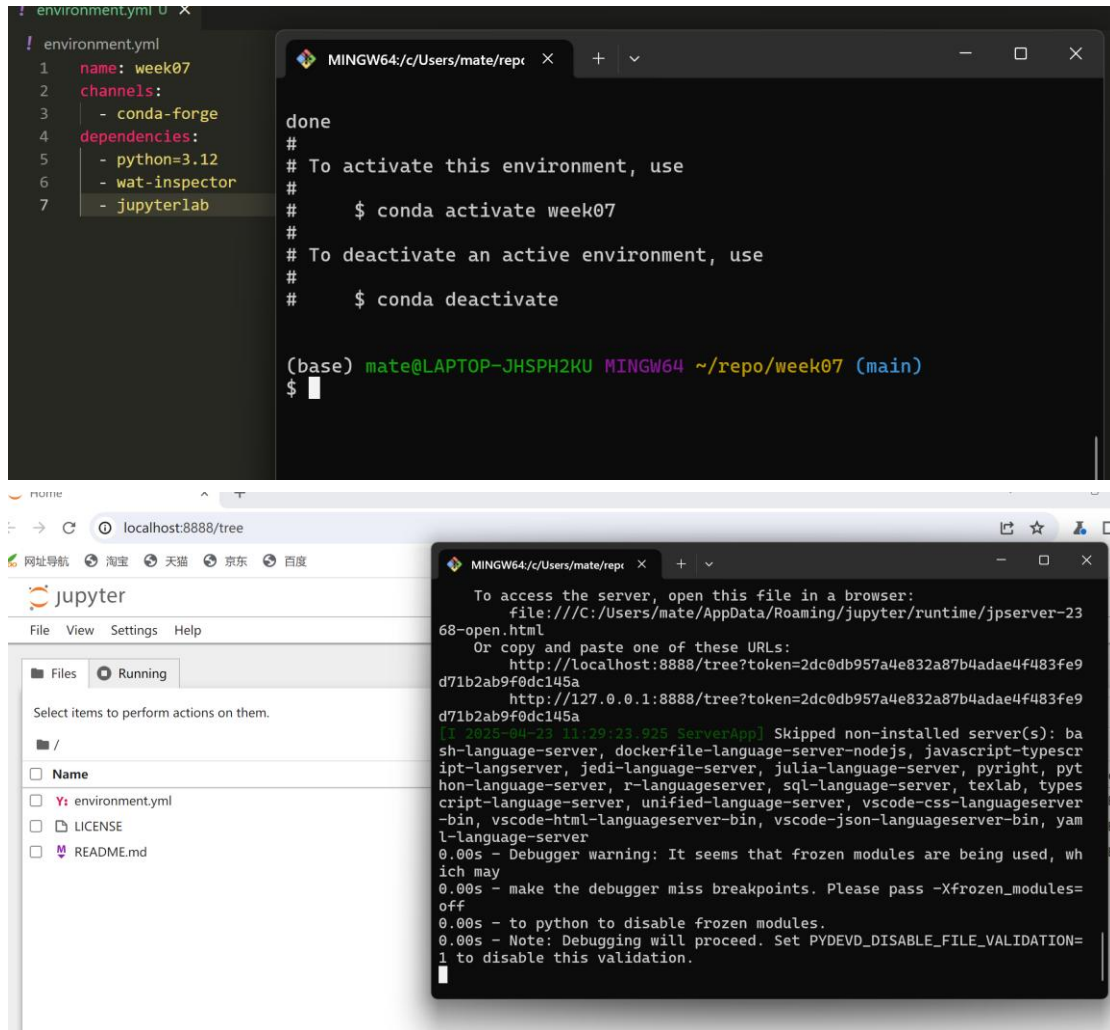


第七周学习报告-数据可视化与交互



浏览器前端 终端后端



```
ter trial-jupyterlab
View Run Kernel
Traceback (most recent call last):
  File "D:\anaconda\Lib\site-packages\tornado\web.py", line 1790, in _execute
    result = await result
    ^^^^^^^^^^^^^^^^^
  File "D:\anaconda\Lib\site-packages\panel\io\jupyter_server_extension.py", line 246, in get
    nb = json.load(f)
    ^^^^^^^^^^^^^^^^^
  File "D:\anaconda\Lib\json\__init__.py", line 293, in load
    return loads(fp.read(),
    ^^^^^^^^^^^^^^^^^
UnicodeDecodeError: 'gbk' codec can't decode byte 0xa2 in position 8133: illegal multibyte sequence
[E 2025-04-24 16:07:10.095 ServerApp] {
  "Host": "localhost:8888",
  "Accept": "text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3;q=0.7",
  "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/135.0.0.0 Safari/537.36 Edg/135.0.0.0"
}
[E 2025-04-24 16:07:10.096 ServerApp] 500 GET /panel-preview/render/trial-jupyterlab.ipynb (fb38e8a30480473ab7cefa563084did7@:1) 167.46ms referer=None
[I 2025-04-24 16:11:09.066 ServerApp] Saving file at /trial-jupyterlab.ipynb
[I 2025-04-24 16:13:53.186 ServerApp] Saving file at /trial-jupyterlab.ipynb
[I 2025-04-24 16:13:59.678 ServerApp] Interrupted...
[IPKernelApp] WARNING | Parent appears to have exited, shutting down.
(week07)
mate@LAPTOP-JHSPH2KU MINGW64 ~/repo/week07 (main)
$
```

复现。

按 `Ctrl+D` 结束前面的 IPython 进程，再

服务器请求 [IPO新股列表](#) 数据，并保存在

```
import tushare as ts

pro = ts.pro_api()
df = pro.new_share()
df.to_parquet("new_share.parquet")
```

其中请求数据函数返回的对象 `df` 是 `pandas`

中的 `DataFrame` 数据按照 [Parquet](#) 格式 (

(serialize) 为字节串 (bytes) 保存到磁盘 (c

询问豆包 (或 DeepSeek 等任何大模型)，并

`new_share` 接口只需要 120 积分，如果你

将数据保存为 `stock_basic.parquet` 文件

以在终端运行下面的命令，从我们开源的

```
2      301636.SZ      301636      泽润新能      ...      0.45      0.000      0.0
0
3      001400.SZ      001400      江顺科技      ...      1.50      5.604      0.0
1
4      301560.SZ      301560      众捷汽车      ...      0.70      5.016      0.0
2
...
1995      600989.SH      730989      宝丰能源      ...      22.00      81.550      0.2
5
1996      300778.SZ      300778      新城市      ...      2.00      5.466      0.02
1997      002953.SZ      002953      日丰股份      ...      1.70      4.526      0.0
3
1998      603697.SH      732697      有友食品      ...      3.10      6.257      0.0
5
1999      300772.SZ      300772      运达股份      ...      2.80      4.792      0.0
4
[2000 rows x 12 columns]

In [7]: df = pro.new_share()

In [8]: type(df)
