

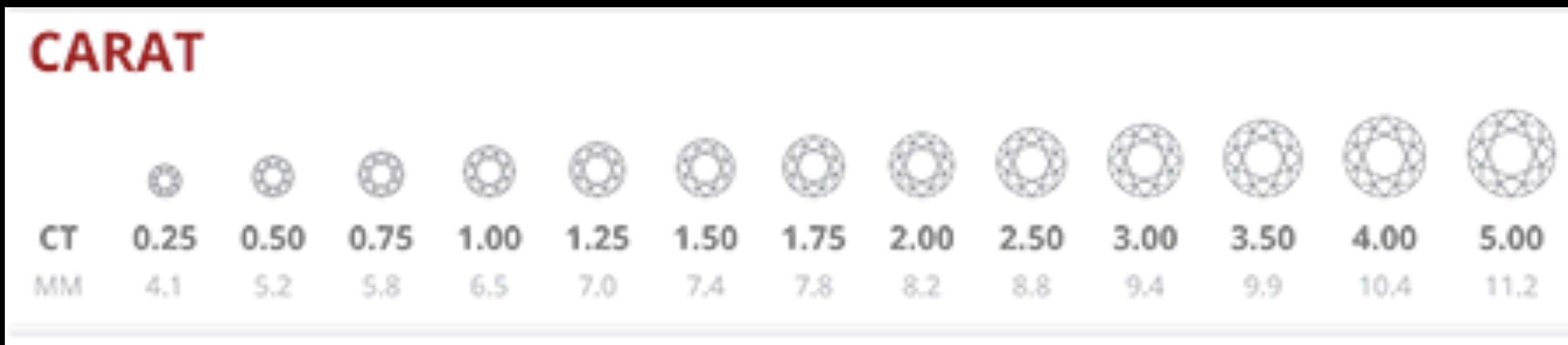
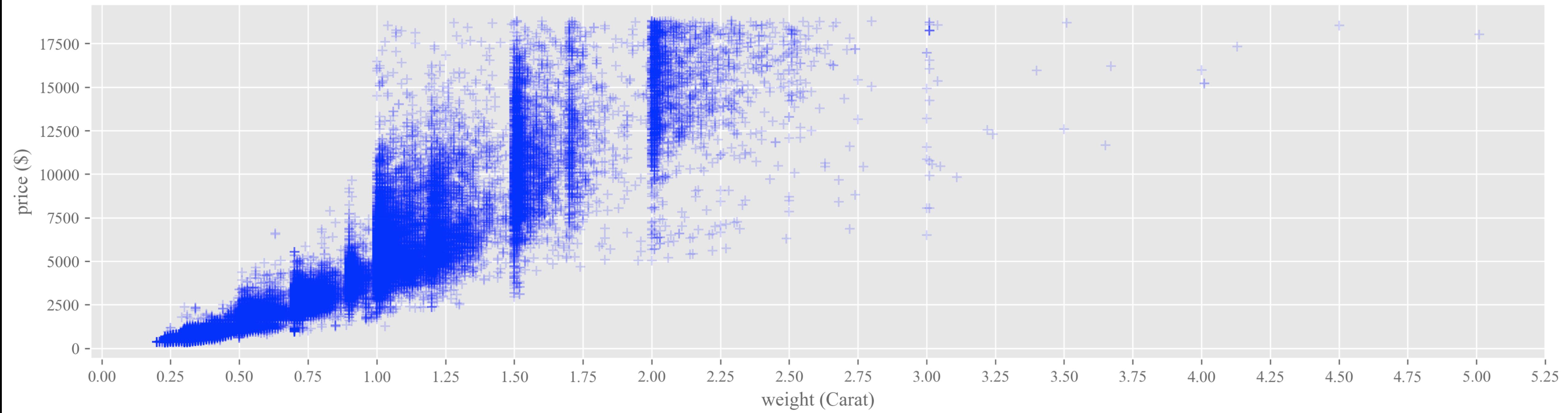
张
安
北
居
游
旅

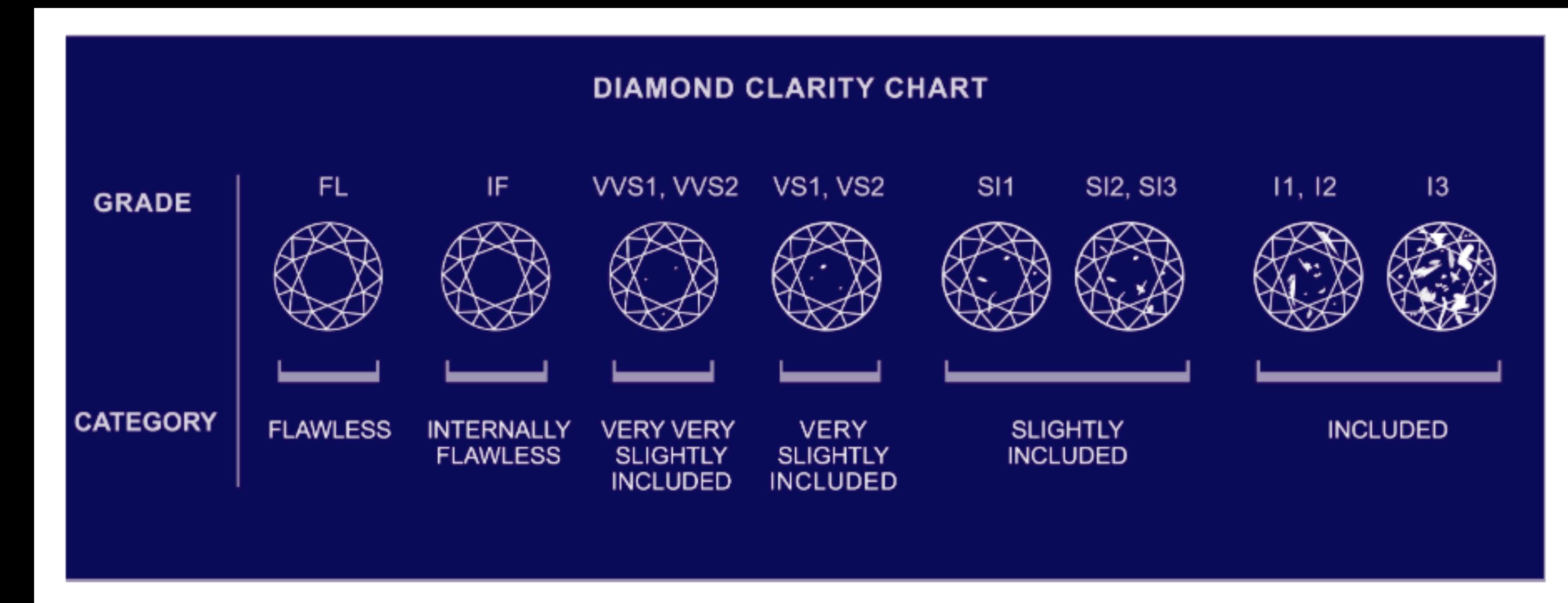
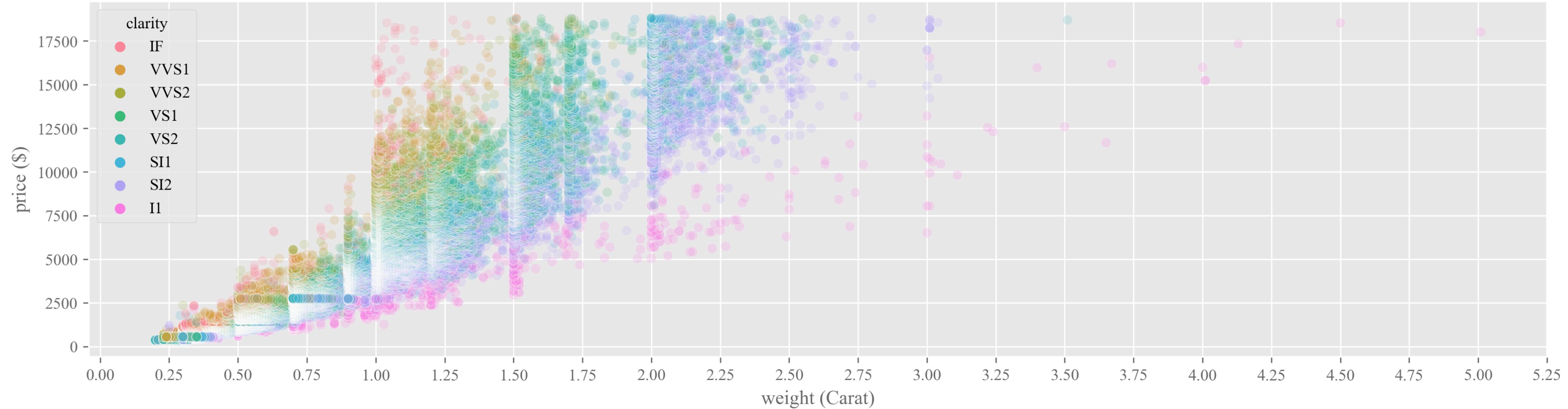


