



# Python with CSV

---



# Index

- CSV format
- CSV library
- Pandas library
- Pandas vs CSV performance
- Conclusion



# What is .csv format



# What is CSV format?



```
name,age,sex
Loc,22,male
An,21,female
Nam,33,male
Chino,11,female
Kokoto,12,male
Stephen,34,male
Curry,14,female
Klay,23,male
Petter,43,male
Parker,34,male
Lily,28,female
Potter,23,female
Nhan,22,male
Michael,21,female
Jackson,33,male
Stark,11,female
```

	name	age	sex
1	Loc	22	male
2	An	21	female
3	Nam	33	male
4	Chino	11	female
5	Kokoto	12	male
6	Stephen	34	male
7	Curry	14	female
8	Klay	23	male
9	Petter	43	male
10	Parker	34	male
11	Lily	28	female
12	Potter	23	female
13	Nhan	22	male
14	Michael	21	female

- CSV stands for "comma separated values"
- Used to store the tabular data

# CSV module in Python

---

# CSV module in Python

- A build-in module in Python
- Provided 2 approaches to read/write .csv file:
  - Default approach
  - Dictionary approach

# csv.reader

---

```
1 with open(file="sample-input.csv", mode="r") as csv_file:
2     csv_reader = csv.reader(csv_file)
3     header = next(csv_reader)
4
5     print("Header:", header)
6     print("Data:")
7     for line in csv_reader:
8         print(line)
9
```

Header: ['name', 'age', 'sex']

Data:

['Loc', '22', 'male']

['An', '21', 'female']

['Nam', '33', 'male']

['Chino', '11', 'female']

['Kokoto', '12', 'male']

['Stephen', '34', 'male']

['Curry', '14', 'female']

# csv.writer

---

```
data = [{"name", "age", "sex"},
        ["Chino", 14, "male"],
        ["Kokoto", 21, "female"],
        ["Tom", 23, "male"]]

with open(file="sample-output.csv", mode="w") as csv_file:
    csv_writer = csv.writer(csv_file)
    csv_writer.writerows(data)
```

sample-output.csv			
Delimiter: ,			
	name	age	sex
1	Chino	14	male
2	Kokoto	21	female
3	Tom	23	male



# csv.DictReader

---

```
1 with open(file="sample-input.csv", mode="r") as csv_file:
2     csv_reader = csv.DictReader(csv_file)
3
4     for line in csv_reader:
5         print(line)
```

```
{'name': 'Loc', 'age': '22', 'sex': 'male'}
{'name': 'An', 'age': '21', 'sex': 'female'}
{'name': 'Nam', 'age': '33', 'sex': 'male'}
{'name': 'Chino', 'age': '11', 'sex': 'female'}
{'name': 'Kokoto', 'age': '12', 'sex': 'male'}
{'name': 'Stephen', 'age': '34', 'sex': 'male'}
{'name': 'Curry', 'age': '14', 'sex': 'female'}
```

# csv.DictWriter

---

```
data = [  
    {'name': "Loc", 'age': 14, 'sex': "male"},  
    {'name': "An", 'age': 21, 'sex': "female"},  
    {'name': "Nam", 'age': 23, 'sex': "male"}  
]  
|  
fieldnames = ["name", "age", "sex"]  
with open(file="sample-output.csv", mode="w") as csv_file:  
    csv_writer = csv.DictWriter(csv_file, fieldnames=fieldnames)  
    csv_writer.writeheader()  
    csv_writer.writerows(data)
```

sample-output.csv			
Delimiter: ,			
	name	age	sex
1	Loc	14	male
2	An	21	female
3	Nam	23	male

# Pandas module in Python

---

# Pandas module library

- Pandas is not a build-in module
- By default, pandas use DataFrame to manipulate csv file

# pandas.read\_csv

---

```
1 import pandas as pd
2
3 data = pd.read_csv("sample-input.csv")
4 print(type(data))
5 print(data)
```

<class 'pandas.core.frame.DataFrame'>

	name	age	sex
0	Loc	22	male
1	An	21	female
2	Nam	33	male
3	Chino	11	female
4	Kokoto	12	male
5	Stephen	34	male
6	Curry	14	female
7	Klay	23	male
8	Petter	43	male
9	Parker	34	male
10	Lily	28	female