Out[8]: pandas.core.frame.DataFrame

In [9]:
```

```
! environment.yml
1 name: week07
2 channels:
3   - conda-forge
4 dependencies:
5   - python=3.12
6   - wat-inspect
7   - jupyterlab
8   - pyarrow
9   - pip
10  - pip:
11    - tushare

IPython: Crepo/week07
In [4]: df = pro.new_share()
In [5]: df
Out[5]:
```

	ts_code	sub_code	name	...	limit_amount	funds	ballot
0	301595.SZ	301595	太力科技	...	0.65	0.000	0.0
0							
1	688755.SH	787755	汉邦科技	...	0.50	0.000	0.0
0							
2	301636.SZ	301636	泽润新能	...	0.45	0.000	0.0
3	001400.SZ	001400	江顺科技	...	1.50	5.604	0.0
1							
4	301560.SZ	301560	众捷汽车	...	0.70	5.016	0.0
2							
...							
1995	600989.SH	730989	宝丰能源	...	22.00	81.550	0.2
5							
1996	300778.SZ	300778	新城市	...	2.00	5.466	0.02
1997	002953.SZ	002953	日丰股份	...	1.70	4.526	0.0
3							
1998	603697.SH	732697	有友食品	...	3.10	6.257	0.0
5							
1999	300772.SZ	300772	运达股份	...	2.80	4.792	0.0
4							

```
[2000 rows x 12 columns]

In [6]:
```

```
IPython: C:\repo\week07
[5415 rows x 10 columns]
In [10]: df = pro.stock_basic()
In [11]: df.columns
AttributeError                                Traceback (most recent call last)
Cell In[11], line 1
----> 1 df.columns

File D:\anaconda\envs\week07\Lib\site-packages\pandas\core\generic.py:6299, in NDFrame.__getattr__(self, name)
    6292 if (
    6293     name not in self._internal_names_set
    6294     and name not in self._metadata
    6295     and name not in self._accessors
    6296     and self._info_axis._can_hold_identifiers_and_holds_name(name)
    6297 ):
    6298     return self[name]
-> 6299 return object.__getattr__(self, name)

AttributeError: 'DataFrame' object has no attribute 'columns'

In [12]: df.columns
Out[12]:
Index(['ts_code', 'symbol', 'name', 'area', 'industry', 'cnsPELL', 'market',
      'list_date', 'act_name', 'act_ent_type'],
      dtype='object')

In [13]:
```

```
In [18]: df.to_parquet("stock_basic.par")
In [19]: ls -lh
驱动器 C 中的卷是 Windows
卷的序列号是 CC81-A8AF
C:\Users\mate\repo\week07 的目录
找不到文件
In [20]: ls
驱动器 C 中的卷是 Windows
卷的序列号是 CC81-A8AF
C:\Users\mate\repo\week07 的目录
2025/04/24 17:05 <DIR> .
2025/04/23 11:21 <DIR> ..
2025/04/23 11:21 1,307 .gitignore
2025/04/23 11:39 <DIR> .ipynb_checkpoints
2025/04/24 16:49 154 environment.yml
2025/04/23 11:21 18,805 LICENSE
2025/04/24 16:53 119,952 new_share.parquet
2025/04/23 11:21 2,239 README.md
2025/04/24 17:05 429,130 stock_basic.par
2025/04/24 16:13 24,919 trial-jupyterlab.ipynb
7 个文件 596,506 字节
3 个目录 136,825,438,208 可用字节

In [21]:
```