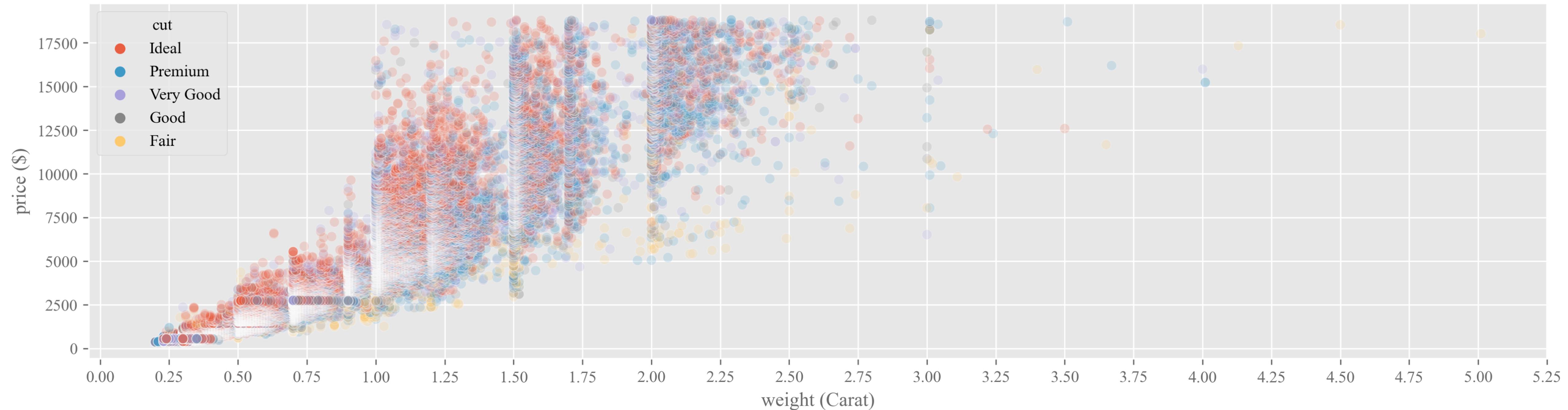
环境配置(environment)

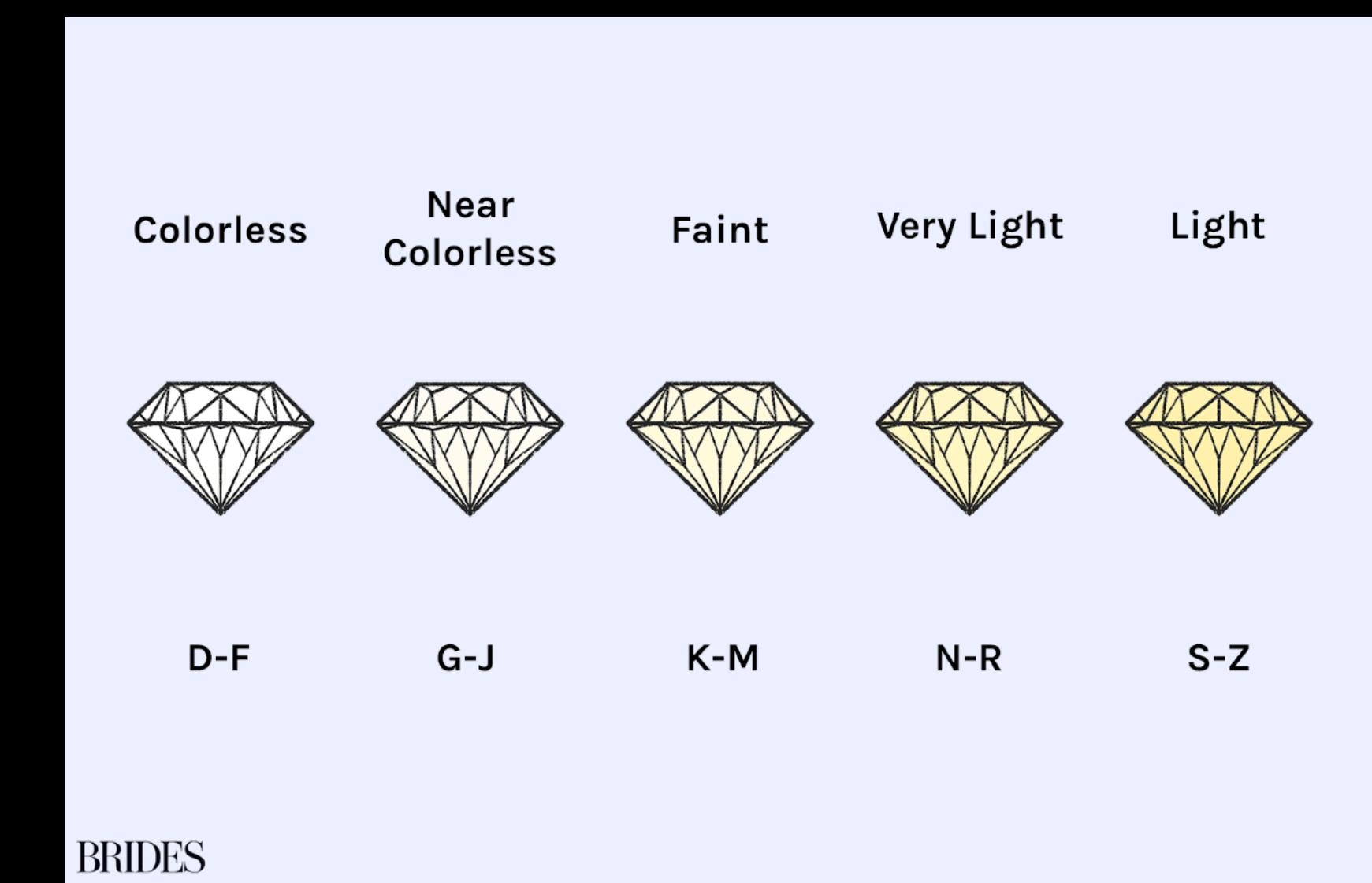
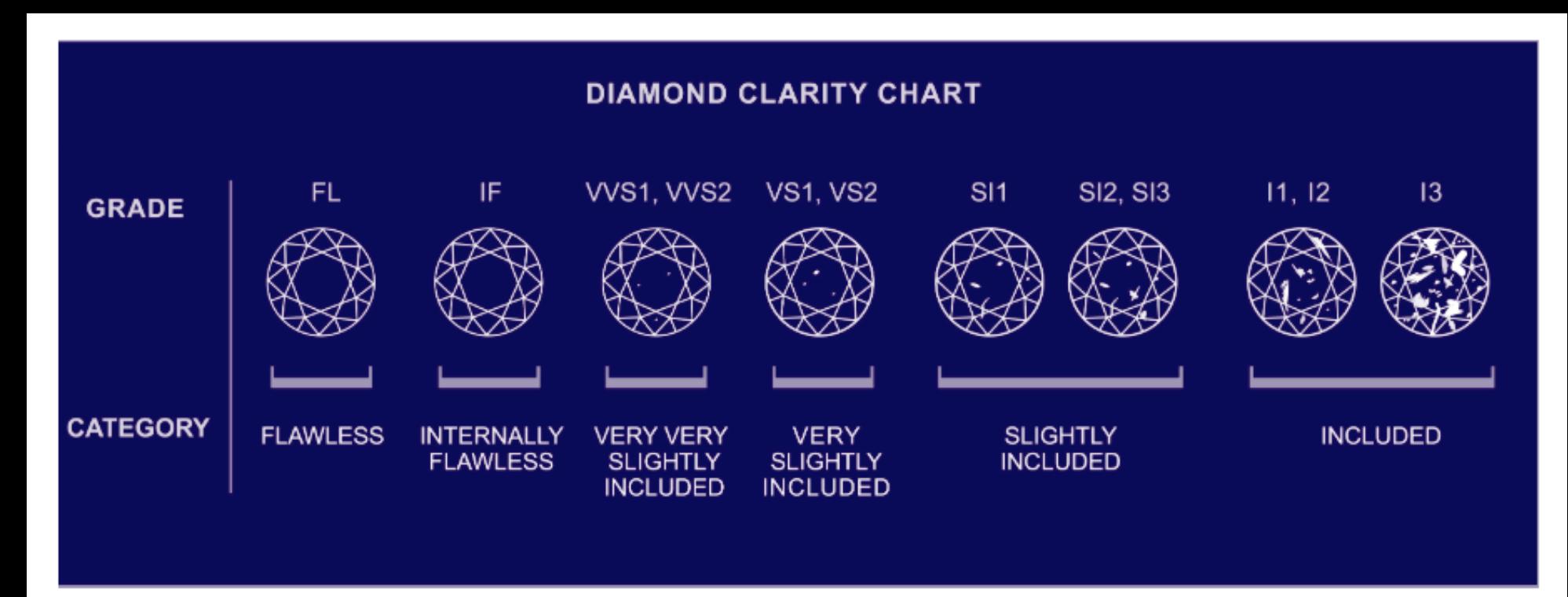
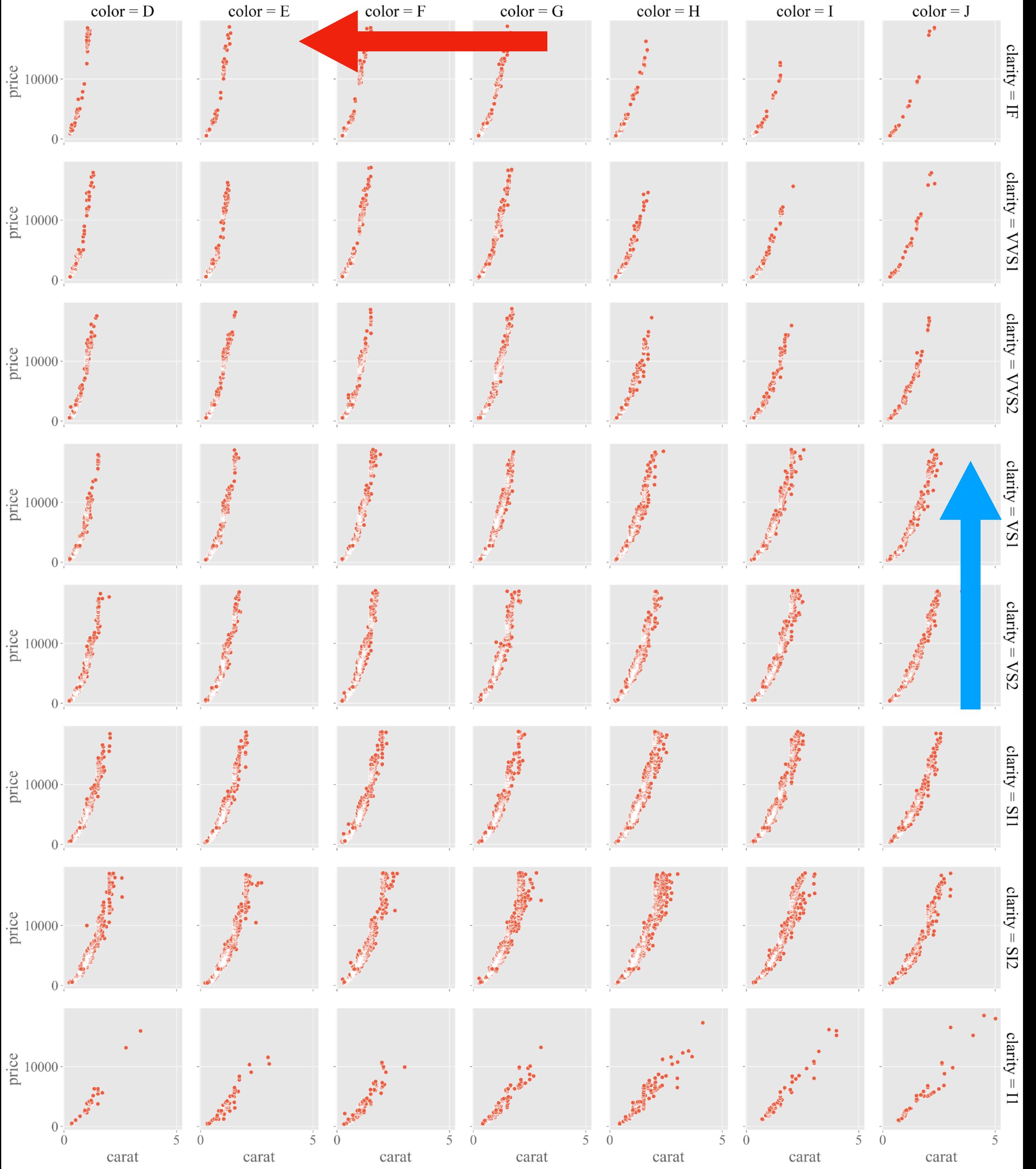
- 使用Anaconda/conda建立多个环境
 - ds: 处理基本的数据处理/可视化
 - r_env: R语言环境
 - dl_torch: pytorch
 - dl_tf: tensorflow
 - 专门的项目配置专门环境：
 - 如：Flask网页；Manim视频；
 - 仿真；强化学习，等等.
 - 策略：
 - 核心包使用conda安装
 - “边缘”包使用conda中的pip安装

```
(base) ~ % conda env list
# conda environments:
#
base          * /opt/anaconda3
dl             /opt/anaconda3/envs/dl
ds             /opt/anaconda3/envs/ds
manim          /opt/anaconda3/envs/manim
manim2         /opt/anaconda3/envs/manim2
r_env          /opt/anaconda3/envs/r_env
spinningup    /opt/anaconda3/envs/spinningup
```



