52456460>> and will run it as-is.
Please report this to the TensorFlow team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH_VERBOSITY=10`) and attach the full output.
Cause: module, class, method, function, traceback, frame, or code object was expected, got cython_function_or_method
To silence this warning, decorate the function with @tf.autograph.experimental.do_not_convert
WARNING:tensorflow:The parameter `return_dict` cannot be set in graph mode and will always be set to `True`.
WARNING:tensorflow:From /Users/kimkwangil/opt/anaconda3/envs/py38pytorch/lib/python3.8/site-packages/tensorflow/python/ops/array_ops.py:5043: calling gather (from tensorflow.python.ops.array_ops) with validate_indices is deprecated and will be removed in a future version.
Instructions for updating:
The `validate_indices` argument has no effect. Indices are always validated on CPU and never validated on GPU.
WARNING:tensorflow:The parameters `output_attentions`, `output_hidden_states` and `use_cache` cannot be updated when calling a model.They have to be set to True/False in the config object (i.e.: `config=XConfig.from_pretrained('name', output_attentions=True)`).
WARNING:tensorflow:The parameter `return_dict` cannot be set in graph mode and will always be set to `True`.
39/39 [=====] - ETA: 0s - loss: 1.2539 - accuracy: 0.4790WARNING:tensorflow:The parameters `output_attentions`, `output_hidden_states` and `use_cache` cannot be updated when calling a model.They have to be set to True/False in the config object (i.e.: `config=XConfig.from_pretrained('name', output_attentions=True)`).
WARNING:tensorflow:The parameter `return_dict` cannot be set in graph mode and will always be set to `True`.
39/39 [=====] - 177s 4s/step - loss: 1.2539 - accuracy: 0.4790 - val_loss: 0.9331 - val_accuracy: 0.6871
Epoch 2/20
39/39 [=====] - 161s 4s/step - loss: 0.8863 - accuracy: 0.7173 - val_loss: 0.7600 - val_accuracy: 0.7742
Epoch 3/20
39/39 [=====] - 161s 4s/step - loss: 0.5830 - accuracy: 0.8393 - val_loss: 0.5819 - val_accuracy: 0.8290
Epoch 4/20
39/39 [=====] - 161s 4s/step - loss: 0.3653 - accuracy: 0.9144 - val_loss: 0.6034 - val_accuracy: 0.8097
Epoch 5/20
39/39 [=====] - 161s 4s/step - loss: 0.2438 - accuracy: 0.9410 - val_loss: 0.6258 - val_accuracy: 0.8323
Epoch 6/20
39/39 [=====] - 161s 4s/step - loss: 0.1565 - accuracy: 0.9580 - val_loss: 0.7288 - val_accuracy: 0.8452
Epoch 7/20
39/39 [=====] - 161s 4s/step - loss: 0.1148 - accuracy: 0.9725 - val_loss: 0.6449 - val_accuracy: 0.8419
Epoch 8/20
39/39 [=====] - 161s 4s/step - loss: 0.0926 - accuracy: 0.9806 - val_loss: 0.6404 - val_accuracy: 0.8452
Epoch 9/20
39/39 [=====] - 161s 4s/step - loss: 0.0559 - accuracy: 0.9903 - val_loss: 0.7248 - val_accuracy: 0.8355
Epoch 10/20
39/39 [=====] - 161s 4s/step - loss: 0.0383 - accuracy: 0.9960 - val_loss: 0.7361 - val_accuracy: 0.8290
Epoch 11/20
39/39 [=====] - 161s 4s/step - loss: 0.0278 - accuracy: 0.9976 - val_loss: 0.7259 - val_accuracy: 0.8290
```





In [59]:

```

1 model_save_path='./models/bert_model.h5'
2
3
4 trained_model = TFBertForSequenceClassification.from_pretrained('bert-base-uncased', num_labels=5)
5 trained_model.compile(loss=loss, optimizer=optimizer, metrics=[metric])
6 trained_model.load_weights(model_save_path)
7
8 preds = trained_model.predict([val_inp, val_mask], batch_size=32)
9 print('preds:', preds)
10 print(type(preds))
11 preds = preds['logits']
12 print('preds_:', preds)
13 # import numpy as np
14 # pred_labels = np.argmax(preds)
15 pred_labels = preds.argmax(axis=1)
16 print('pred_labels', pred_labels)
17 print('val_label', val_label)
18
19 f1 = f1_score(val_label, pred_labels, average='micro')
20 print('F1 score', f1)
21 print('Classification Report')
22 target_names = ['left', 'lean_left', 'center', 'lean_right', 'right']
23 print(classification_report(val_label, pred_labels, target_names=target_names))
24
25 print('Training and saving built model.....')

```

All model checkpoint layers were used when initializing TFBertForSequenceClassification.

Some layers of TFBertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initialized: ['classifier']  
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

WARNING:tensorflow:The parameters `output\_attentions`, `output\_hidden\_states` and `use\_cache` cannot be updated when calling a model.They have to be set to True/False in the config object (i.e.: `config=XConfig.from\_pretrained('name', output\_attentions=True)`).

WARNING:tensorflow:The parameter `return\_dict` cannot be set in graph mode and will always be set to `True`.  
preds: TFSequenceClassifierOutput(loss=None, logits=array([-0.49090445, 3.2869735, -1.013723, -1.1928202, -0.56186223],

```

[ 2.469226 ,  0.00602326, -1.1767939 , -0.8681813 , -0.42022336],
[-0.17980987,  0.7214129 , -0.22669274, -1.2690276 ,  2.09991 ],
...
[ 0.6358856 ,  2.750293 , -1.5966513 , -1.0756427 , -0.36031324],
[ 2.7467015 , -1.3513494 , -0.7086528 , -0.6813501 , -0.1758186 ],
[ 2.412053 , -0.24241978, -1.2128685 , -0.7496565 , -0.6367665 ],
dtype=float32), hidden_states=None, attentions=None)
<class 'transformers.modeling_tf_outputs.TFSequenceClassifierOutput'>
preds_: [[-0.49090445 3.2869735 -1.013723 -1.1928202 -0.56186223]
[ 2.469226  0.00602326 -1.1767939 -0.8681813 -0.42022336]
[-0.17980987  0.7214129 -0.22669274 -1.2690276  2.09991 ]
...
[ 0.6358856  2.750293 -1.5966513 -1.0756427 -0.36031324]
[ 2.7467015 -1.3513494 -0.7086528 -0.6813501 -0.1758186 ]
[ 2.412053 -0.24241978 -1.2128685 -0.7496565 -0.6367665 ]]
```

```

pred_labels [1 0 4 1 0 0 0 1 1 0 1 1 1 4 1 0 0 1 1 1 0 1 1 1 1 1 1 1 0 1 1 1 0 1
1 1 0 0 1 0 1 1 1 1 0 0 1 4 1 1 1 4 1 1 1 1 4 1 1 2 1 0 1 0 1 0 4 0
1 0 1 4 0 4 1 4 0 1 1 1 1 0 0 1 0 0 1 0 1 0 1 0 1 0 0 1 1 4 1 1 0
1 1 0 1 1 4 1 1 1 4 1 1 1 0 1 0 0 1 1 0 0 1 1 1 1 0 0 1 1 1 1 0 1 0 0 0 0
1 1 1 1 0 1 1 4 1 1 1 4 0 4 1 1 1 1 4 0 1 0 4 1 4 0 1 1 0 1 1 1 4 1 1 1
1 0 2 1 1 0 4 1 1 1 0 2 0 1 1 1 0 1 1 1 0 0 0 0 1 0 1 0 1 1 1 4 1 1 4 1
1 1 1 1 1 1 4 1 1 1 1 0 4 0 1 1 1 0 1 0 0 1 1 1 1 0 1 1 1 1 0 0 1 4 1
1 0 1 1 0 1 1 1 1 0 1 1 1 0 1 1 0 1 4 1 0 4 1 0 0 1 4 1 1 0 1 0 1 0 0 0 0
1 1 4 1 1 1 1 0 4 1 1 1 0 0]
```