1. **Parquet 格式特点**

列式存储

* Parquet 是一种列式存储文件格式。在存储数据时，它会将同一列的数据存储在一起。例如，在一个包含用户信息（如用户 ID、姓名、年龄、地址等字段）的表格中，Parquet 会把所有用户 ID 存放在一个区域，姓名存放在另一个区域，依此类推。

* 这种存储方式的优势在于，当进行查询操作时，如果只需要查询某一列或某几列的数据，例如统计所有用户的年龄分布情况，系统可以直接读取存储年龄的那一列数据，而无需读取整个行的数据。这大大减少了数据读取量，提高了查询效率，尤其是在处理大数据场景下，数据量庞大时，这种优势更为明显。

*****高效的压缩和编码****

* Parquet 支持多种压缩和编码算法，如 Snappy、Gzip、LZO 等。通过压缩，可以有效减少数据的存储空间。例如，对于文本类型的列数据，使用合适的压缩算法可以使其存储空间大幅减少。

* 同时，编码算法（如字典编码、游程编码等）可以进一步优化数据存储。字典编码是将列中重复出现的值用一个较短的代码来代替，从而节省存储空间。例如，一个列中有大量重复的字符串值“apple”，在字典编码后，可以用一个代码（如 1）来代替所有“apple”，在解码时再将其还原。

*****数据类型丰富****

* 它支持多种数据类型，包括基本数据类型（如整数、浮点数、字符串等）、复杂数据类型（如嵌套数据结构、列表、映射等）。这使得 Parquet 格式能够很好地存储和表达各种复杂的数据模型。

*****与大数据生态系统高度集成****

* Parquet 是为大数据处理而设计的，与 Hadoop、Spark、Presto 等大数据处理工具和框架无缝集成。在 Spark 中，可以很方便地将数据读取为 Parquet 格式的数据集，并进行各种分布式计算操作。例如，在 SparkSQL 中，可以通过简单的 SQL 语句查询 Parquet 格式的数据，而无需进行繁琐的数据转换。

2. **CSV 格式特点**

*****简单的文本格式****

* CSV（Comma - Separated Values）格式是一种简单的文本文件格式。其文件内容是以行为单位，每行表示一条记录，行内的各个字段之间用逗号(,)分隔。例如，“1,John,25,New York”表示一条包含用户 ID、姓名、年龄和地址四个字段的记录。

*****易于阅读和编辑****

* 由于是纯文本格式，不需要特殊的软件就可以阅读和编辑。在简单的文本编辑器（如 Notepad、Notepad++ 等）或者电子表格软件（如 Excel）中都可以方便地打开和修改 CSV 文件。这使得它在数据交换和简单数据处理场景中非常方便。

*****存储效率低****

* 相比 Parquet 的列式存储，CSV 是行式存储。在存储大量数据时，尤其是包含很多列的数据时，存储效率较低。而且它没有像 Parquet 那样高效的压缩机制，存储空间占用较大。

*****数据类型支持有限****

* CSV 本身不直接存储数据类型信息。在读取 CSV 文件时，需要根据实际数据内容来推断数据类型，这可能会导致一些数据类型转换的错误或者不确定性。例如，一个字段可能既包含数字字符串（如“123”）又包含文本字符串（如“abc”），在读取时很难准确判断其数据类型。