## pandas.DataFrame.to\_csv

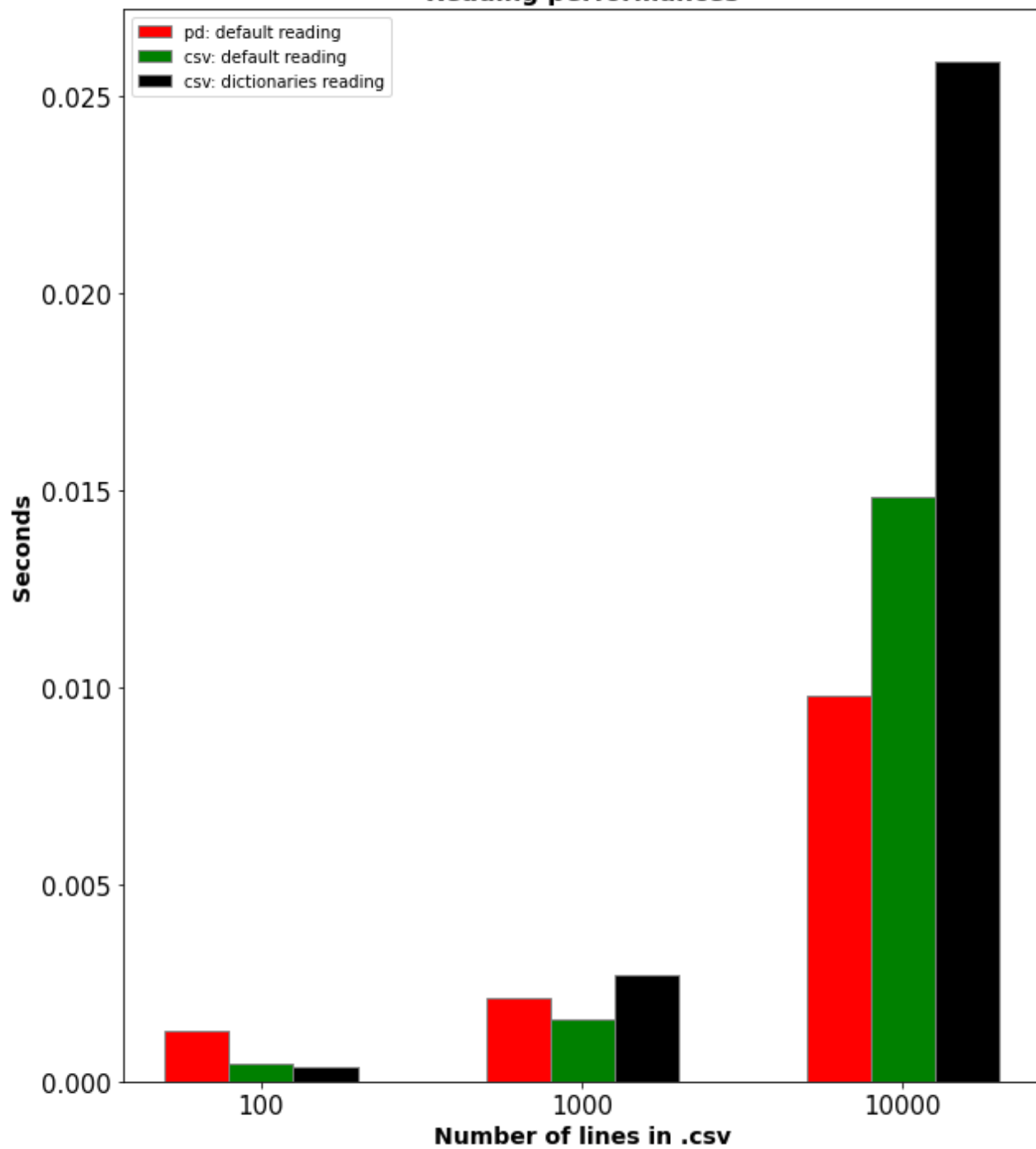
```
1 import pandas as pd
2
3 data = pd.DataFrame(
4     {
5         'name': ['Loc', 'An', 'Nam'],
6         'age': [22, 19, 25],
7         'sex': ['male', 'male', 'female']
8     }
9 )
10 data.to_csv("sample-output.csv", index=False)
11
```

sample-output.csv				
Delimiter: ,				
	name	age	sex	
1	Loc	22	male	
2	An	19	male	
3	Nam	25	female	

