Introduction

For quite some time, I've been intrigued by the algorithms capable of predicting weather forecasts and have been on the lookout for a suitable dataset to delve into this field. However, my search led me to a common hurdle—most resources were either paid or provided insufficient data to work with.

Finally, my persistence paid off when I stumbled upon a website that generously offered an extensive archive of weather data for various cities worldwide, all available in open access. That's the **weather diary for schoolchildren by Gismeteo**.

Gismeteo is one of the most popular CIS weather websites. The site was created on December 12, 2000, however weather forecasts have been published with forecasts from the 1990s on BBS and FIDONet. The site provides an API for professional meteorologists.

Frankly speaking, as a forecast, Gismeteo predictions very often are totally unreliable. But as far as they are one of the pioneers in weather forecasts I guess they at least have skills to reliably take readings from the weather stations and other sources.

Let's have fun scraping it ...

Description

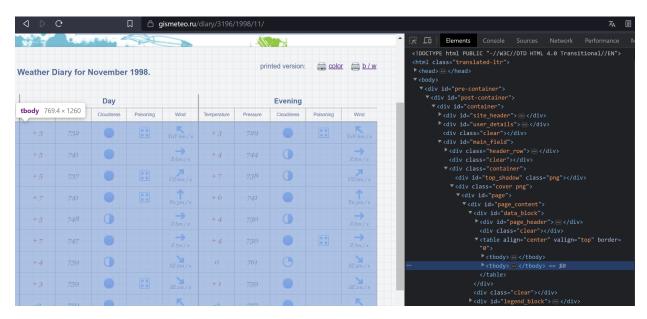
In the picture below there is an example of a weather diary for school children webpage (translated by browser + partially manually).

There is a block of selectors where users can select the country, city, month and year of the diary.

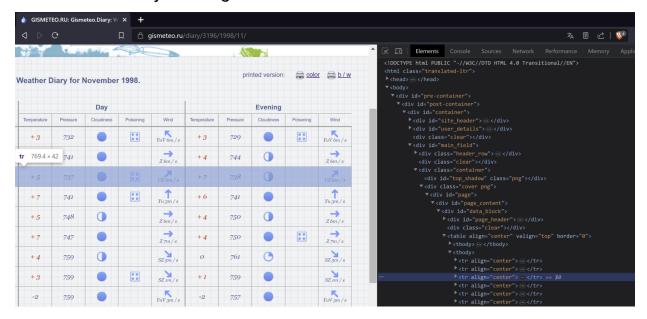


Next there is a table of observations for the selected time period. The data is organized in such a way that every row represents the date and includes the same 5 parameters (Temperature, Pressure, Cloudiness, Weather conditions, Wind) for both day and evening.

Inspection



The data of our interest is located in the element. Every row of the table is defined by **>** tag.



Every data cell is defined by tag.

Day of month, Temperature and Pressure stored as text into **** tag.

Cloudiness and Weather conditions are organized in the form of a gif icon with a unique filename for each factor.

Wind cell contains both icon and text inside for wind direction and wind speed.

```
▼
▼
 ▼<font style="vertical-align: inherit;">
    <font style="vertical-align: inherit;">3</font>
 ▼
 ▼<font style="vertical-align: inherit;">
    <font style="vertical-align: inherit;">+ 5</font>
   </font>
 ▼<font style="vertical-align: inherit;">
    <font style="vertical-align: inherit;">737</font>
   </font>
 <img src="//st6.gismeteo.ru/static/diary/img/dull.png" class="label_icon label_small screen_icon">
   <img src="//st7.gismeteo.ru/static/diary/img/dull-bw.gif" class="label_icon label_small print_icon">
▼
   <img src="//st8.gismeteo.ru/static/diary/img/rain.png" class="label_icon label_small screen_icon">
   <img src="//st4.gismeteo.ru/static/diary/img/rain-bw.png" class="label_icon label_small print_icon">
 <img src="//st5.gismeteo.ru/static/diary/img/w5.gif" class="screen_icon">
    <img src="//st6.gismeteo.ru/static/diary/img/w5-bw.gif" class="print_icon">
   ▼<font style="vertical-align: inherit;">
      <font style="vertical-align: inherit;">UZ 6m / s</font>
    </font>
   </span>
```

The url for every single page (selected month/year) looks like: gismeteo.ru/diary/{cityID}/{year}/{month}

Scraper mechanics

Users define 3 parameters: cityID, start date and finish date.

Scraper creates a list of urls for every month between start and finish date: gismeteo.ru/diary/{cityID}/{year}/{month}.

Then scraper grabs the table cells line by line and appends it to the dataframe.

The first cell **Day** scraper converts into a **Date** adding {month} and {year} parameters.

Next two cells **Temperature** and **Pressure** parsed as text.

Cloudiness and **Weather conditions** (like snow, rain, storm, hail) are recognized by the uniqueness of the .gif icon name and stored in the dataframe as factors.