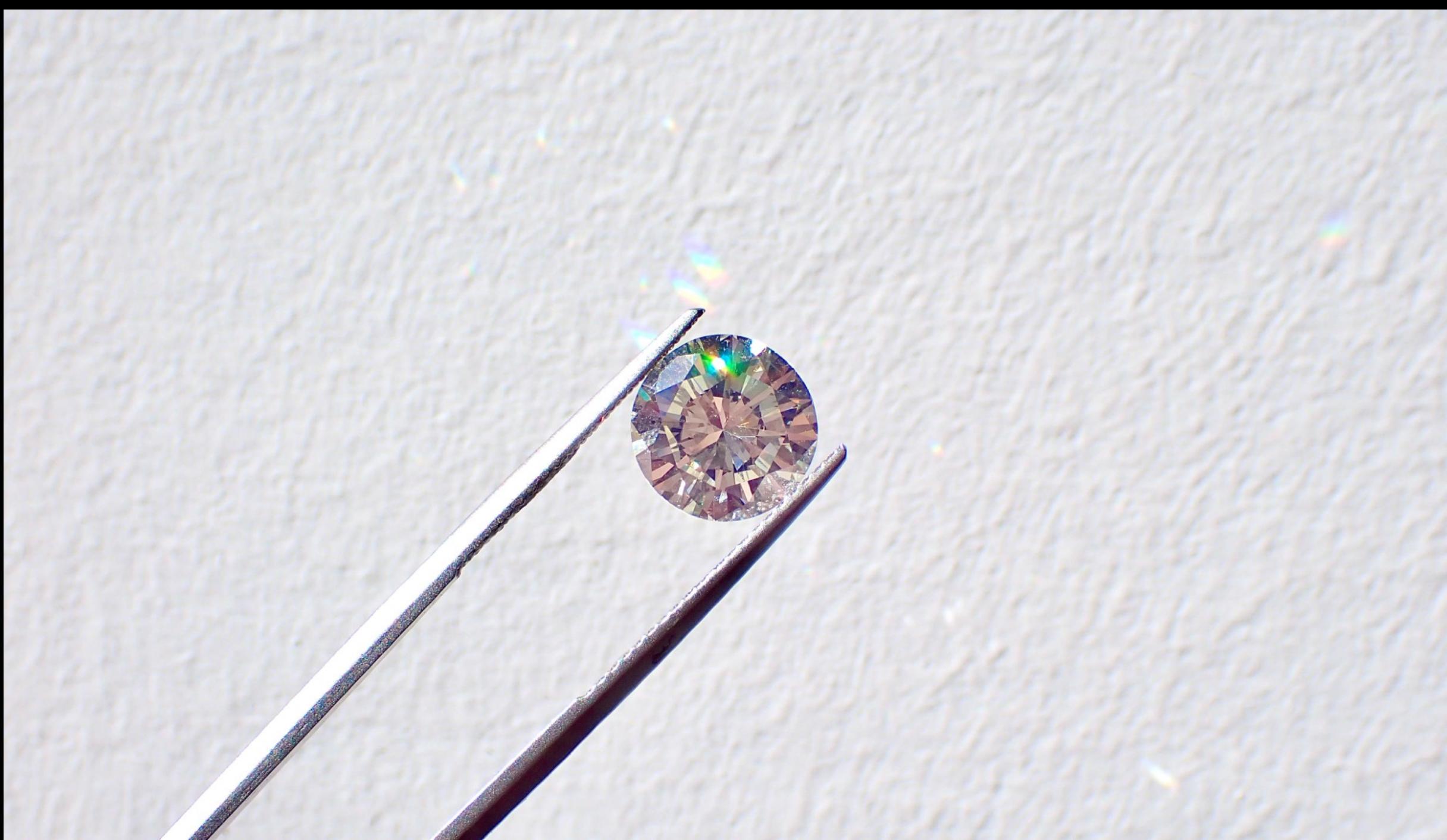






新建项目(new project)

- 最简单的初始文件夹结构
 - notebook: 存储jupyter代码(.ipynb)
 - scripts: 存储python代码 (.py)
 - data_input: 输入数据
 - data_output: 输出数据



Name
▼ diamonds
▼ data_input
diamonds.csv
> data_output
▼ notebook
EDA.ipynb
▼ scripts
preprocess.py

使用IDE开发(environment)

- Jupyter-Lab
 - 主要处理notebook中的代码
- PyCharm
 - Community version 社区版本 免费
 - 主要处理python文件中的代码
 - 匹配conda环境
- 策略：
 - 同时使用两个软件打开一个项目

开启Jupyter-lab

```
(base) ~ % cd /Documents/bilibili/diamonds/diamonds  
[base] diamonds % conda activate ds  
(ds) diamonds % jupyter-lab
```

The screenshot shows the Jupyter Lab interface. On the left is a file browser with a sidebar containing icons for File, Edit, View, Run, Kernel, Tabs, Settings, and Help. The file browser shows a directory structure under '/diamonds/' with files: data_input, data_output, notebook, and scripts. A search bar at the top of the file browser allows filtering by name. The main area contains two tabs: 'EDA.ipynb' (active) and 'preprocess.py'. The 'EDA.ipynb' tab shows a Python 3 (ipykernel) session with the following code and output:

```
[1]: import seaborn as sns
import pandas as pd
import os

[2]: print(sns.get_dataset_names())
['anagrams', 'anscombe', 'attention', 'brain_networks', 'car_crashes', 'diamonds', 'dots', 'dowjones', 'exercise', 'flights', 'fmri', 'geyser', 'glue', 'healthexp', 'iris', 'mpg', 'penguins', 'planets', 'seacie', 'taxis', 'tips', 'titanic']

[3]: df_diamonds = sns.load_dataset('diamonds')

[4]: df_diamonds.sample(10)
   carat      cut  color clarity depth  table  price     x     y     z
0  2.3224    Ideal     G     SI1   62.4   57.0  11230  7.33  7.41  4.60
1  0.7213     Fair     H     SI1   64.6   55.0   4191  6.29  6.25  4.05
2  0.45506   Very Good    D    VS2   62.6   58.0   1682  5.14  5.18  3.23
3  0.46695     Good     G   VVS2   64.1   56.0   1791  5.06  5.02  3.23
4  0.14506     Good     G    VS2   63.5   59.0   5860  6.23  6.33  3.99
5  0.41819    Ideal     G     SI2   59.1   56.0   1250  5.58  5.52  3.28
6  0.5189     Good     H    VS2   63.6   55.0   3775  6.14  6.09  3.89
7  0.27327   Very Good    E     SI2   62.4   61.0  17871  8.10  8.19  5.08
8  0.17826     Fair     H    VS2   65.6   56.0   7188  7.08  4.70  4.00
9  0.2691     Good     I    VS2   63.7   61.0   3246  6.10  6.06  3.87
```

The DataFrame output shows 10 rows of diamond data with columns: carat, cut, color, clarity, depth, table, price, x, y, z.

```
[5]: # store the dataset
df_diamonds.to_csv("../data_input/diamonds.csv")
```

The bottom status bar indicates 'Mode: Command' and 'Ln 1, Col 10' for the EDA.ipynb file.