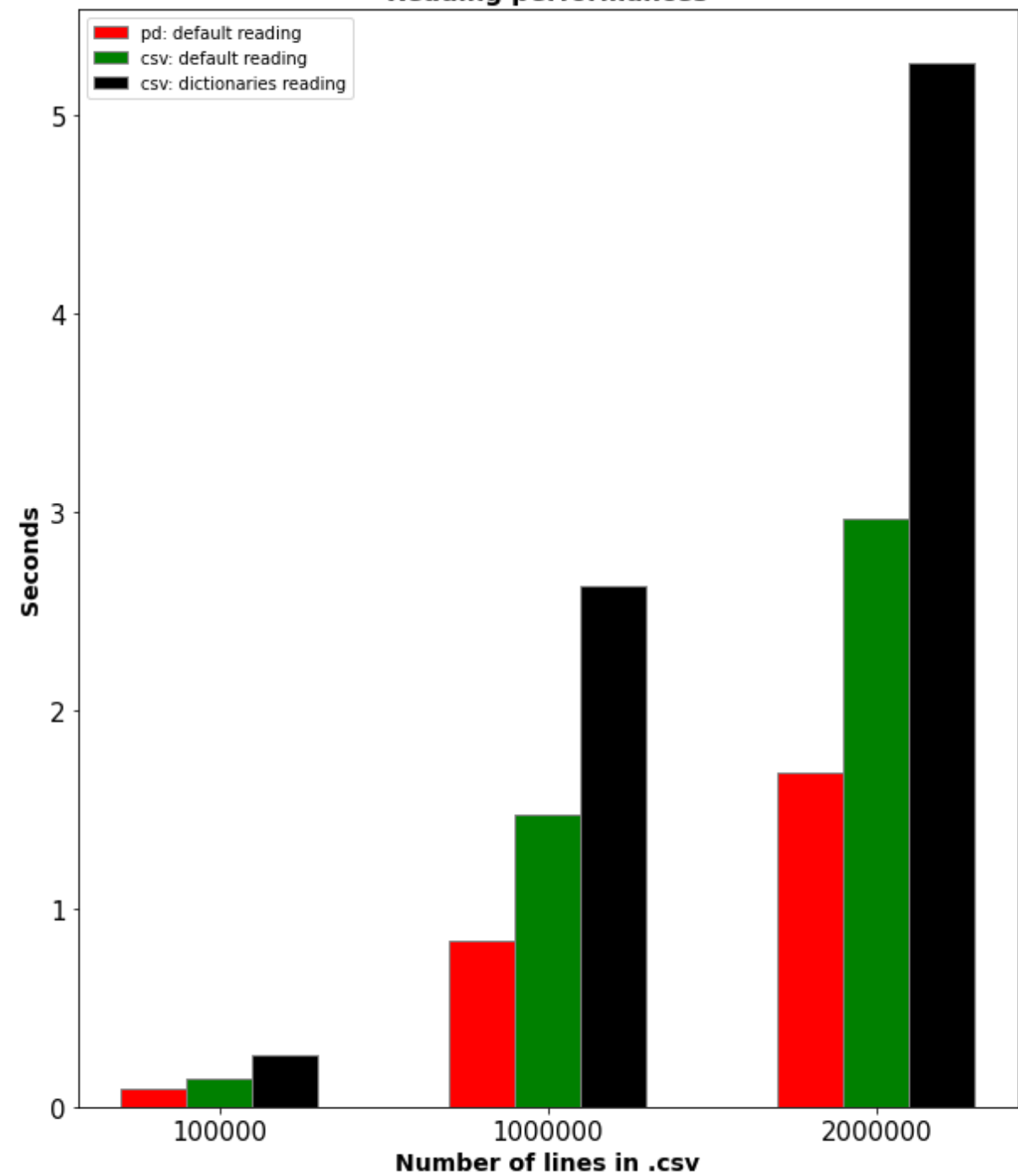
# Pandas vs CSV performance

---

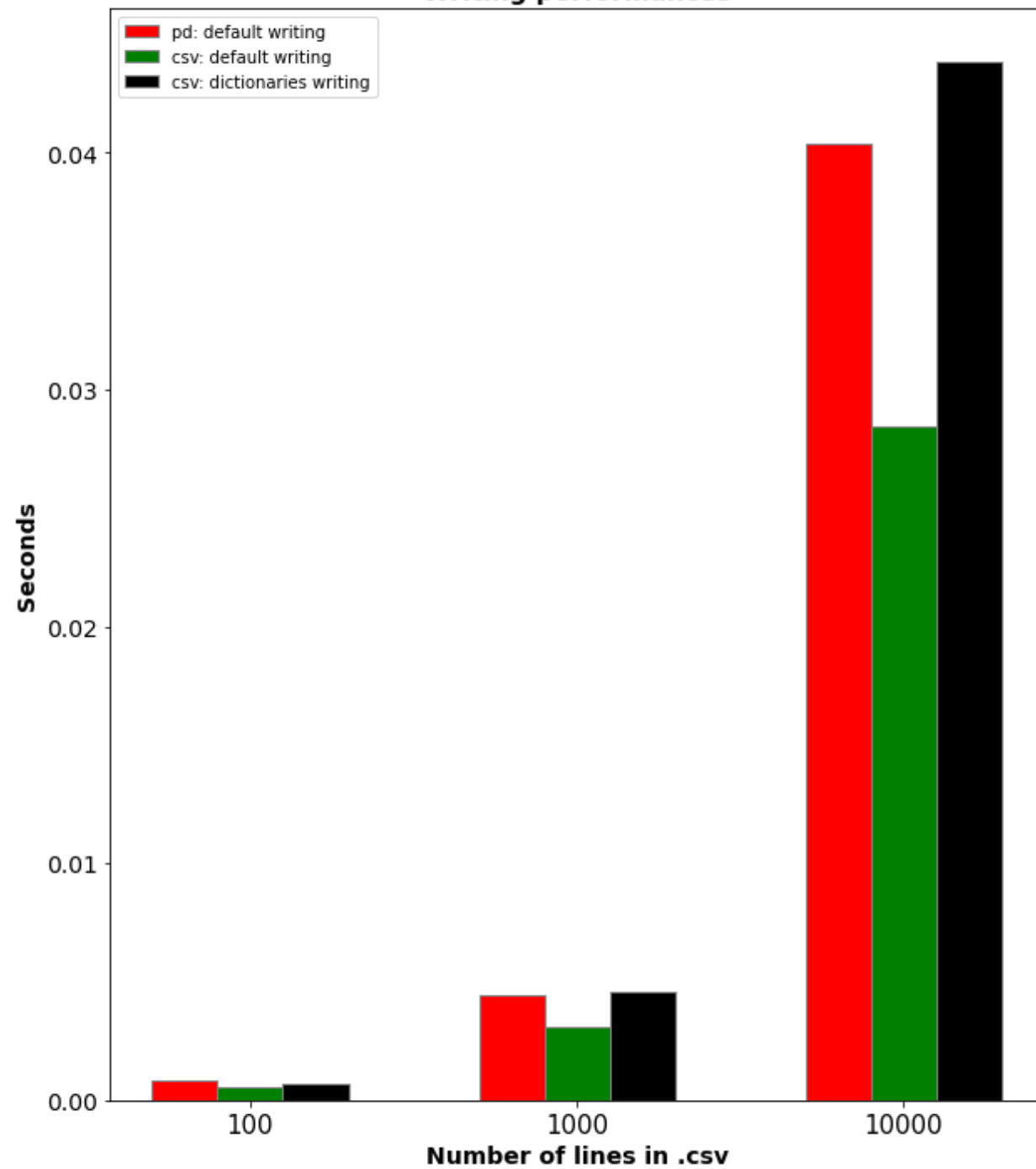
# Reading performances



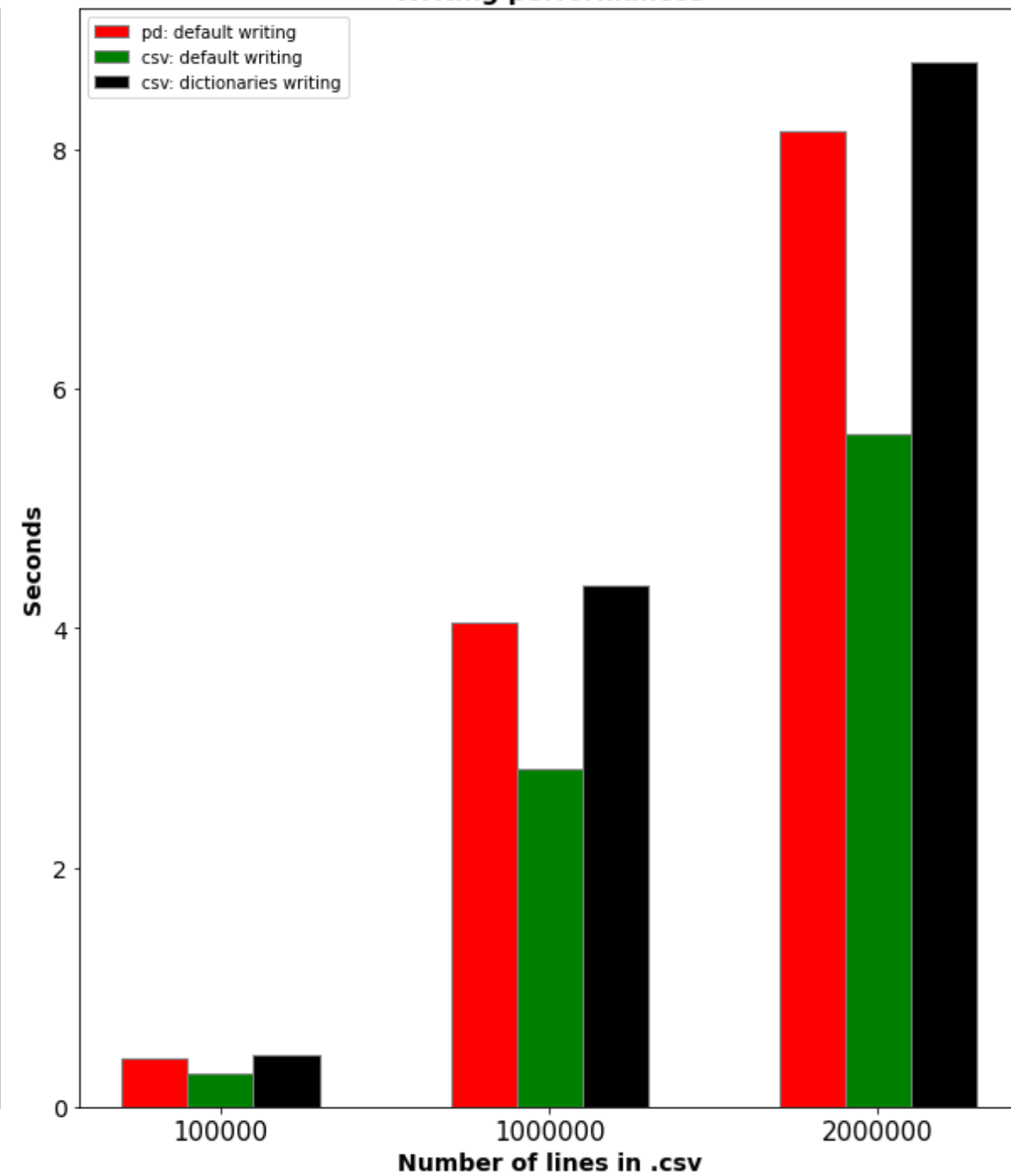
# Reading performances