Last cell **Wind** when parsing is splitted into 2 cells: **Wind direction** and **Wind speed**.

When all the table is parsed scraper continues with the next url from the list and appends new data to the dataframe.

Scraper stops when all the links between start and finish dates are parsed and saves dataframe to .csv file.

Scrapers output

	I== .						I					
Date	D Temperature	D Pressure	D Cloudiness	D Weather Condition		D Wind Speed	N Temperature	N Pressure	N Cloudiness	N Weather Condition		N Wind Speed
1.4.1997	+15		3 Sun		3	8м/с	+8		Sun		3	8м/с
2.4.1997	+11		2 Dull		C3	5м/с	+7		Sun		C3	5м/с
3.4.1997	+13		3 Sun		Ю	6м/с	+8		Mostly Clouds		Ю	6м/с
4.4.1997	+4		Mostly Clouds		C3	9м/с	+1		Sun/Clouds		C3	9м/с
5.4.1997	+5		Mostly Clouds		C3	5м/с	0		Dull	Snow	C3	5м/с
6.4.1997	+2		3 Dull		С	7м/с	0		Sun		С	7м/с
7.4.1997	+1		Mostly Clouds		С	5м/с	-1		Sun		С	5м/с
8.4.1997	+5		Mostly Clouds		C3	6м/с	+2		Mostly Clouds		C3	6м/с
9.4.1997	+5		Dull		С	2м/с	+4		Dull		С	2м/с
10.4.1997	+9		2 Dull		C3	6м/с	+5		Dull	Rain	C3	6м/с
11.4.1997	+4		Mostly Clouds		C3	17м/с	+1		Dull		C3	17м/с
12.4.1997	-	0 747	7 Dull		С	7м/с	-1		Dull	Snow	С	7м/с
13.4.1997	+3	756	Mostly Clouds		C3	5м/с	0	756	Sun		C3	5м/с
14.4.1997	+4	745	Dull	Rain	Ю3	6м/с	+5	742	Dull	Rain	ЮЗ	6м/с
15.4.1997	+2	747	Mostly Clouds		C3	8м/с	-1	750	Sun/Clouds		C3	8м/с
16.4.1997	+1	750	Dull	Rain	С	4м/с	+1	750	Dull		С	4м/с
17.4.1997	+11	746	6 Mostly Clouds		С	8м/с	+5	746	Dull		С	8м/с
18.4.1997	+7	744	Dull		C3	5м/с	+3	743	Dull	Rain	C3	5м/с
19.4.1997	+5	747	Mostly Clouds		C3	6м/с	0	749	Sun		C3	6м/с
20.4.1997	+2	750	Mostly Clouds		СВ	3м/с	0	751	Mostly Clouds		СВ	3м/с
21.4.1997	+5	750	Mostly Clouds		С	2м/с	-1	750	Sun/Clouds		С	2м/с
22.4.1997	+7		Mostly Clouds		С	2м/с	+1		Sun		С	2м/с
23.4.1997	+7	752	2 Mostly Clouds		С3	8м/с	+2	754	Sun		С3	8м/с
24.4.1997	+14		Mostly Clouds		юз	9м/с	+9	750	Sun		юз	9м/с
25.4.1997	+9		Dull	Rain	3	4м/с	+5	751	Sun/Clouds		3	4м/с
26.4.1997	+11		3 Sun		СВ	Зм/с	+5		Sun		СВ	3м/с
27.4.1997	+15		Sun/Clouds		ЮВ	6м/с	+10		Sun		ЮВ	6м/с
28.4.1997	+18		Mostly Clouds		ю	7м/с	+13		Mostly Clouds		ю	7m/c
29.4.1997	+16		5 Dull		юз	5м/с	+10		Mostly Clouds		юз	5м/с
30.4.1997	+17) Sun/Clouds		3	5м/с	+11		Sun/Clouds		3	5м/с
1.5.1997	+13		Mostly Clouds		C3	7m/c	+10		Dull		C3	7м/с
2.5.1997	+17		Mostly Clouds		3	9m/c	+14		Mostly Clouds		3	9м/с
3.5.1997	+16		B Dull		C3	10m/c	+8		Mostly Clouds		C3	10m/c
4.5.1997	+16		Sun/Clouds		3	6m/c	+8		Sun		3	6м/с
5.5.1997	+24				юз	5m/c	+16				юз	5M/c
			Mostly Clouds						Mostly Clouds			
6.5.1997	+23		Mostly Clouds		Ю	6м/с	+17		Sun/Clouds		Ю	6м/с
7.5.1997	+12		Mostly Clouds		3	5м/с	+7		Sun/Clouds		3	5м/с
8.5.1997	+9		Dull	Rain	CB	7м/с	+12		Dull		CB	7м/с
9.5.1997	+13		Mostly Clouds		3	7м/с	+8		Sun		3	7м/с
10.5.1997	+17		2 Mostly Clouds		ЮЗ	5м/с	+10		Sun		ЮЗ	5м/с
11.5.1997	+22		Sun/Clouds		Ю	4м/с	+15		Sun		Ю	4м/с
12.5.1997	+24		Mostly Clouds		Ю	Зм/с	+17		Sun		Ю	Зм/с