3. **Parquet 格式适用领域**

*****大数据分析和处理****

* 在处理海量数据时，如数据仓库中的数据存储和查询，Parquet 的列式存储和高效压缩编码能够大大提升性能。例如，在电信行业，对用户的通话记录、流量使用记录等海量数据进行分析时，Parquet 格式可以快速地过滤和查询特定列的数据，如统计某个地区用户的流量使用高峰时段。

*****数据湖存储****

* 数据湖是一种存储大量原始数据的存储库，Parquet 格式是数据湖中常用的数据存储格式之一。它可以存储不同类型的数据源（如日志数据、传感器数据等）生成的结构化和半结构化数据，并且方便后续的数据分析和挖掘。例如，在物联网场景下，存储来自各种传感器的大量时间序列数据，Parquet 格式可以有效地组织和存储这些数据，便于后续的时序分析。

*****分布式计算框架中的数据存储****

* 在 Spark 等分布式计算框架中，Parquet 格式的数据可以被高效地分布式存储和处理。因为 Parquet 的文件格式支持数据的分片存储，能够很好地适应分布式文件系统（如 HDFS）的存储架构，使得数据可以在多个计算节点上并行处理。

4. **CSV 格式适用领域**

*****数据交换和共享****

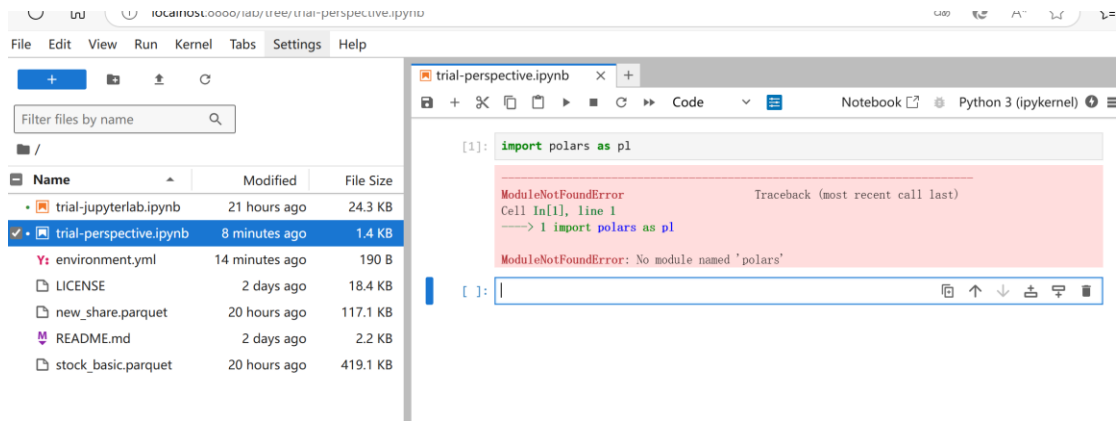
* CSV 格式常用于不同系统之间进行数据交换。例如，在一个企业的不同部门之间，或者企业与合作伙伴之间，当需要共享一些简单的数据（如客户名单、产品列表等）时，CSV 格式是一种简单易用的选择。因为大多数数据库系统和数据分析软件都支持导入和导出 CSV 文件。

*****简单的数据处理和分析****

* 对于一些小型的、简单的数据处理任务，如在 Excel 中进行数据的排序、筛选、简单的统计分析等，CSV 格式可以直接在电子表格软件中打开和使用。例如，一个小商店可以使用 CSV 格式记录每天的销售数据，并在 Excel 中进行数据的初步分析，如计算每日销售额、统计销售量等。

*****数据导入导出的中间格式****

* 在数据从一个系统迁移到另一个系统，或者从一个软件工具导入到另一个软件工具时，CSV 格式常作为中间格式。例如，将一个旧的数据库系统中的数据导出为 CSV 文件，然后将其导入到新的数据库系统或者数据分析软件中。