```

val_label [1 0 4 1 0 0 0 2 1 1 4 1 1 4 1 1 0 1 1 0 0 1 1 1 1 0 1 1 1 3 1 1 1 0 1
0 0 0 0 0 1 0 1 1 0 0 1 4 1 1 1 0 1 4 1 0 4 1 1 2 1 0 1 0 1 0 4 0
1 0 1 4 0 4 1 4 0 1 4 1 1 1 1 0 1 1 0 0 1 3 1 1 0 0 0 1 0 0 1 1 4 1 1 0
1 1 0 1 1 0 1 1 1 4 0 1 1 3 1 0 1 1 1 0 0 1 1 1 2 0 0 1 1 1 1 2 1 0 0 0 0
```

# 인공지능 AI 의 political compass accuracysms 20회 반복 학습한 결과 83% 달성!

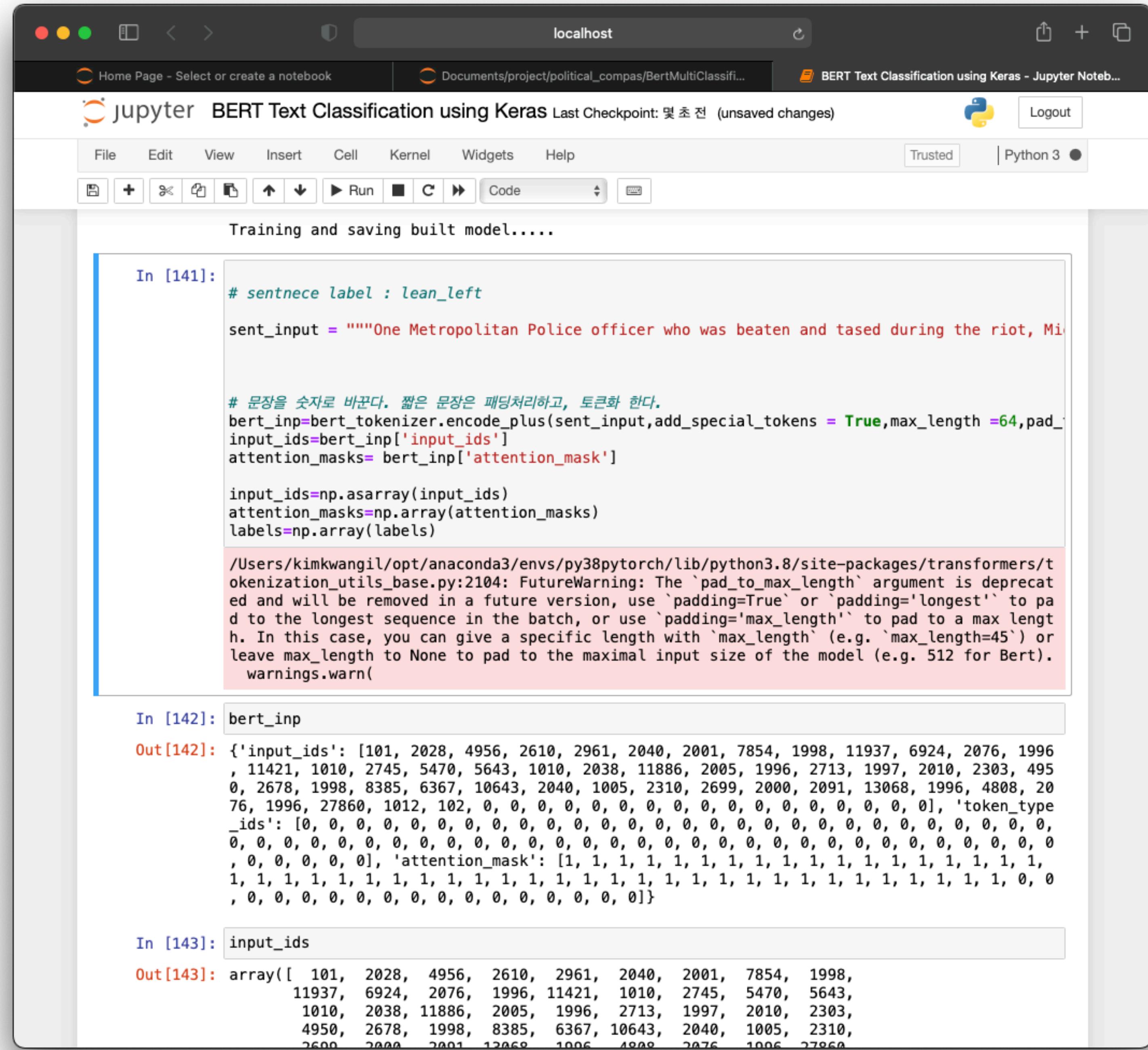
F1 score 0.8290322580645161

Classification Report

	precision	recall	f1-score	support
left	0.84	0.75	0.79	99
lean_left	0.81	0.95	0.88	161
center	1.00	0.38	0.55	8
lean_right	0.00	0.00	0.00	10
right	0.87	0.84	0.86	32
accuracy			0.83	310
macro avg	0.71	0.58	0.61	310
weighted avg	0.81	0.83	0.81	310

Training and saving built model.....