PyCharm

New UI What's New Features Learn Pricing Download

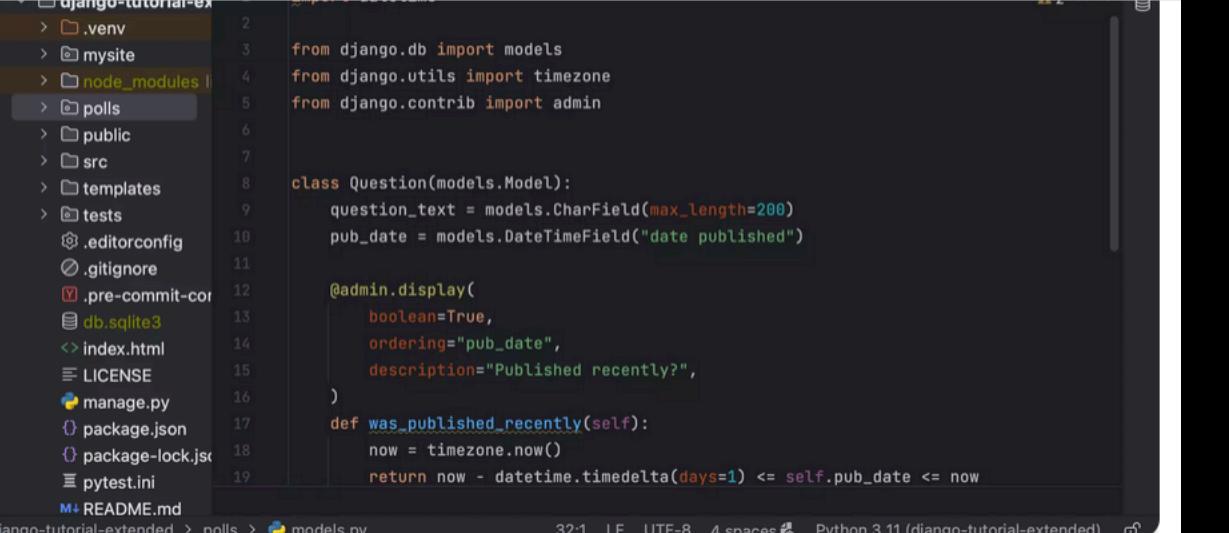
PyCharm Professional

The Python IDE for Professional Developers

Download .dmg ▾

Free 30-day trial

Select an installer for Intel or Apple Silicon



Version: 2023.2 System requirements Other versions
Build: 232.8660.197 Installation instructions Third-party software
25 July 2023

We value the vibrant Python community, and that's why we proudly offer the PyCharm Community Edition for free, as our open-source contribution to support the Python ecosystem.

PyCharm Community Edition

The IDE for Pure Python Development

Download .dmg ▾

Free, built on open source

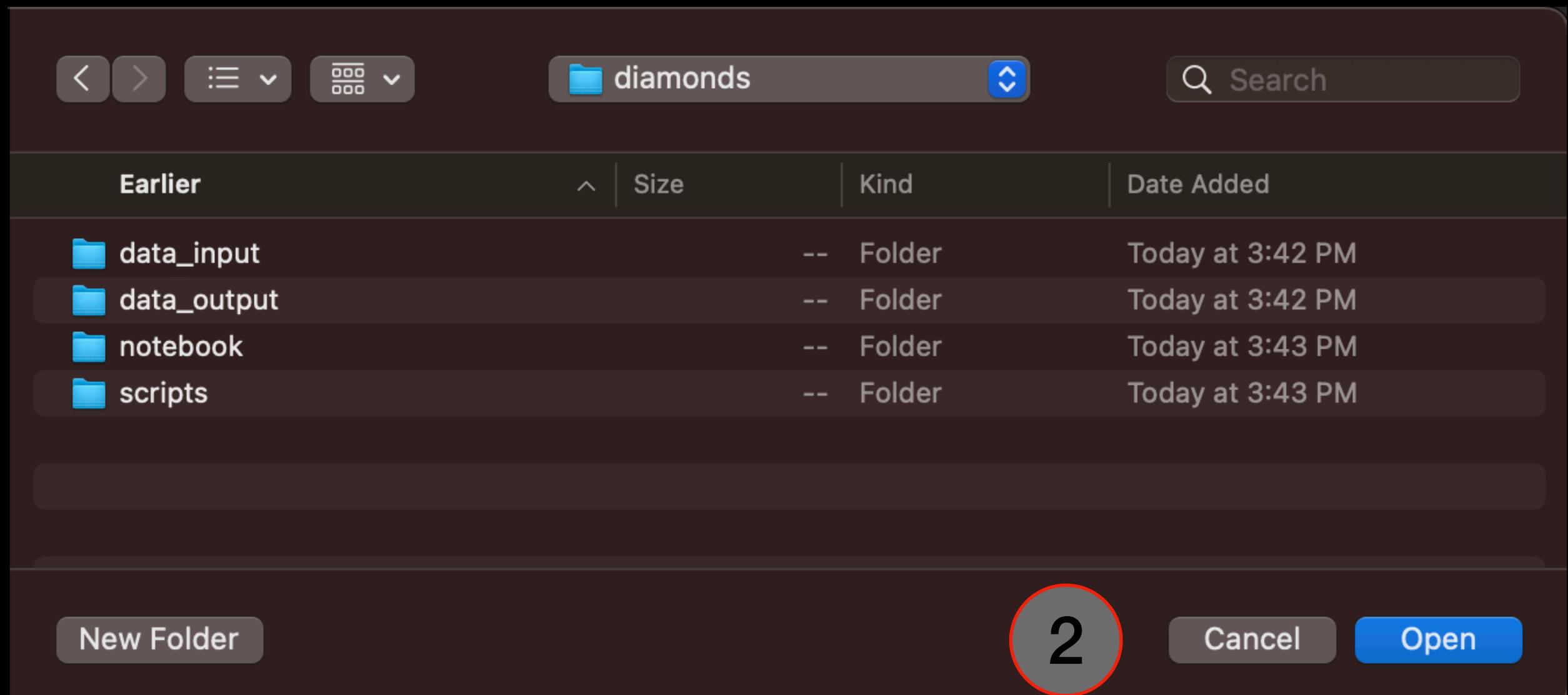
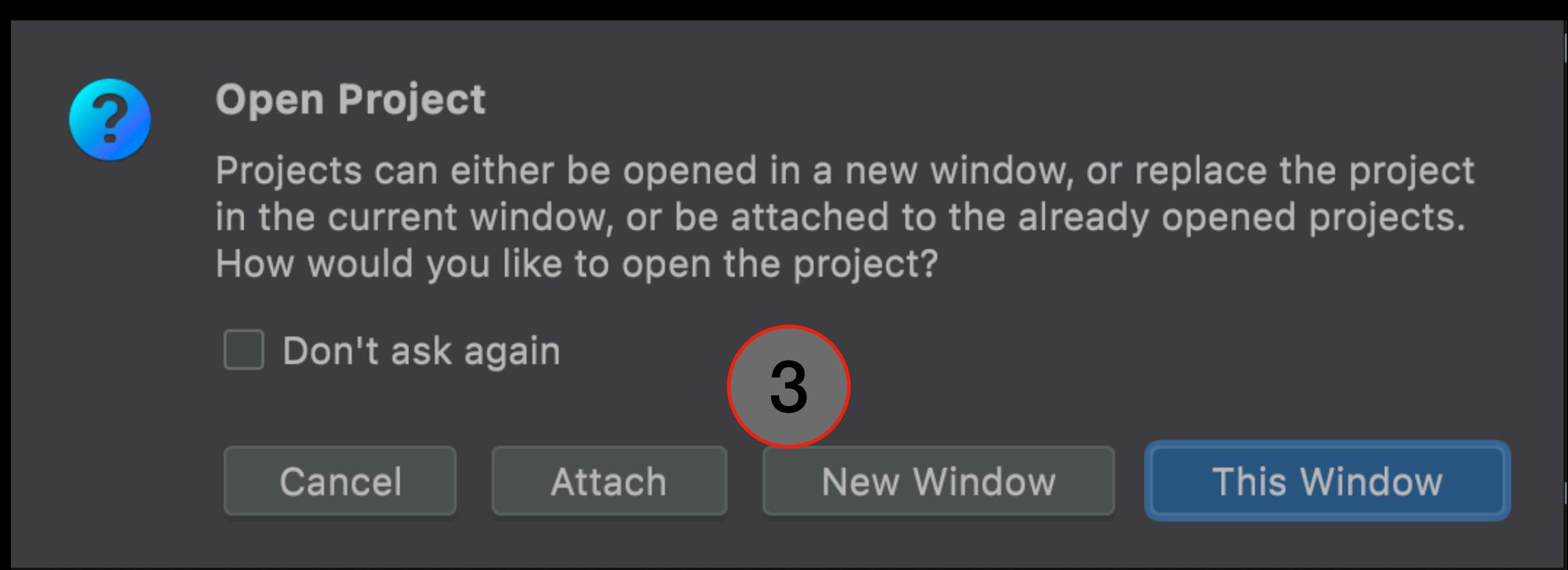
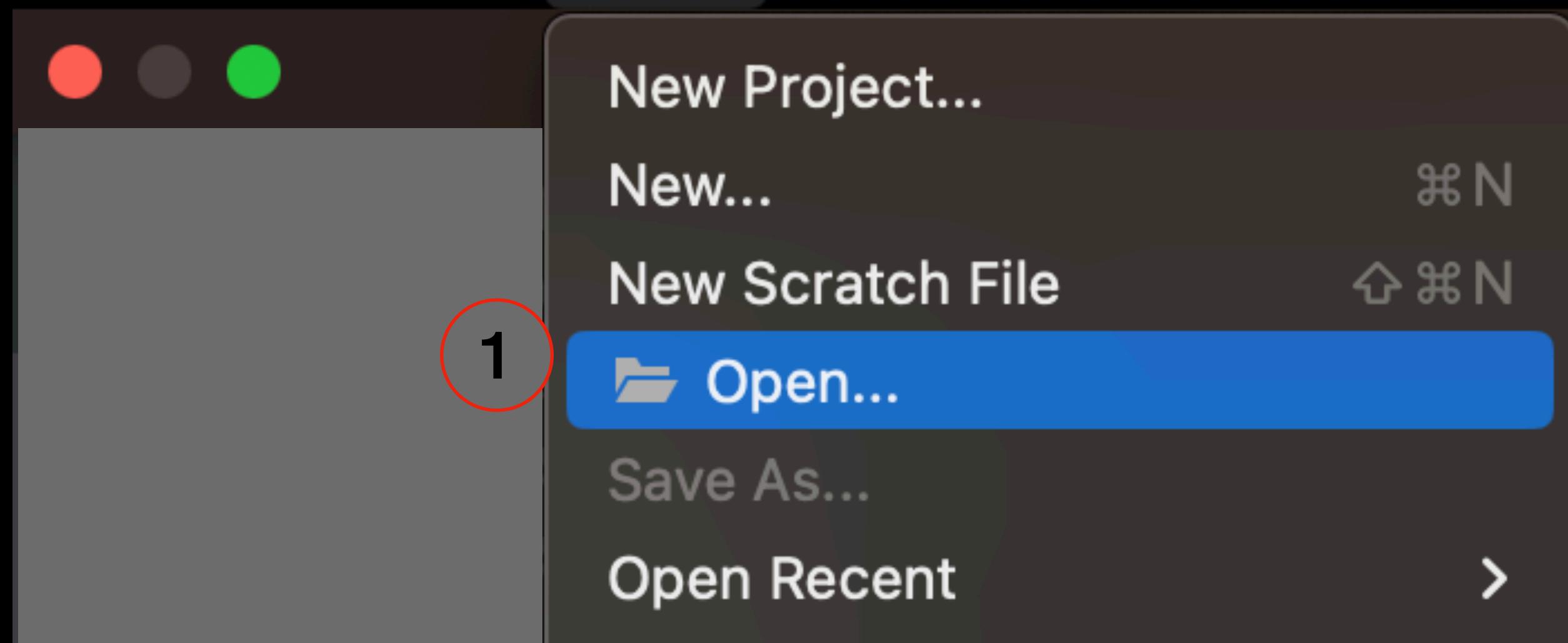
Select an installer for Intel or Apple Silicon