遇到问题，没有安装 polars，但是我已经在 environment.yml 里面输入了 polars

```

webencodings              0.5.1
websocket-client           1.8.0
Werkzeug                  3.0.3
whatthepatch              1.0.2
wheel                      0.44.0
widgetsnbextension        3.6.6
win-inet-pton             1.1.0
wrapt                      1.14.1
xarray                    2023.6.0
xlwings                   0.32.1
xyzservices                2022.9.0
yapf                       0.40.2
yarl                       1.11.0
zict                       3.0.0
zipp                       3.17.0
zope.interface             5.4.0
zstandard                  0.23.0

(base) mate@LAPTOP-JHSPH2KU MINGW64 ~/repo/week07 (main)
$ conda activate week07
(week07)
mate@LAPTOP-JHSPH2KU MINGW64 ~/repo/week07 (main)
$ conda install -c conda-forge polars
Channels:
 - conda-forge
 - https://repo.anaconda.com/pkgs/main
 - https://repo.anaconda.com/pkgs/r
 - https://repo.anaconda.com/pkgs/msys2
Platform: win-64
Collecting package metadata (repodata.json): -

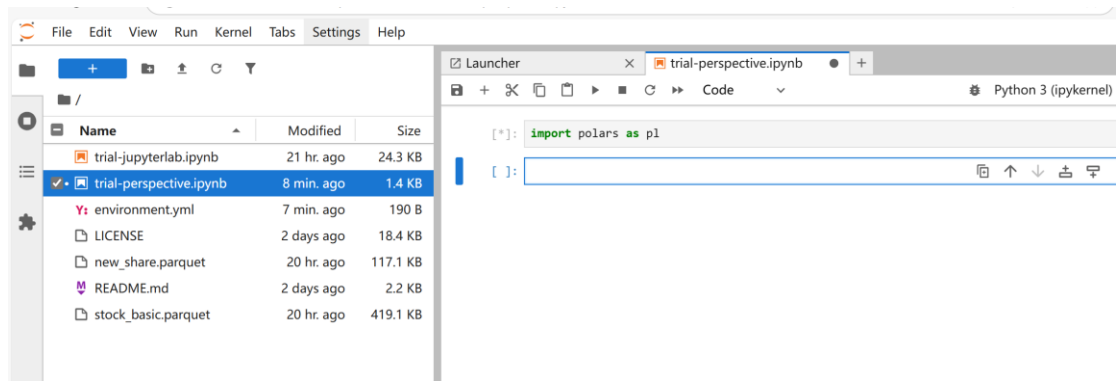
```