```
/Users/kimkwangil/opt/anaconda3/envs/py38pytorch/lib/python3.8/site-packages/sklearn/metrics/_classification.py:124
8: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted s
amples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
/Users/kimkwangil/opt/anaconda3/envs/py38pytorch/lib/python3.8/site-packages/sklearn/metrics/_classification.py:124
8: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted s
amples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
/Users/kimkwangil/opt/anaconda3/envs/py38pytorch/lib/python3.8/site-packages/sklearn/metrics/_classification.py:124
8: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted s
amples. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
```





localhost

Home Page - Select or create a notebook | Documents/project/political\_compas/BertMultiClassifi... | BERT Text Classification using Keras - Jupyter Noteb...

jupyter BERT Text Classification using Keras Last Checkpoint: 3분 전 (unsaved changes) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [147]:

```
1 plt.hist(pred_labels)
2
3 plt.title("Political Compass AI")
4
5
6 plt.show()
```

Out[147]:

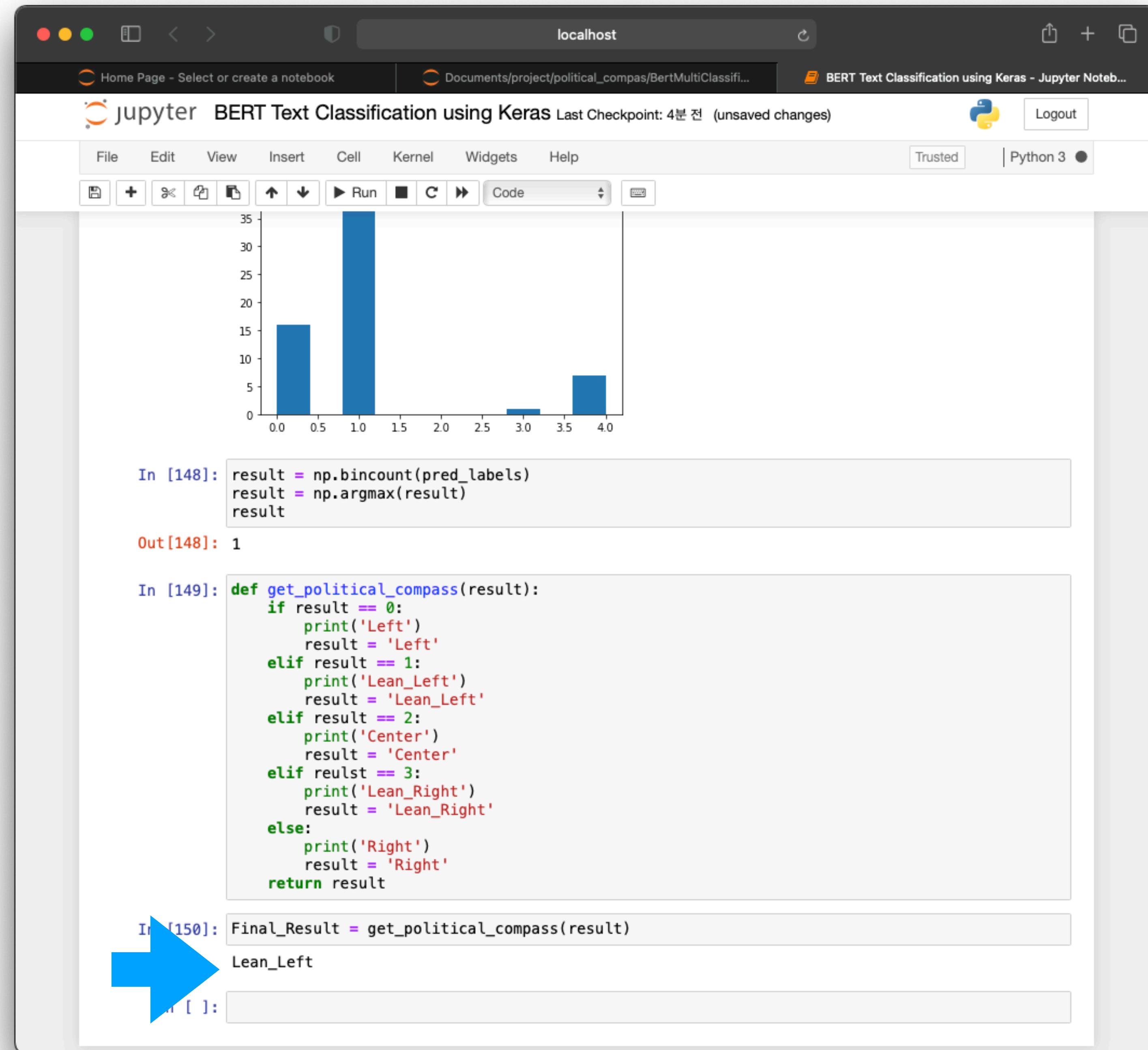
In [148]:

```
result = np.bincount(pred_labels)
result = np.argmax(result)
result
```

Out[148]: 1

In [149]:

```
def get_political_compass(result):
    if result == 0:
        print('Left')
        result = 'Left'
    elif result == 1:
        print('Lean_Left')
        result = 'Lean_Left'
    elif result == 2:
        print('Center')
        result = 'Center'
    elif result == 3:
        print('Lean_Right')
        result = 'Lean_Right'
    else:
        print('Right')
        result = 'Right'
```



## Precision(정밀도)

정밀도란 모델이 True라고 분류한 것 중에서 실제 True인 것의 비율입니다. 즉, 아래와 같은 식으로 표현할 수 있습니다.

$$(Precision) = \frac{TP}{TP + FP}$$

## Recall(재현율)

재현율이란 실제 True인 것 중에서 모델이 True라고 예측한 것의 비율입니다.

**Positive 정답률, PPV(Positive Predictive Value)**, 예를 들어 날씨 예측 모델이 맑다로 예측했는데, 실제 날씨가 맑았는지를 살펴보는 지표

통계학에서는 **sensitivity**으로, 그리고 다른 분야에서는 **hit rate**라는 용어로도 사용합니다. 실제 날씨가 맑은 날 중에서 모델이 맑다고 예측한 비율을 나타낸 지표인데, 정밀도(Precision)와 True Positive의 경우를 다르게 바라보는 것입니다. 즉, Precision이나 Recall은 모두 실제 True인 정답을 모델이 True라고 예측한 경우에 관심이 있으나, 바라보고자 하는 관점만 다릅니다. Precision은 모델의 입장에서, 그리고 Recall은 실제 정답(data)의 입장에서 정답을 정답이라고 맞춘 경우를 바라보고 있습니다. 다음의 경우를 생각해보겠습니다.