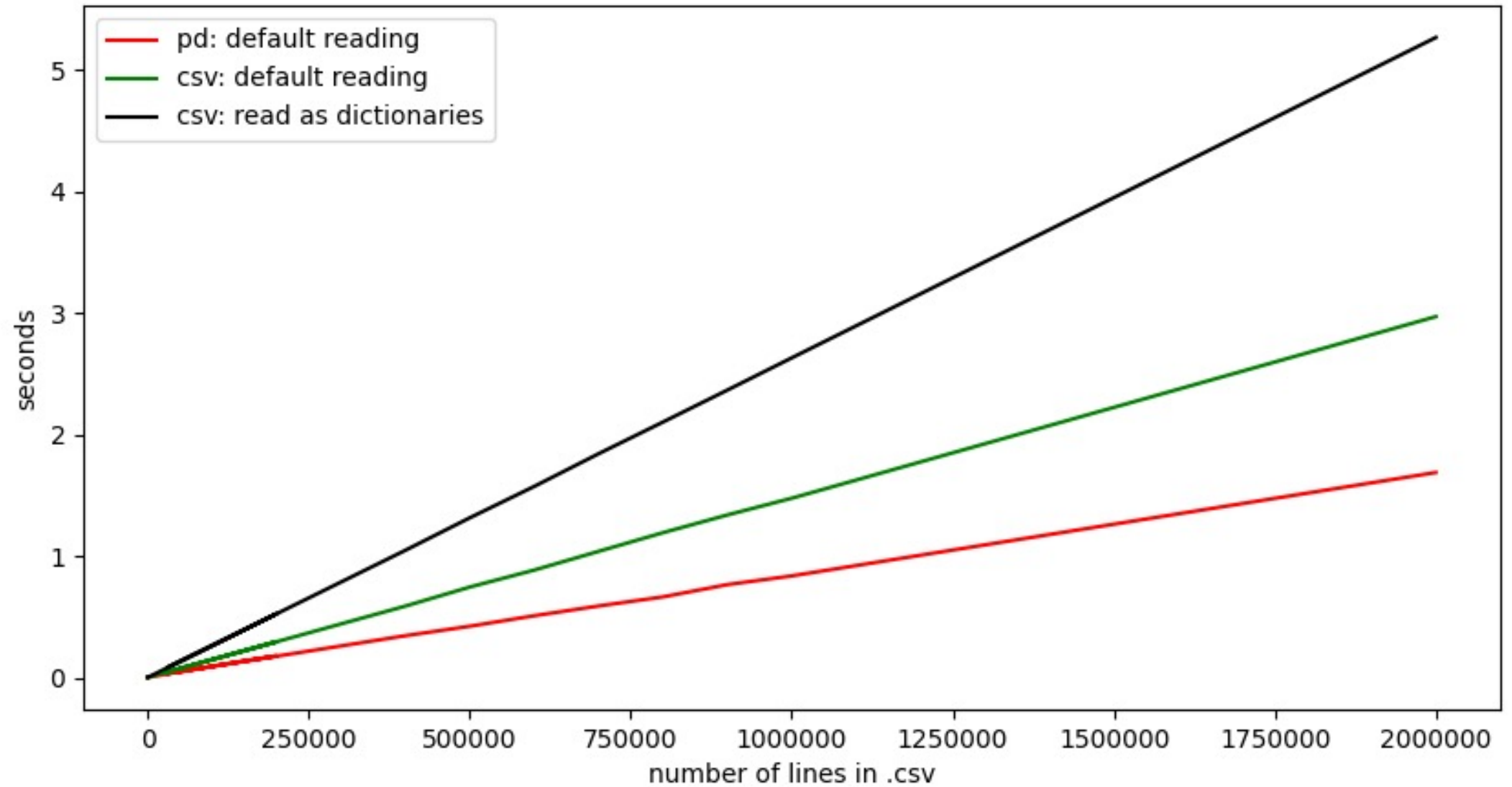
### Writing performances



### Writing performances

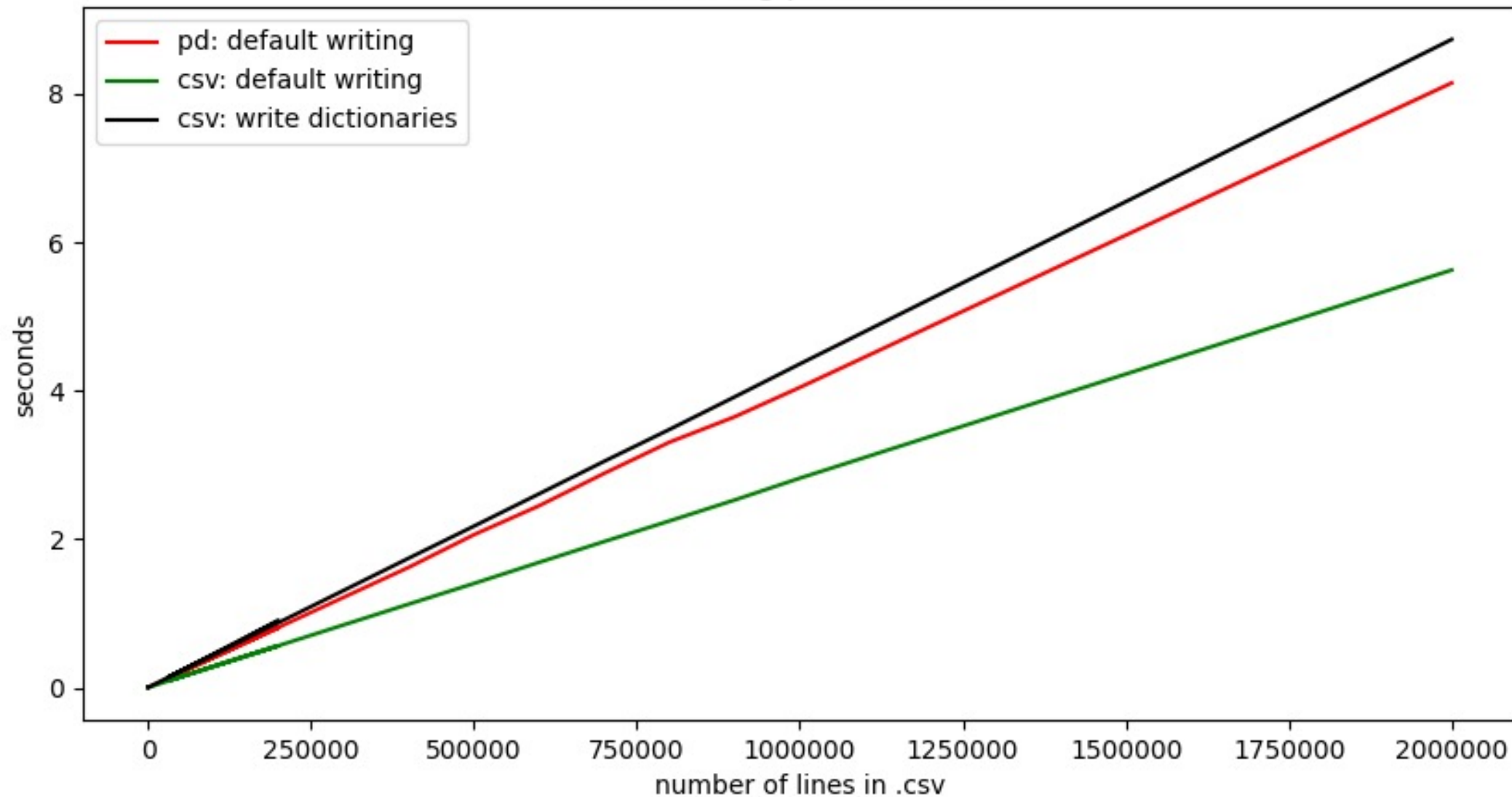


CSV reading performances





CSV writing performances



# Reading performances

- `pandas.read_csv` is slower than `csv.reader` when the csv file is less than 5000 lines
- For csv file with more than 5000 lines, `pandas.read_csv` is significantly faster than `csv.reader`, and therefore faster than `csv.DictReader`

# Writing performances

- `pandas.DataFrame.to_csv` is slower than `csv.writer`
- `pandas.DataFrame.to_csv` is faster than `csv.DictWriter`

# pandas vs csv module

pandas	csv module
Is not a built-in module → increases project dependencies	Is a built-in module of Python
Appropriate for small to big data	Appropriate for small size csv
Can be applied for .csv, excel, json, html, sql file	Only applied for .csv format