1





File Edit View Navigate C



This composite screenshot illustrates the process of changing the Python interpreter in PyCharm across three panels.

Panel 1: The PyCharm main menu is shown with the "Settings..." option highlighted (circled in red). The "File" menu is also visible.

Panel 2: The "Preferences" dialog is open, showing the "Project: diamonds > Python Interpreter" section. The "Python Interpreter" dropdown is highlighted (circled in red).

Panel 3: The "Python Interpreter" dropdown is expanded, showing a list of available interpreters. The "Python 3.8 (ds) /opt/anaconda3/envs/ds/bin/python" option is selected (highlighted with a blue border and circled in red).

The "Python Interpreter" dropdown list contains the following packages and their versions:

Package	Version	Latest version
abseil-cpp	20210324.2	
absl-py	1.4.0	
affine	2.3.1	
aiosignal	1.3.1	
alabaster	0.7.12	
altair	4.2.0	
altgraph	0.17.2	
anyio	3.5.0	
appdirs	1.4.4	
applaunchservices	0.2.1	
appnope	0.1.3	
argh	0.26.2	
argon2-cffi	21.3.0	
argon2-cffi-bindings	21.2.0	
arrow	1.2.2	
arrow-cpp	6.0.1	
astroid	2.11.2	
astunparse	1.6.3	
async_generator	1.10	
atomicwrites	1.4.0	
attrs	21.4.0	
autopep8	1.6.0	
aws-c-auth	0.6.8	

联合快速开发(fast development)

✓ 准备步骤：

✓ 同时使用两个软件打开一个项目

✓ 新建python/notebook文件

✓ 导入数据

• 开发步骤(以Explorative Data Analysis为例)：

✓ 在Notebook中进行各种实验 (如可视化)

• 将可视化代码写成函数

• 将函数迁移至python文件中

• 在notebook中添加路径，导入包和函数

• 将notebook调整至动态加载模式

联合快速卡发(fast development)

EDA.ipynb

Markdown Python 3 (ipykernel)

```
[1]: import seaborn as sns  
import pandas as pd  
import os  
  
[1]: 1  
[1]: import sys  
sys.path.append("../scripts")  
import preprocess as pre  
  
[3]: # use this method so everytime you changed your Python file (.py)  
# the files are re-imported again at the time you  
# running the next cell in Jupyter Notebook (.ipynb)  
%load_ext autoreload  
%autoreload 2  
  
[1]: 4  
[1]: pre.hello_world()  
Hello World!
```

preprocess.py

```
1 import os  
2  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns  
5 import numpy as np  
6 import pandas as pd  
7  
7 2  
7 def hello_world():  
8     print("Hello World!")  
9  
10  
11  
12 #-----below are methods loading the diamond dataset-----  
13 def load_data_from_seaborn():  
14     # the diamond dataset is contained in the seaborn package  
15     df = sns.load_dataset("diamonds")  
16     return df  
17  
18 def load_data_from_local(pth="../data_input/diamonds.csv"):  
19     df = pd.read_csv(pth, index_col=0)  
20     return df
```