$$(Recall) = \frac{TP}{TP + FN}$$

"어떤 요소에 의해, 확실하지 않은 날을 예측할 수 있다면 해당하는 날에만 맑은 날이라고 예측하면 되겠다."

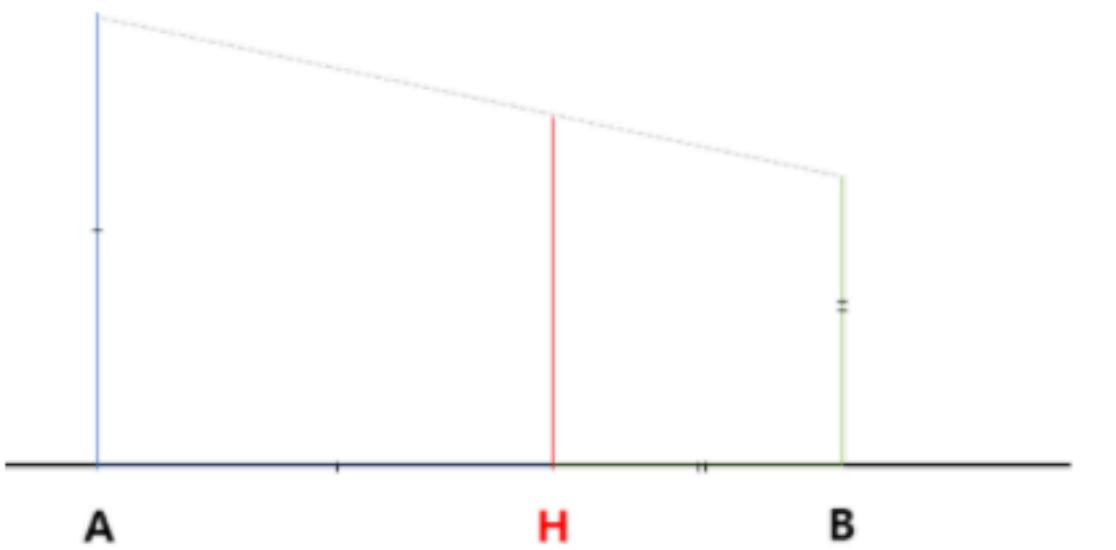
이 경우에는 확실하지 않은 날에는 아예 예측을 하지 않고 보류하여 FP의 경우의 수를 줄여, Precision을 극도로 끌어올리는 일종의 편법입니다. 예를 들어 한 달 30일 동안 맑은 날이 20일이었는데, 확실한 2일만 맑다고 예측한다면, 당연히 맑다고 한 날 중에 실제 맑은 날(Precision)은 100%가 나오게 됩니다. 하지만 과연, 이러한 모델이 이상적인 모델일까요?

따라서, 우리는 실제 맑은 20일 중에서 예측한 맑은 날의 수도 고려해 보아야 합니다. 이 경우에는 Precision만큼 높은 결과가 나오지 않습니다. Precision과 함께 Recall을 함께 고려하면 실제 맑은 날들(즉, 분류의 대상이 되는 정의역, 실제 data)의 입장에서 우리의 모델이 맑다고 예측한 비율을 함께 고려하게 되어 제대로 평가할 수 있습니다. Precision과 Recall은 상호보완적으로 사용할 수 있으며, 두 지표

$$(F1\text{-}score) = 2 \times \frac{1}{\frac{1}{Precision} + \frac{1}{Recall}} = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

F1 score는 Precision과 Recall의 조화평균입니다.

조화평균은 기하학적으로 다음과 같이 표현할 수 있습니다. 서로 다른 길이의 A, B와 이 두 길이의 합만큼 떨어진 변(AB)으로 이루어진 사다리꼴을 생각해봅시다. 이 AB에서 각 변의 길이가 만나는 지점으로부터 맞은 편의 사다리꼴의 변으로 내린 선분이 바로 조화평균을 나타냅니다.



<Fig5. 조화평균의 기하학적 의미>

기하학적으로 봤을 때, 단순 평균이라기보다는 작은 길이 쪽으로 치우치게 된, 그러면서 작은 쪽과 큰 쪽의 사이의 값을 가진 평균이 도출됩니다. 이렇게 조화평균을 이용하면 산술평균을 이용하는 것보다, 큰 비중이 끼치는 bias가 줄어든다고 볼 수 있습니다. 즉, F1-score는 아래와 같이 생각할 수 있습니다.