检查是否安装 polars 发现未安装 重新安装 polars

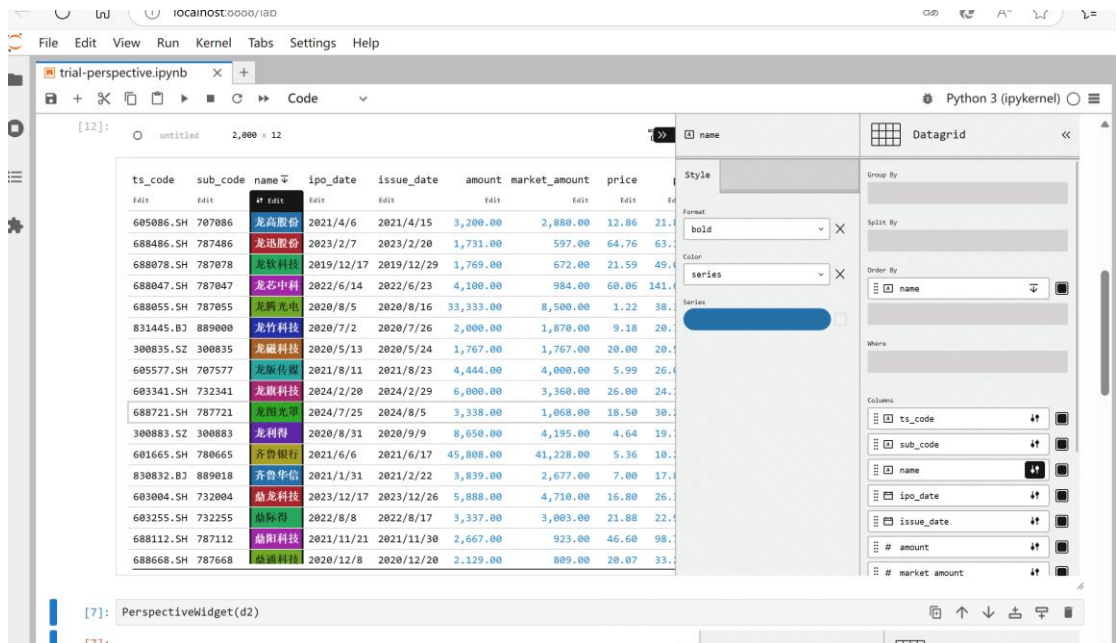
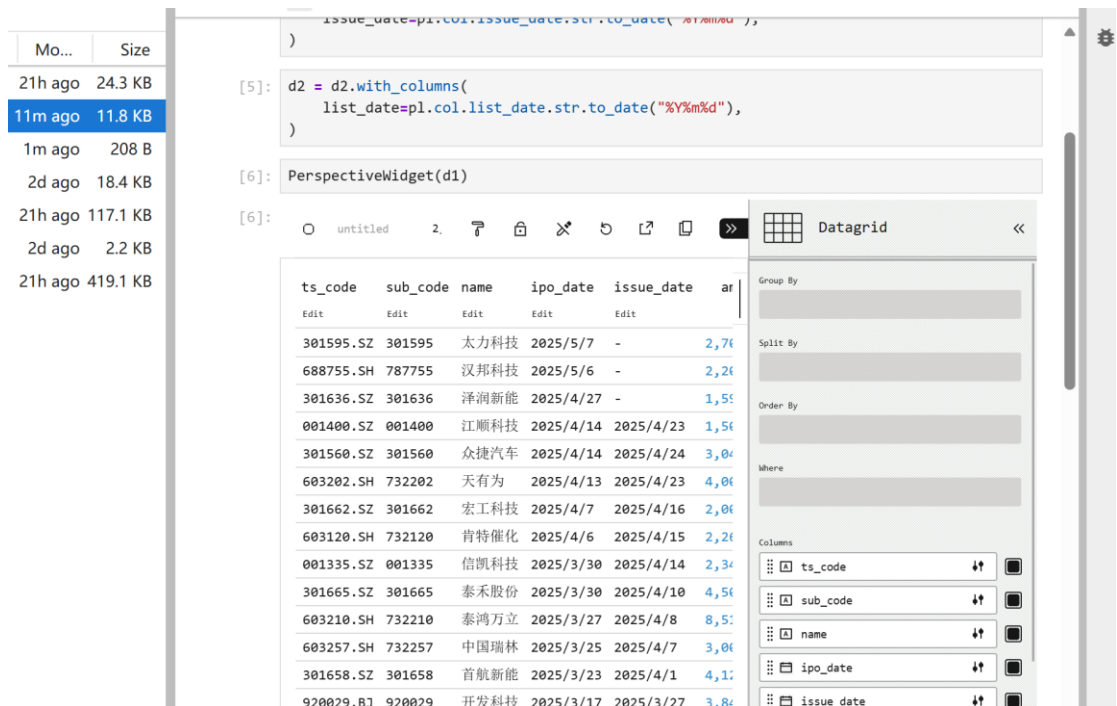
```

nbconvert                  7.16.6
nbformat                   5.10.4
nest_asyncio               1.6.0
notebook_shim              0.2.4
numpy                      2.2.5
overrides                  7.7.0
packaging                  25.0
pandas                     2.2.3
pandocfilters              1.5.0
parso                      0.8.4
perspective-python         3.0.1
pickleshare                0.7.5
pip                        25.0.1
pkgutil_resolve_name       1.3.10
platformdirs               4.3.7
polars                     1.27.1
prometheus_client          0.21.1
prompt_toolkit             3.0.51
psutil                     7.0.0
pure_eval                  0.2.3
pyarrow                    19.0.1
pycparser                  2.22
Pygments                   2.19.1
PySocks                    1.7.1
python-dateutil            2.9.0.post0
python-json-logger         2.0.7
pytz                       2025.2
pywin32                    307
pywinpty                   2.0.15
PyYAML                     6.0.2

```

成功



Python 3 (ipykernel)

Code

market amount

[7]: PerspectiveWidget(d2)

untitled 5 (5,415) x 1 (17)

ts_code

5,415

3,183

1,381

265

586

industry2

Attributes

Expression

1 "industry"

DELETE COLUMN

RESET

SAVE

Datagrid

Group By

market

Split By

Order By

Where

Columns

count ts_code

All Columns

NEW COLUMN

industry2

act_ent_type

[7]: PerspectiveWidget(d2)

untitled 9 (5,415) x 3 (17)

ts_code industry industry2

5,415

265

265

2,282

1,696

586

2,868

1,487

1,381

电气设备

汽车配件

汽车配件

半导体

汽车配件

半导体

元器件

元器件

软件服务

Style

Foreground

Bar

Color Range

Max Value

107

Background

Disabled

Style

Decimal

Minimum Integer Digits

1

Rounding Increment

Auto

Notation

Standard

Datagrid

Group By

exchange

market

Split By

Order By

Where

Columns

count ts_code

distinct industry

All Columns

NEW COLUMN

act_ent_type

[13]: PerspectiveWidget(d2)

untitled 4 (5,415) x 8 (17)