3 保存

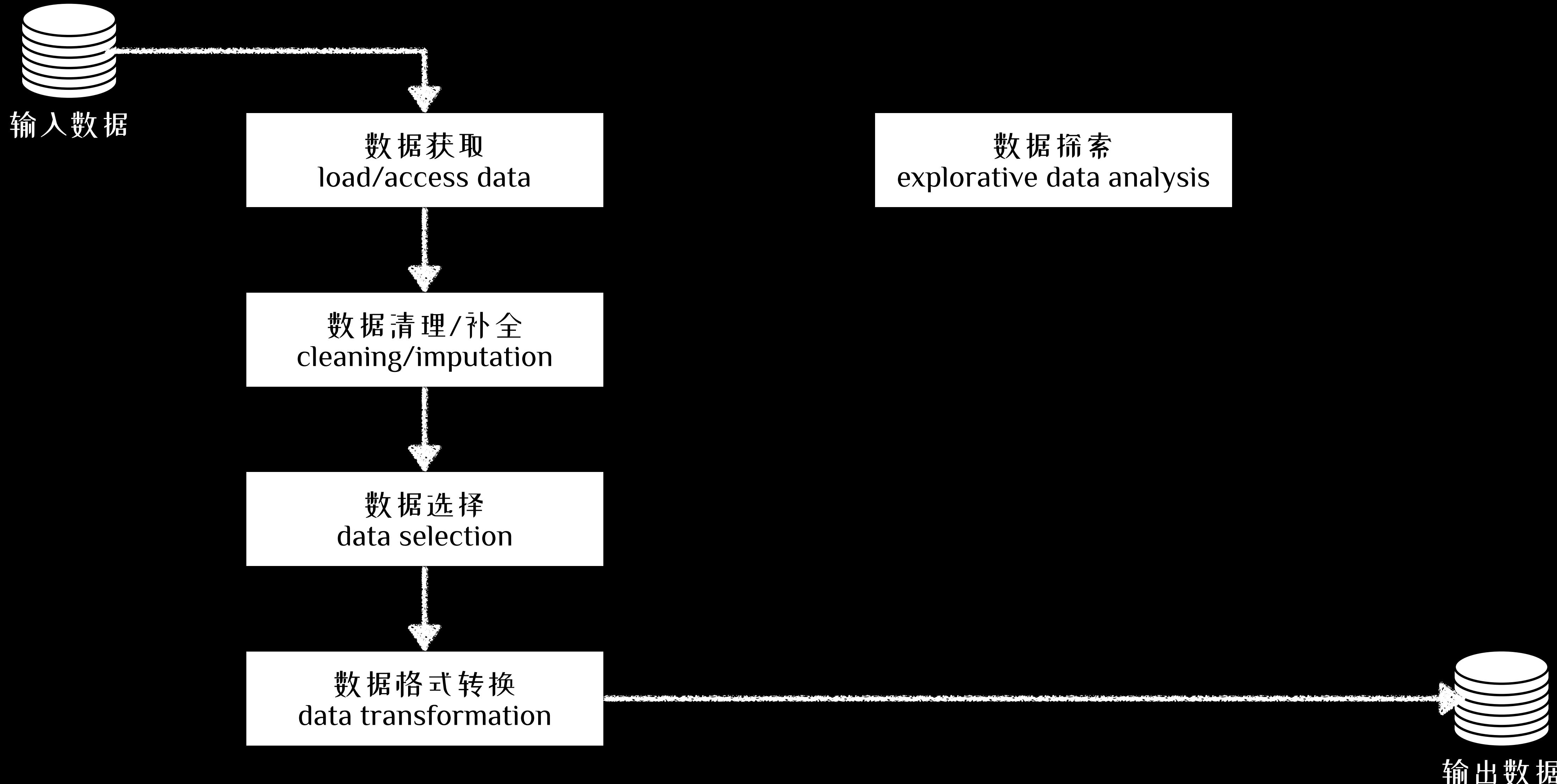


数据科学项目的具体步骤

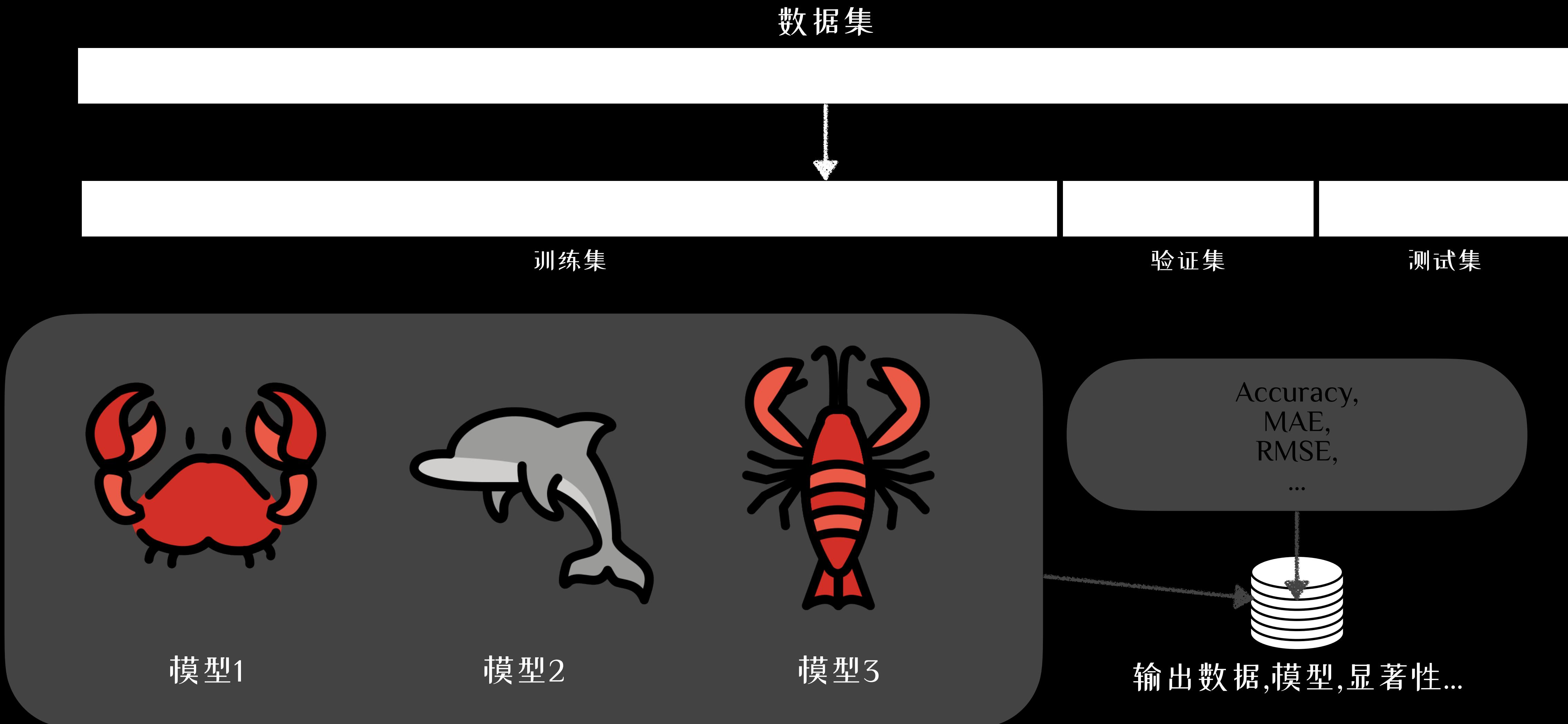
- 数据管道 (data pipeline)
 - 管道1: 数据预处理 (pre-process)
 - 管道2: 建立/测试模型 (机器学习流程) (analysis)
 - 管道3: 评估并汇报结果 (post-analysis)



管道1/步骤1: 数据预处理

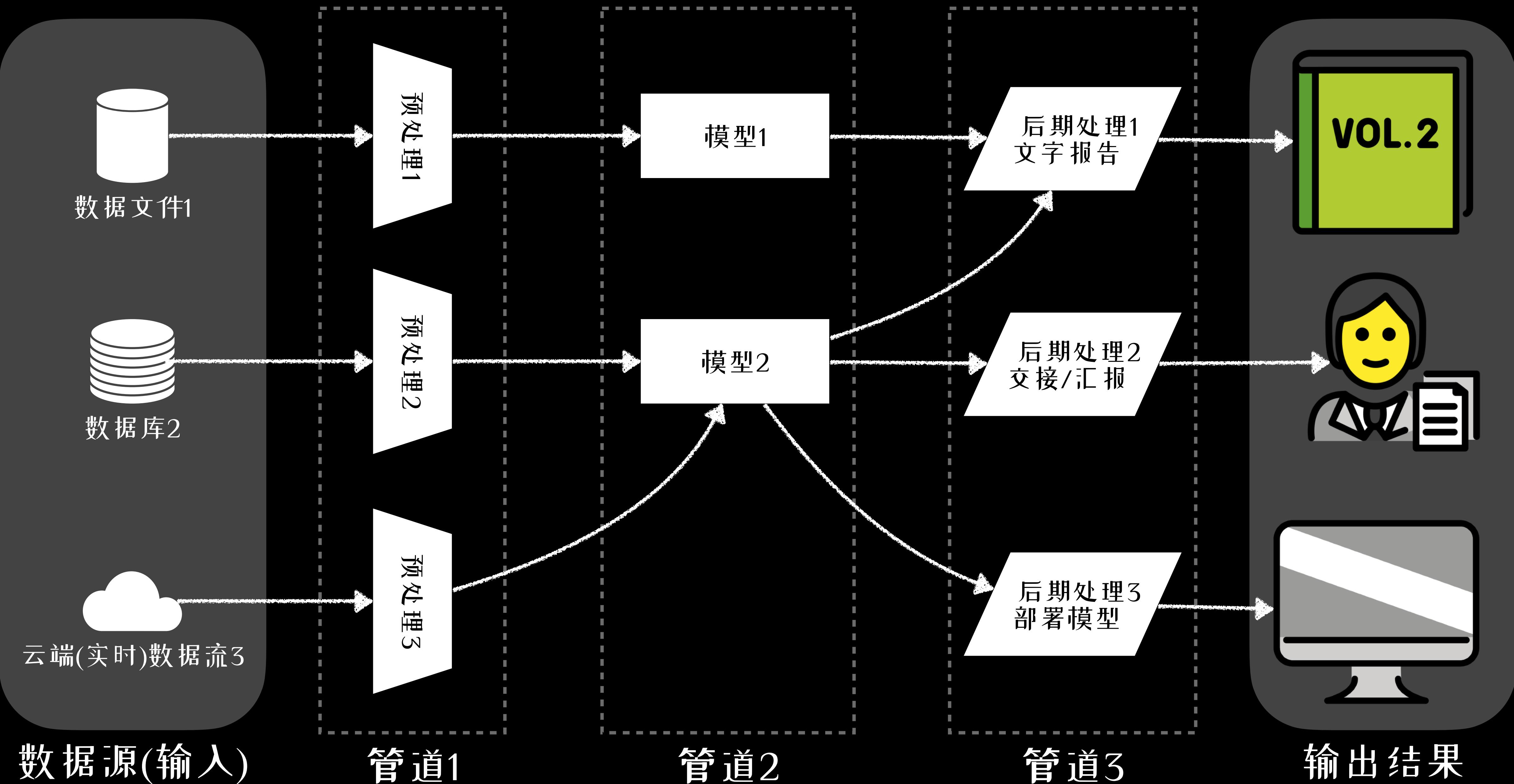


管道2/步骤2: 建立/测试模型(例:机器学习流程)



管道3/步骤3:评估并汇报/转交成果

- 目标1：将结果呈现给他人
 - 上级，同事，甲方等等
 - 画图：需要高质量
 - 文字：写报告，论文，PPT，网页...
- 目标2：提交最终/部署模型



需要尽量遵循的原则

- 原子化
- 面向对象(过程标准化)
- 配置文件
- 保存全部结果
- 实验可重复
- 数据与代码分离
- 代码规范/充分注释原则
- git/GitHub管理
- 充分利用多个程序入口
- 测试原则

```
# -----below are methods for visualization-----
def set_plot_style():
    plt.style.use('ggplot')
    # set font to Times New Roman
    plt.rcParams["font.family"] = "serif"
    plt.rcParams["font.serif"] = ["Times New Roman"]
    plt.rcParams["font.size"] = 20
    # set pixels per inches
    plt.rcParams['figure.dpi'] = 300

def plot_price_carat(df, ax=None):
    if ax is None:
        fig, ax = plt.subplots(1, 1, dpi=300, figsize=(16, 4))
    ax.plot(df.carat, df.price, "b+", alpha=0.2)
    ax.set_xlabel("weight (Carat)")
    ax.set_ylabel("price ($)")
    plt.xticks(np.arange(0, 5.5, 0.25))
    return ax

def plot_price_carat_seaborn(df, ax=None, hue="clarity"):
    if ax is None:
        fig, ax = plt.subplots(1, 1, dpi=300, figsize=(16, 4))
    # ax.plot(df.carat, df.price, "b+", alpha=0.2)
    sns.scatterplot(data=df, x="carat", y="price",
                     hue=hue, ax=ax, alpha=0.2)
    ax.set_xlabel("weight (Carat)")
    ax.set_ylabel("price ($)")
    plt.xticks(np.arange(0, 5.5, 0.25))
    return ax

def plot_price_carat_facet(df):
    sea = sns.FacetGrid(
        df, row="clarity", col="color",
        margin_titles=True)
    sea.map(sns.scatterplot, "carat", "price")
```