13.5.1997	+19		Dull		С	4м/с	+14		Sun		С	4м/с
14.5.1997	+26		Sun/Clouds		Ю	3м/с	+19		Sun/Clouds		Ю	Зм/с
17.5.1997	+21	754	Sun		В	4m/c	+14	755	Sun/Clouds		В	4м/с
18.5.1997	+23		Sun/Clouds		С	2м/с	+16	748	Sun		С	2м/с
19.5.1997	+20	743	B Dull	Rain	С	2м/с	+15	744	Dull		С	2м/с
20.5.1997	+17	744	Mostly Clouds		С	2м/с	+12	744	Mostly Clouds		С	2м/с
21.5.1997	+15	743	B Dull	Rain	3	4м/с	+13	744	Sun		3	4м/с
22.5.1997	+15	745	Dull		С	4м/с	+10	748	Dull		С	4м/с
23.5.1997	+11	753	Mostly Clouds		С	2м/с	+8	753	Dull		С	2м/с
24.5.1997	+12	756	Mostly Clouds		С	2м/с	+8	757	Mostly Clouds		С	2м/с
25.5.1997	+14	761	Mostly Clouds		СЗ	5м/с	+7	761	Sun		СЗ	5м/с
26.5.1997	+15	756	Mostly Clouds		С	2м/с	+9	753	Mostly Clouds		С	2м/с
27.5.1997	+11	749	Dull		С	2м/с	+8	748	Dull	Rain	С	2м/с
28.5.1997	+10	752	2 Dull	Storm	СЗ	2м/с	+9	754	Dull		СЗ	2м/с
29.5.1997	+14	753	B Dull		3	5м/с	+8	752	Dull	Rain	3	5м/с
30.5.1997	+9		2 Dull		СЗ	6м/с	+7		Dull	Rain	СЗ	6м/с
31.5.1997	+13		3 Dull	Storm	3	4м/с	+9		Dull		3	4м/с
1.6.1997	+12		3 Dull	Rain	СЗ	3м/с	+10		Dull	Rain	СЗ	3м/с
2.6.1997	+15		Mostly Clouds		С	4m/c	+11		Mostly Clouds		С	4м/с
3.6.1997	+13		Dull		3	5m/c	+10		Mostly Clouds		3	5м/с
4.6.1997	+18		Mostly Clouds		3	3m/c	+13		Sun		3	Зм/с
5.6.1997	+21		Mostly Clouds		No wind	No wind	+14		Sun/Clouds		No wind	No wind
6.6.1997	+22		Sun/Clouds		B	4M/C	+17		Sun		B	4m/c
7.6.1997	+22		Sun/Clouds		В	5m/c	+16		Sun/Clouds		В	5m/c
8.6.1997	+17		B Dull	Storm	С	2m/c	+15		Sun/Clouds		С	2M/C
9.6.1997	+17		5 Dull	Storm	3	3M/c	+15		Mostly Clouds		3	2M/C 3M/C
	+17			Gioini	C3		+16				C3	
10.6.1997			Mostly Clouds			5m/c			Sun			5M/c
11.6.1997	+21		Sun/Clouds		С	4m/c	+16		Sun		С	4m/c
12.6.1997	+24		Sun/Clouds		C	2м/с	+17		Sun/Clouds		C	2м/с
13.6.1997	+25		Mostly Clouds		Ю3	3м/с	+18		Mostly Clouds		Ю3	3м/с
14.6.1997	+16		Dull	Storm	3	5м/с	+16		Mostly Clouds	Storm	3	5м/с
15.6.1997	+19	748	Dull	Storm	3	7м/с	+15	750	Mostly Clouds		3	7м/с

Scrapers performance

The best results in terms of program speed are observed in Scrapy-based

solutions - this is due to the fact that the framework works asynchronously,

simultaneously processing several links.

The second fastest solution showed BeautifulSoup - all links are processed

sequentially, but very quickly individually - only one http request is required to

analyze data for one month.

The worst solution for such a task is performed by Selenium. The low speed

of work is due to the fact that the solution completely emulates a browser,

which means it does a lot of unnecessary work: it loads additional css, js files;

execute javascript code; sends additional http requests (for example, to the

yadro.ru site, which is connected to gismeteo.ru)

Total running time for the whole dataset without page limits (using VPN):

Beautiful Soup: 1307 seconds

Scrapy: 33.5 seconds

Selenium: 6054 seconds

Data analysis

* for the whole period from 1997 to 2023