主板 创业板 北交所 科创板

ts_code list_date ts_code list_date ts_code list_date ts_code list_date

1,253 1,253 1,062 1,062 265 265 586 586

-

BSE

SSE

SZSE

834 834

419 419

1,062 1,062

-

-

-

-

Datagrid

Group By

exchange

market

Split By

Order By

ts_code

Where

list_date >= 2012/08/01

Columns

act_ent_type

act_name

File Edit View Run Kernel Tabs Settings Help

trial-perspective.ipynb Python 3 (ipykernel)

[11]: polars.dataframe.frame.DataFrame

[12]: PerspectiveWidget(d1)

[12]:

ts_code	sub_code	name	ipo_date	issue_date	amount	market_amount	上网发行比例	price_cent	price	pe	limit_amount	funds	ballot	configure
600905.SH	730905	三峡能源	2021/5/30	2021/6/9	857,100	625,451	72.97%	¥265	2.65	28.87	257.10	227.13	1.28	
601728.SH	780728	中国电信	2021/8/8	2021/8/19	1,057,477	380,402	35.97%	¥453	4.53	20.18	311.80	479.04	0.96	
601816.SH	780816	京沪高铁	2020/1/5	2020/1/15	628,563	234,379	37.29%	¥488	4.88	23.39	94.20	0.00	0.79	
601916.SH	780916	浙商银行	2019/11/13	2019/11/25	255,000	186,081	72.97%	¥494	4.94	9.39	76.50	0.00	0.69	
003816.SZ	003816	中国广核	2019/8/11	2019/8/25	504,986	184,252	36.49%	¥249	2.49	14.60	75.70	111.67	0.60	
601658.SH	780658	邮储银行	2019/11/27	2019/12/9	594,799	181,026	30.43%	¥550	5.50	9.58	170.60	0.00	1.26	
600968.SH	730968	海油发展	2019/6/13	2019/6/25	186,510	167,859	90.00%	¥204	2.04	22.93	55.90	38.05	0.61	
600938.SH	730938	中国海油	2022/4/11	2022/4/20	299,000	99,257	33.20%	¥1,080	10.80	24.07	78.00	322.92	0.43	
601077.SH	780077	渝农商行	2019/10/13	2019/10/28	135,700	99,024	72.97%	¥736	7.36	9.26	40.70	0.00	0.35	
601825.SH	780825	沪农商行	2021/8/3	2021/8/18	96,444	86,800	90.00%	¥890	8.90	10.65	28.90	85.84	0.21	
003035.SZ	003035	南网能源	2021/1/7	2021/1/18	75,758	68,182	90.00%	¥140	1.40	22.88	22.70	10.61	0.18	
600989.SH	730989	宝丰能源	2019/4/29	2019/5/15	73,336	66,002	90.00%	¥1,112	11.12	22.07	22.00	81.55	0.25	
600918.SH	730918	中泰证券	2020/5/19	2020/6/2	69,686	62,718	90.00%	¥438	4.38	48.08	20.90	0.00	0.24	
600925.SH	730925	苏能股份	2023/3/16	2023/3/28	68,889	62,000	90.00%	¥618	6.18	22.49	20.60	42.57	0.24	
688538.SH	787538	和辉光电	2021/5/17	2021/5/27	308,366	56,310	18.26%	¥265	2.65	0.00	77.75	81.72	0.30	
601778.SH	780778	晶科科技	2020/5/5	2020/5/18	59,459	53,513	90.00%	¥437	4.37	16.58	17.80	26.82	0.22	
001286.SZ	001286	陕西能源	2023/3/28	2023/4/9	75,000	52,500	70.00%	¥960	9.60	14.86	22.50	72.00	0.19	
001227.SZ	001227	兰州银行	2022/1/4	2022/1/16	56,957	51,261	90.00%	¥357	3.57	22.97	17.05	20.33	0.14	
601096.SH	780096	盛华源	2023/12/12	2023/12/21	66,879	46,815	70.00%	¥170	1.70	31.54	20.05	11.37	0.25	