需要尽量遵循的原则

- 原子化
- 面向对象(过程标准化)
- 配置文件
- 保存全部结果
- 实验可重复
- 数据与代码分离
- 代码规范/充分注释原则
- git/GitHub管理
- 充分利用多个程序入口
- 测试原则

```
107 # I like to write such processing object to organize the data pipeline
108 class DiamondPreprocess:
109     def __init__(self, config_dict=None):
110         if config_dict is None:
111             config_dict = {
112                 "source": "local", # either "local" or "seaborn"
113                 "carat_range": None,
114                 "cuts_types": [],
115                 "colors": [],
116             }
117
118         self.config = config_dict
119         self.df = None
120
121     def step0_load_dataset(self):
122         source = self.config["source"]
123         if source == "seaborn":
124             self.df = load_data_from_seaborn()
125         elif source == "local":
126             self.df = load_data_from_local()
127         else:
128             print("the source is not recognized")
129             # consider to raise an Error if you want
130             # raise Exception(f"the data source {source} is not recognized")
131
132             # after loading the dataset, set the column types correctly
133             set_categorical_columns()
134
135     def step1_view_basic_info(self):
136         print_diamond_basic_info(df)
137
138     def step2_data_selection(self):
139         # step 1-1. filter by carat
140         carat_lb, carat_ub = *self.config["carat_range"]
141         df = filt_by_carat(df, carat_lb, carat_ub)
142
143         # step 1-2. filter by cuts
144         cuts_types = self.config["cuts"]
145         df = filt_by_cuts(cuts_types)
146
147         # step 2-3. filter by colors
148         colors = self.config["colors"]
149         df = filt_by_price(df)
150
151         return df
152
153     def step3_data_cleaning(self):
154         pass
155
156     def step4_data_transformation(self):
157         pass
158
159     def step5_export_outputs(self):
160         self.df.to_csv("../data_output/diamonds_of_interest.csv")
161
162     def pipeline(self):
163         # just chain the methods step by step
164         self.step0_load_dataset()
165         self.step1_data_selection()
166         self.step2_data_cleaning()
167         self.step4_data_transformation()
168         self.step5_export_outputs()
```

```
1 dataset:
2   script_path: ./datasets/cifar10_keras.py
3 model:
4   script_path: ./models/optimized.py
5 optimizer:
6   script_path: ./optimizers/adam_keras.py
7   initial_lr: 0.0001
8 train:
9   script_path: ./train/train_keras.py
10 artifacts_path: ./artifacts/cifar10_opt/
11 batch_size: 64
12 epochs: 1000
13 data_augmentation:
14   samplewise_center: False
15   samplewise_std_normalization: False
16   rotation_range: 0
17   width_shift_range: 0.1
18   height_shift_range: 0.1
19   horizontal_flip: True
20   vertical_flip: False
21   zoom_range: 0
22   shear_range: 0
23   channel_shift_range: 0
24   featurewise_center: False
25   zca_whitening: False
26 evaluate:
27   batch_size: 1000
28   augmentation_factor: 32
29 data_augmentation:
30   samplewise_center: False
31   samplewise_std_normalization: False
32   rotation_range: 0
33   width_shift_range: 0.15
34   height_shift_range: 0.15
35   horizontal_flip: True
36   vertical_flip: False
37   zoom_range: 0
38   shear_range: 0
39   channel_shift_range: 0
40   featurewise_center: False
41   zca_whitening: False
```

my_configuration.yaml hosted with ❤ by GitHub

yaml文件

```
  "dataset": {
    "script_path": ".../datasets/cifar10_keras.py"
  },
  "model": {
    "script_path": ".../models/optimized.py"
  },
  "optimizer": {
    "script_path": ".../optimizers/adam_keras.py",
    "initial_lr": 0.0001
  },
  "train": {
    "script_path": ".../train/train_keras.py",
    "artifacts_path": ".../artifacts/cifar10_opt/",
    "batch_size": 64,
    "epochs": 1000,
    "data_augmentation": {
      "samplewise_center": false,
      "samplewise_std_normalization": false,
      "rotation_range": 0,
      "width_shift_range": 0.1,
      "height_shift_range": 0.1,
      "horizontal_flip": true,
      "vertical_flip": false,
      "zoom_range": 0,
      "shear_range": 0,
      "channel_shift_range": 0,
      "featurewise_center": false,
      "zca_whitening": false
    }
  },
  "evaluate": {
    "batch_size": 1000,
    "augmentation_factor": 32,
    "data_augmentation": {
      "samplewise_center": false,
      "samplewise_std_normalization": false,
      "rotation_range": 0,
      "width_shift_range": 0.15,
      "height_shift_range": 0.15,
      "horizontal_flip": true,
      "vertical_flip": false,
      "zoom_range": 0,
      "shear_range": 0,
      "channel_shift_range": 0,
      "featurewise_center": false,
      "zca_whitening": false
    }
  }
}
```

JSON文件

```
{'dataset': {'script_path': '../datasets/cifar10_keras.py'},
'model': {'script_path': '../models/optimized.py'},
'optimizer': {'script_path': '../optimizers/adam_keras.py',
'initial_lr': 0.0001},
'train': {'script_path': '../train/train_keras.py',
'artifacts_path': '../artifacts/cifar10_opt/',
'batch_size': 64,
'epochs': 1000,
'data_augmentation': {'samplewise_center': False,
'samplewise_std_normalization': False,
'rotation_range': 0,
'width_shift_range': 0.1,
'height_shift_range': 0.1,
'horizontal_flip': True,
'vertical_flip': False,
'zoom_range': 0,
'shear_range': 0,
'channel_shift_range': 0,
'featurewise_center': False,
'zca_whitening': False}},
'evaluate': {'batch_size': 1000,
'augmentation_factor': 32,
'data_augmentation': {'samplewise_center': False,
'samplewise_std_normalization': False,
'rotation_range': 0,
'width_shift_range': 0.15,
'height_shift_range': 0.15,
'horizontal_flip': True,
'vertical_flip': False,
'zoom_range': 0,
'shear_range': 0,
'channel_shift_range': 0,
'featurewise_center': False,
'zca_whitening': False}}}}
```

Python 基本数据结构：字典Dictionary

需要尽量遵循的原则

- 原子化
- 面向对象
- 配置文件
- 保存全部结果
- 实验可重复
- 数据与代码分离
- 代码规范/充分注释原则
- git/GitHub管理
- 充分利用多个程序入口
- 测试原则

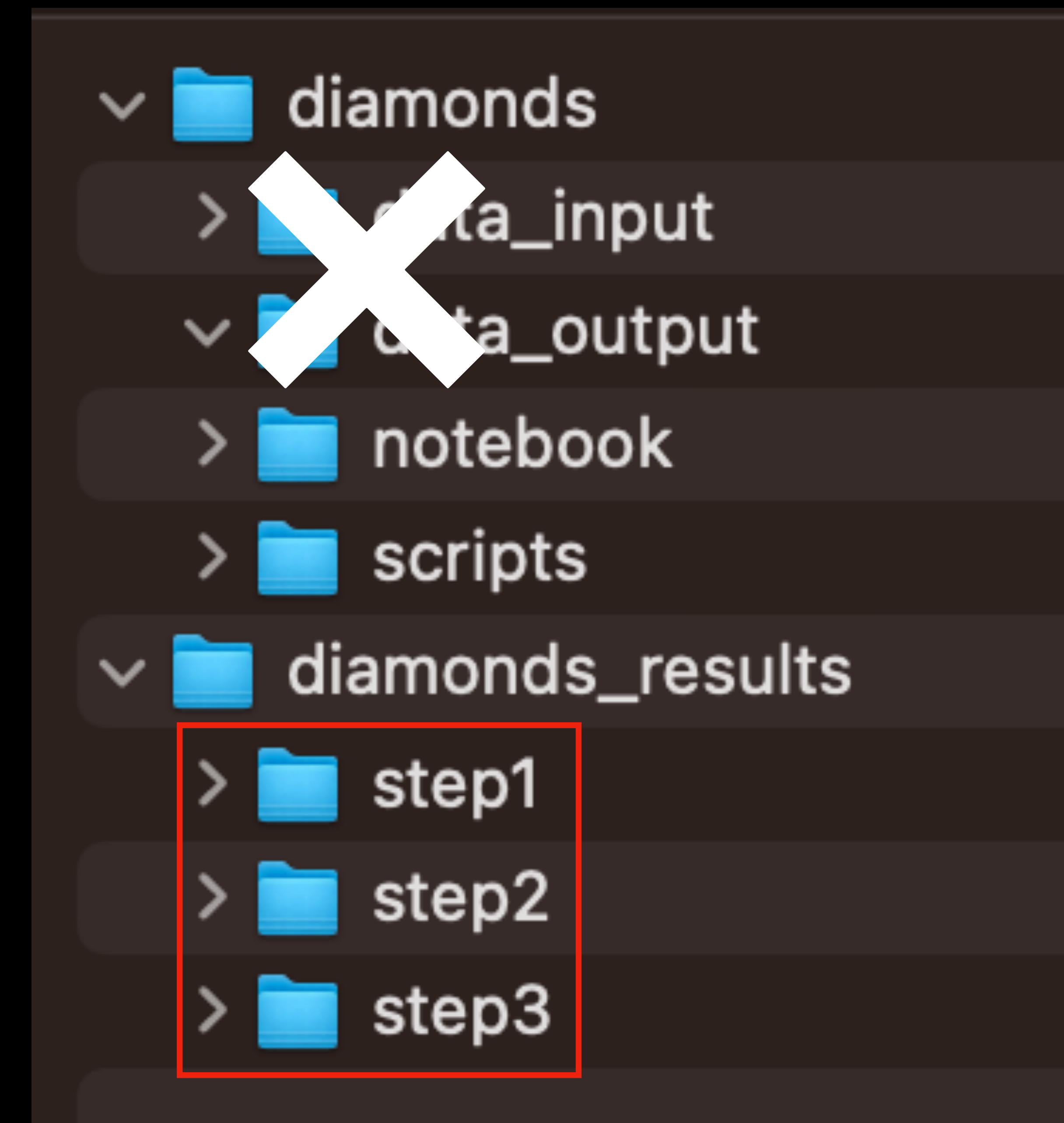
`np.random.seed(0)` makes the random numbers predictable

```
>>> numpy.random.seed(0) ; numpy.random.rand(4)
array([ 0.55,  0.72,  0.6 ,  0.54])
>>> numpy.random.seed(0) ; numpy.random.rand(4)
array([ 0.55,  0.72,  0.6 ,  0.54])
```

```
>>> numpy.random.rand(4)
array([ 0.42,  0.65,  0.44,  0.89])
>>> numpy.random.rand(4)
array([ 0.96,  0.38,  0.79,  0.53])
```