File Edit View Run Kernel Tabs Settings Help

trial-perspective.ipynb Python 3 (ipykernel)

```
[17]: print(config)

{"version": "3.6.0", "plugin": "Datagrid", "plugin_config": {"columns": {}, "edit_mode": "READ_ONLY", "scroll_lock": false, "columns_config": {"name": {"format": "bold", "string_color_mode": "foreground"}, "price": {"number_bg_mode": "gradient", "bg_gradient": 557.8, "market_amount": {"number_format": {"minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "price_cent": {"number_format": {"style": "currency", "currency": "CNY", "minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "amount": {"number_bg_mode": "gradient", "bg_gradient": 1057477.0, "number_format": {"minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "上网发行比例": {"number_fg_mode": "disabled", "number_bg_mode": "gradient", "bg_gradient": 1.0, "number_format": {"style": "percent", "minimumFractionDigits": 2.0, "maximumFractionDigits": 2.0}}, "settings": true, "theme": "Pro Light", "title": null, "group_by": [], "split_by": [], "sort": [{"market_amount": "desc"}], "filter": [], "expressions": {"上网发行比例": "market_amount" / "amount", "price_cent": "price" * 100}, "columns": ["ts_code", "sub_code", "name", "ipo_date", "issue_date", "amount", "market_amount", "上网发行比例", "price_cent", "price", "pe", "limit_amount", "funds", "ballot"], "aggregates": {}}}}

[18]: from pathlib import Path

[20]: config2 = Path("C:/Users/mate/Downloads/untitled.config.json").read_text(encoding="utf8")

[22]: print(config2)

{"version": "3.6.0", "plugin": "Datagrid", "plugin_config": {"columns": {}, "edit_mode": "READ_ONLY", "scroll_lock": false, "columns_config": {"name": {"format": "bold", "string_color_mode": "foreground"}, "price": {"number_bg_mode": "gradient", "bg_gradient": 557.8, "market_amount": {"number_format": {"minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "price_cent": {"number_format": {"style": "currency", "currency": "CNY", "minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "amount": {"number_bg_mode": "gradient", "bg_gradient": 1057477.0, "number_format": {"minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "上网发行比例": {"number_fg_mode": "disabled", "number_bg_mode": "gradient", "bg_gradient": 1.0, "number_format": {"style": "percent", "minimumFractionDigits": 2.0, "maximumFractionDigits": 2.0}}, "settings": true, "theme": "Pro Light", "title": null, "group_by": [], "split_by": [], "sort": [{"market_amount": "desc"}], "filter": [], "expressions": {"上网发行比例": "market_amount" / "amount", "price_cent": "price" * 100}, "columns": ["ts_code", "sub_code", "name", "ipo_date", "issue_date", "amount", "market_amount", "上网发行比例", "price_cent", "price", "pe", "limit_amount", "funds", "ballot"], "aggregates": {}}}}
```

codebeautify.org/jsonviewer

1 {"version": "3.6.0", "plugin": "Datagrid", "plugin_config": {"columns": {}, "edit_mode": "READ_ONLY", "scroll_lock": false, "columns_config": {"name": {"format": "bold", "string_color_mode": "foreground"}, "price": {"number_bg_mode": "gradient", "bg_gradient": 557.8, "market_amount": {"number_format": {"minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "price_cent": {"number_format": {"style": "currency", "currency": "CNY", "minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "amount": {"number_bg_mode": "gradient", "bg_gradient": 1057477.0, "number_format": {"minimumFractionDigits": 0.0, "maximumFractionDigits": 2.0}, "上网发行比例": {"number_fg_mode": "disabled", "number_bg_mode": "gradient", "bg_gradient": 1.0, "number_format": {"style": "percent", "minimumFractionDigits": 2.0, "maximumFractionDigits": 2.0}}, "settings": true, "theme": "Pro Light", "title": null, "group_by": [], "split_by": [], "sort": [{"market_amount": "desc"}], "filter": [], "expressions": {"上网发行比例": "market_amount" / "amount", "price_cent": "price" * 100}, "columns": ["ts_code", "sub_code", "name", "ipo_date", "issue_date", "amount", "market_amount", "上网发行比例", "price_cent", "price", "pe", "limit_amount", "funds", "ballot"], "aggregates": {}}}}

Ln: 1 Col: 1165 size: 1.15 KB

File URL

☒ Auto Update ☐ Big Num

Tree Viewer

2 Tab Space

Beautify

Minify Validate

to XML to CSV

选择一个节点...

- object {14}
- version: 3.6.0
- plugin: Datagrid
- plugin_config {3}
- columns {0}
- (清空 object)
- edit_mode: READ_ONLY
- scroll_lock: false
- columns_config {6}
- name {2}
- format: bold
- string_color_mode: foreground
- price {2}
- number_bg_mode: gradient
- bg_gradient: 557.8
- market_amount {1}
- number_format {2}
- minimumFractionDigits: 0
- maximumFractionDigits: 2