需要尽量遵循的原则

- 原子化
- 面向对象(过程标准化)
- 配置文件
- 保存全部结果
- 实验可重复
- 数据与代码分离
- 代码规范/充分注释原则
- git/GitHub管理
- 充分利用多个程序入口
- 测试原则



需要尽量遵循的原则

- 原子化
- 面向对象(过程标准化)
- 配置文件
- 保存全部结果
- 实验可重复
- 数据与代码分离
- 代码规范/充分注释原则
- git/GitHub管理
- 充分利用多个程序入口
- 测试原则

The screenshot shows a code editor window titled "Solver.py". The code implements a quadratic equation solver:

```
import math

class Solver:
    def demo(self):
        d = def demo(self,a,b,c):
            if d >= 0:
                disc = math.sqrt(d)
                root1 = (- b + disc) / (2 * a)
                root2 = (- b - disc) / (2 * a)
                print(root1, root2)
                return
            else:
                raise Parameter 'a' unfilled
```

The code contains several PEP 8 style guide violations, which are highlighted by yellow callout boxes:

- "def demo(self,a,b,c):" has a "missing whitespace after ','" violation.
- "(- b + disc) / (2 * a)" has a "whitespace before '('" violation.
- "Parameter 'a' unfilled" appears three times, once for each parameter 'a', 'b', and 'c'.
- "too many blank lines (3)" is indicated by a callout pointing to the empty lines between the class definition and the method implementation.

The status bar at the bottom of the editor shows "4:9 CRLF UTF-8". A green and yellow decorative graphic with the letters "PC" is visible in the bottom right corner.

总结

- 核心代码的质量高，则整体犯错的机会会少很多
- 要在开发效率（速度）与可靠性之间找到一个平衡点
- 没有完美的数据流，需要单独花精力重构代码
- 针对你的专业，可以不断积累你的代码库

彩蛋 1.

```
[ds] -MBP ~ % pip install seaborn
Requirement already satisfied: seaborn in /opt/anaconda3/envs/ds/
Requirement already satisfied: matplotlib>=2.2 in /opt/anaconda3/
Requirement already satisfied: scipy>=1.0 in /opt/anaconda3/envs/
Requirement already satisfied: pandas>=0.23 in /opt/anaconda3/env
Requirement already satisfied: kiwisolver>=1.0.1 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages (0.1
Requirement already satisfied: seaborn in /opt/anaconda3/envs/ds/lib/python3.8/site-packages (0.1
Requirement already satisfied: scipy>=1.0 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages (f
Requirement already satisfied: pandas>=0.23 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages
Requirement already satisfied: numpy>=1.15 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages
Requirement already satisfied: matplotlib>=2.2 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages (from seaborn) (3.5.1)
Requirement already satisfied: packaging>=20.0 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages (from matplotlib>=2.2->seaborn) (21.3)
Requirement already satisfied: pyparsing>=2.2.1 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages (from matplotlib>=2.2->seaborn) (3.0.7)
Requirement already satisfied: cycler>=0.10 in /opt/anaconda3/envs/ds/lib/python3.8/site-packages (from matplotlib>=2.2->seaborn) (0.11.0)
```

彩蛋 1. Google CoLab

```
!pip install torch-scatter  
!pip install torch-sparse  
!pip install torch-cluster  
!pip install torch-spline-conv  
!pip install torch-geometric  
!pip install torch-geometric-temporal
```



安装一些PyTorch与GNN有关的包

彩蛋 1. Google Colab...

```
from google.colab import drive  
drive.mount('/content/drive')
```

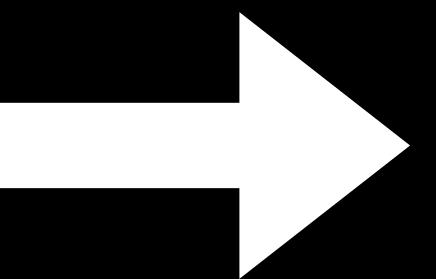
Mounted at /content/drive

Permit this notebook to access your Google Drive files?

This notebook is requesting access to your Google Drive files. Granting access to Google Drive will permit code executed in the notebook to modify files in your Google Drive. Make sure to review notebook code prior to allowing this access.

No thanks

[Connect to Google Drive](#)



“My Drive”

/content/drive

```
# create foo.txt file in your Google Drive  
# with open('/content/drive/My Drive/foo.txt', 'w') as f:  
#   f.write('Hello Google Drive!')
```

```
# print the text file  
# !cat /content/drive/My\ Drive/foo.txt
```

```
# disconnect from Google Drive  
# drive.flush_and_unmount()
```

My Drive ▾

Type ▾

People ▾

Modified ▾

Name

Owner

Last modified ▾

File size

📄 foo.txt

👤 me

Jun 1, 2022 me

19 bytes

安装python包，并保存到Google Drive

```
[6] # How do I install a library permanently in Colab?  
# https://stackoverflow.com/questions/55253498/how-do-i-install-a-library-permanently-in-colab  
import os  
from google.colab import drive  
nb_path = '/content/notebooks'  
  
# create a symbolic link pointing to src named dst  
# os.symlink('/content/drive/My Drive/Colab Notebooks', nb_path)  
  
[7] import sys  
sys.path.insert(0,nb_path)  
  
,  
[8] # install a package at the preferred location  
# (only install for once)  
!pip install --target=$nb_path torch-scatter --upgrade  
!pip install --target=$nb_path torch-sparse --upgrade  
!pip install --target=$nb_path torch-cluster --upgrade  
!pip install --target=$nb_path torch-spline-conv --upgrade  
!pip install --target=$nb_path torch-geometric --upgrade  
!pip install --target=$nb_path torch-geometric-temporal --upgrade
```

安装python包到Google Drive 路径

My Drive > Colab Notebooks

Type People Modified

Name	Owner	Last modified	File size	⋮
torch	me	Jun 1, 2022 me	—	⋮
Cython-0.29.30.dist-info	me	Jun 1, 2022 me	—	⋮
torch_geometric_temporal-0.52.0.dist-info	me	Jun 1, 2022 me	—	⋮
torch_geometric_temporal	me	Jun 1, 2022 me	—	⋮
pyximport	me	Jun 1, 2022 me	—	⋮
Cython	me	Jun 1, 2022 me	—	⋮
decorator-4.4.2.dist-info	me	Jun 1, 2022 me	—	⋮
networkx-2.6.3.dist-info	me	Jun 1, 2022 me	—	⋮
torch-1.11.0.dist-info	me	Jun 1, 2022 me	—	⋮
networkx	me	Jun 1, 2022 me	—	⋮
share	me	Jun 1, 2022 me	—	⋮
caffe2	me	Jun 1, 2022 me	—	⋮
typing_extensions-4.2.0.dist-info	me	Jun 1, 2022 me	—	⋮

再次执行. . .

```
[6] # How do I install a library permanently in Colab?  
# https://stackoverflow.com/questions/55253498/how-do-i-install-a-library-permanently-in-colab  
import os  
from google.colab import drive  
nb_path = '/content/notebooks'  
  
# create a symbolic link pointing to src named dst  
# os.symlink('/content/drive/My Drive/Colab Notebooks', nb_path)
```

```
[7] import sys  
sys.path.insert(0,nb_path)
```

```
[8] # Install a package in the preferred location  
# !pip install torch-scatter  
!pip install torch-sparse  
!pip install torch-cluster  
!pip install torch-spline-conv  
!pip install -target=$nb_path torch-geometric-temporal --upgrade  
!pip install -target=$nb_path torch-scatter --upgrade  
!pip install -target=$nb_path torch-sparse --upgrade  
!pip install -target=$nb_path torch-cluster --upgrade  
!pip install -target=$nb_path torch-spline-conv --upgrade  
!pip install -target=$nb_path torch-geometric-temporal --upgrade
```

彩蛋 2. Split - Apply - Combine 策略

- “分而治之” /Divide and Conquer
- Split-apply-combine
- 找到每个打磨级别中最便宜的钻石iamond

- 1. 分组(Groupby)
- 2. 分析(Apply)
- 3. 合并(Combine)

```
def get_cheapest_by_group(dfg):  
    idx = np.argmin(dfg.price.tolist())  
    return dfg.iloc[idx, :]
```

```
cheapest = df_diamonds.groupby(["cut"]).apply(get_cheapest_by_group)  
cheapest
```

	carat	cut	color	clarity	depth	table	price	x	y	z
	cut									
	Fair	0.22	Fair	E	VS2	65.1	61.0	337	3.87	3.78
	Good	0.23	Good	E	VS1	56.9	65.0	327	4.05	4.07
	Ideal	0.23	Ideal	E	SI2	61.5	55.0	326	3.95	3.98
	Premium	0.21	Premium	E	SI1	59.8	61.0	326	3.89	3.84
	Very Good	0.24	Very Good	J	VVS2	62.8	57.0	336	3.94	3.96

Split - Apply - Combine策略(大数 据)

- Spark: MapReduce
 - 问题：统计水果的数量

